# LCP: A Low-Communication Parallelization Method
# for Fast Neural Network Inference in Image Recognition

Ramyad Hadidi, Bahar Asgari, Jiashen Cao, Younmin Bae, Da Eun Shim, Hyojong Kim,
Sung-Kyu Lim, Michael S. Ryoo†, Hyesoon Kim

Georgia Institute of Technology, †Stony Brook University

***Abstract —*** *Deep neural networks (DNNs) have inspired new studies in myriad edge applications with robots, autonomous agents, and Internet-of-things (IoT) devices. However, performing inference of DNNs in the edge is still a severe challenge, mainly because of the contradiction between the intensive resource requirements of DNNs and the tight resource availability in several edge domains. Further, as communication is costly, taking advantage of other available edge devices by using data- or model-parallelism methods is not an effective solution. To benefit from available compute resources with low communication overhead, we propose the first DNN parallelization method for reducing the communication overhead in a distributed system. We propose a low-communication parallelization (LCP) method in which models consist of several almost-independent and narrow branches. LCP offers close-to-minimum communication overhead with better distribution and parallelization opportunities while significantly reducing memory footprint and computation compared to data- and model-parallelism methods. We deploy LCP models on three distributed systems: AWS instances, Raspberry Pis, and PYNQ boards. We also evaluate the performance of LCP models on a customized hardware (tailored for low latency) implemented on a small edge FPGA and as a 16mW 0.107mm$^2$ ASIC @7nm chip. LCP models achieve a maximum and average speedups of 56x and 7x, compared to the originals, which could be improved by up to an average speedup of 33x by incorporating common optimizations such as pruning and quantization.*

## 1. Introduction & Motivation

The advancements of deep neural networks (DNNs) have made revolutionary changes in domains such as robotics [1–5], unmanned aerial vehicles (UAVs) [6–9], and Internet-of-things (IoT) [10–15]. In these domains, such as smart homes/cities/offices (*e.g.*, connected cameras, gaming consoles, TVs, routers) or collaborative robots/drones (*e.g.*, disaster relief [16–18], agriculture [19, 20], mining [21], construction [22], mapping [18, 23]), (i) ensuring an acceptable accuracy is enough (*e.g.*, detecting human sound in a disaster area with either 87% or 90% accuracy necessitates more investigation); (ii) the network of devices is standalone (*i.e.*, Internet connection is not available/necessary); and (iii) the network has a unified ownership and hence communication among devices is not hazardous (*e.g.*, IoTs at home, robots at warehouse). In such domains, executing inference in-the-edge could enable several features; however, performing the heavy inference computations locally is still a challenge. *This paper enables performing DNN inference locally with efficient distribution and parallelization in the edge environments.*

The widespread approach to deal with the heavy inference computations of DNNs is to offload the requests and private data to high-performance servers of cloud providers [24, 25]. However, cloud-based offloading is not always available (*e.g.*, no Internet access) and often relies on unreliable network latency. Furthermore, with the exponential increase in the number of edge devices [26] and the scale of raw collected data, centralized cloud-based approaches might not scale [27, 28]. Privacy concerns [29–31] and personalization are other main driving forces for in-the-edge inference. However, the challenge is that performing DNN inference locally in the edge demands high compute and memory resources that contradict the energy, computational, and economical profile of edge devices [32, 33].

**The Current Approach:** The current approach for enabling local DNN inference while adhering to edge devices computational and economical profile is to locally distribute inference computations by taking advantage of the existing surrounding devices such as idle IoT devices. The distribution is based on data- or model-parallelism methods [34, 35]. In data parallelism, the entire model is duplicated on each device for performing *separate inferences*. Hence, the system needs several live and concurrent inputs to be efficient without real-time jitter. Simply put, data parallelism only increases throughput. In model parallelism, the model is divided and distributed across several devices for *the same inference*.

**The Key Challenge:** The communication overhead and the inherent inter-layer data dependency limits effective parallelism. Therefore, an ideal parallelization method for edge devices, must minimize the communication overhead, while yielding low memory and computation footprints per node. However, none of the current distribution methods jointly reduce memory usage, computations, and communication (see Table 1). §2 presents a detailed description.

**Our Solution:** To address the aforementioned challenge, we propose a low-communication parallelization (LCP) method that enables the following: *(i) Reduces Communication:* LCP models replace a single, wide, and deep model with several narrow ones that only communicate for input and pre-final activations. Thus, their communication load is low with distributions (see Table 1). *(ii) Reduces Compute & Memory Footprints Per Node:* LCP models have fewer connections than those of the original ones, so their number of parameters and computational demands are also lower than those of the peer model-parallelism versions, shown in Table 1. *(iii) Allows Inter-Layer Parallelism:* Narrow branches in LCP models are independent of each other, which enables inter-layer parallelism. This is in contrast to model parallelism

**Table 1: Methods for distributing inference computations.**

| | Data Parallelism | Model Parallelism | Target | LCP |
|---|---|---|---|---|
| **Memory Per Device** | DNN | $\frac{1}{n}$DNN | $\frac{1}{n}$DNN | $\leq \frac{1}{n}$DNN |
| **Communication Per Inference** | IN/OUT | Intermediates +IN/OUT | IN/OUT | $\approx$ IN/OUT |
| **Computation Per Device** | DNN | $\frac{1}{n}$DNN | $\frac{1}{n}$DNN | $\leq \frac{1}{n}$DNN |

DNN: Metrics associated with the entire model; $n$: Number of devices.

that only allows intra-layer parallelism due to the single-chain dependency between consecutive layers. *(iv): Recovers Accuracy with no Additional Parameters:* After splitting the model into branches, to recover a possible accuracy loss, LCP may slightly fatten the branches. However, since it reduces unnecessary communications, the overall parameters after fattening are still fewer than the original one.

LCP is orthogonal and an addition to current techniques such as weight pruning [32] and quantization [36] that reduce the computational demand of models. LCP models offer distribution/parallelism opportunities for distributed computing, whereas current techniques apply accuracy/performance trade-offs to single-node models. Such techniques can be applied to each branch of our method, as shown in §4.3. Thus, LCP complement such techniques rather than compete with them.

**Experiments Overview: (1)** We generate and evaluate LCP models based on image-recognition DNNs on MNIST [37], CIFAR10/100 [38], Flower102 [39], and ImageNet [40] datasets (total of 53 training results), covering all MLPerf [41] image-recognition models. **(2)** To evaluate the execution performance of our method, we conduct real-world implementations on three distributed systems with up to 10 Raspberry Pis (RPis), two PYNQ boards, and up to eight AWS instances. RPis are chosen because they represent the de facto choice for several robotic and edge use cases and they are readily available [42–46]. **(3)** We also evaluate the performance of LCP on customized hardware. Because, besides tailoring models based on hardware limitations, the architecture of hardware could be tailored to better achieve the goal of fast inference. To this end, we slightly modify the architecture of TPU [47] to make it latency-optimized for edge applications rather than throughput-optimized, and implement it on a small Xilinx FPGA. **(4)** To further investigate area and power efficiency of our tailored hardware for integrating with edge devices, we implement an ASIC chip in ASAP 7 nm [48].

**Contributions:** Our contributions are as follows:

- We propose the first DNN parallelization method to reduce the communication overhead for distributed inference.
- We generate LCP models, with inter-layer parallelism for fast inference at small memory and computation footprints.
- We investigate the impact of hardware/software co-design on inference performance, by tailoring the hardware of TPU [47] for optimizing single-batch inference latency, and implement it on a small FPGA and as a tiny 0.107mm$^2$ low-power chip consuming only 16mW.

## 2. Challenges

We first explain inevitable resource limitation for executing DNNs causing the single device Pareto frontier. Then, we summarize current distribution methods and their limitations, causing straggler problem and limited scope of parallelism.

**Resource Limitation & Pareto Frontier:** DNN models consist of several layers, the computations of which are based on custom weights that are learned during the training phase with back-propagation. In the inference, feed-forward computations are performed on batched inputs, and learned parameters stay constant. The most compute- and data-intensive layers [49] are fully connected and convolution layers.[1] Figure 1 shows the number of multiply-accumulate operations and parameter size in several DNN models. As shown, generally newer models encapsulate more parameters and perform more computations for better and more generalized feature understanding than their predecessors. *In short, this trend of modern models will inevitably surpass the capabilities of any resource-constrained device.*
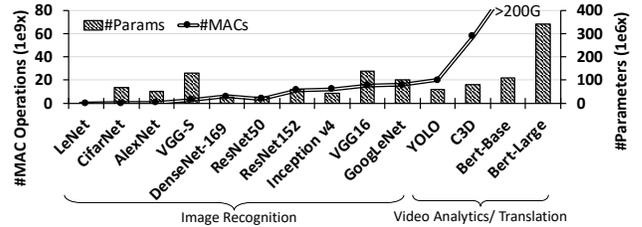


**Figure 1: DNNs #MAC operations/inference and parameters.**

The capabilities of resource-constrained platforms are limited. Figure 2 depicts latency per image using state-of-the-art image recognition models on RPi [51]. All implementations heavily utilize modern machine learning optimizations such as pruning [32], quantization, low-precision inference [36,52,53], and handcrafted models [54]. Additionally, the models are highly optimized for ARMv8 architectures using the ELL compilation tool [50]. However, achieving higher execution performance is impossible on a single device due to the Pareto frontier. As seen, the latency for high-accuracy models is longer than 400ms, and generally, latencies are longer than
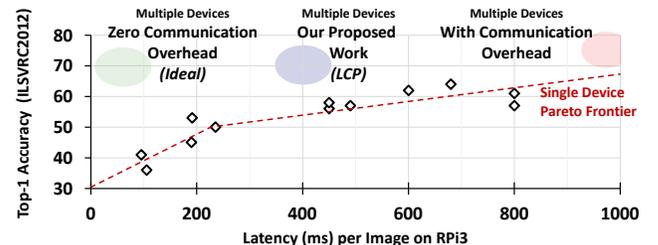


**Figure 2: Latency-Accuracy Pareto Frontier – Single device: Latency per image on RPi3 for ILSVRC models with the optimized platform-specific compilation ELL [50] tool [51]. Multiple devices: Breaking the single device Pareto frontier, but with significant communication overhead.**

---

[1]Since this paper focuses on visual models, we only introduced the layers in such models. For future work, we aim to include other types of DNNs.
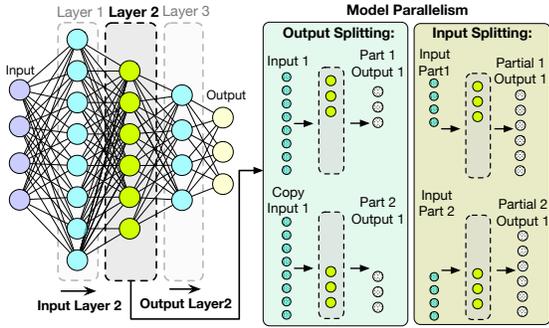
**Figure 3: Model parallelism for a fully connected layer.**

100ms. In addition, the data shown in the figure is only for image-recognition models; DNNs in other domains are already surpassing these models in size and complexity. Fitting such an exponentially increasing computation on a single device, especially for edge devices, is a limiting factor for executing DNNs in the edge. *In other words, even after applying all optimization techniques for DNNs, the single device Pareto frontier limits the widespread applicability of DNNs in several edge domains necessitating distribution and parallelization.*

**Current Distribution Methods: (1)** Data parallelism (Figure 4a) parallelizes the computations of independent inputs [34, 35]. Data parallelism does not apply to the edge because: It (i) serves several independent requests, the number of which is limited in an edge environment; (ii) does not reduce *end-to-end latency* per inference and only increases throughput. Latency is important in several applications in the edge; and (iii) does not change the computation and memory footprints per node (Table 1).

**(2)** Model-parallelism (Figure 4b) divides the inference computations for the same request [34, 35]. This method divides the computations within layer(s) while keeping dependencies intact. Depending on the type of layer, the dividing could take several forms. Figure 3 presents a simple example for distributing a fully connected (fc) layer, illustrating two extremes of model parallelism: Input and output splitting [14]. In output splitting, producing each output(s) is divided among the devices. In input splitting, the input is split and each device computes all parts of the output that are dependent on their received input. As shown in Figure 3, each method has communication overhead (transmission of the input to all nodes or partial sums to a final node for summation). New model-parallelism methods is also crafted by
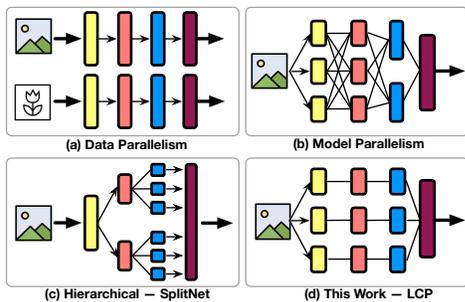


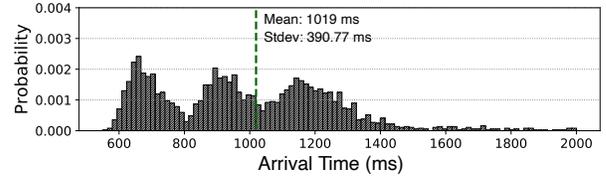**Figure 4: Overview of distribution/parallelism methods.**



**Figure 5: Histogram of prediction latencies on a six RPi system executing AlexNet with model parallelism (§4.2).**

mixing these two extremes, but they similarly suffer from the same discussed overhead. Several model-parallelism methods also exist for convolution layers by using matrix-matrix multiplication [55, 56]. Model parallelism does not change the interconnection of a model. *Hence, although model parallelism reduces the compute and memory footprint per node; the single-chain dependency between consecutive layers limits the parallelism scope within a single inference and causes communication overhead.*

**(3)** SplitNet [57], shown in Figure 4c, gradually splits the model in a tree-structured style *manually* based on the dataset semantics, extracted in intermediate to final layers. Therefore, SplitNet (i) splits only intermediate to final layers, (ii) is invariant to the number devices, (iii) creates imbalanced workload because of its dependency on semantics, (iv) results in tree-style connections, incurring high communication overhead, and (v) enforces a new splitting when dataset changes.

**Communication Overhead & Limited Parallelism:** Current distribution methods have a high communication overhead and limited scope of parallelism which stems from the single-chain dependency between consecutive layers. High communication induces the straggler problem, in which a system is lagged by its slowest node. Specifically, since edge devices usually use a wireless network, the latency deviations are high. As an example, Figure 5 depicts the histogram of prediction latencies on a distributed IoT system consisting of six RPis executing AlexNet with model parallelism. The computing time is bounded to 500ms, but the average delay is $\approx 2\times$ longer (and $\approx 4\times$ for tail latency). To gain perspective, Figure 6a shows VGG-S with model parallelism and its communication overhead. As seen, dependencies enforce a strongly interconnected network among the nodes. Although several techniques such as compression could alleviate the cost of communication, still the number of connections remains constant. *Therefore, an ideal distribution method for*
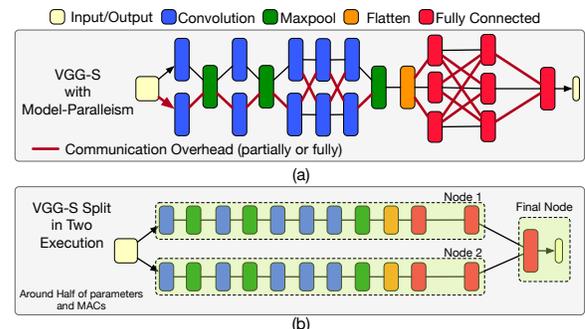


**Figure 6: VGG-S (a) model parallelism and (b) LCP versions.**

*edge devices besides yielding low memory and computation footprints per node must consider communication overhead.*

The single-chain dependency between consecutive layers limits the available parallelism that could be harvested by the aforementioned methods. The limitation is that after the computations of a single/few layer(s) are done, the intermediate results must be merged before being forwarded to the next layer. Such merging acts as a global barrier, which similar to parallel programming, limits the gained performance speedup. *In summary, with parallel execution on multiple devices, ideally, we could pass the frontier in Figure 2. However current distribution methods are limited by the communication overhead and the inherent inter-layer data dependency. The next section proposes LCP models, which significantly reduce communication and allow inter-layer parallelism.*

## 3. LCP For Fast Inference

To address challenges, we propose LCP method, which replaces a single, wide, and deep model with several narrow branches that only communicate for input and pre-final activation (Figure 4d). Figure 6b shows an example of a two-branch LCP model for VGG-S. This section first explains the design procedure of LCP models and discusses their key features enabling low-communication parallelization. The second part focuses on tailoring a systolic architecture for edge computing.

### 3.1. Tailoring Models

**Design Procedure:** Figure 7 describes the design procedure of LCP models. We start by inputting the DNN model and its per-layer memory and computation footprints. Similarly, we input the specification of the hardware, such as memory size, computation capability, and any overhead associated with executing a DNN on our hardware. For instance, several DNN frameworks have a memory overhead because of the framework. A splitter procedure, described in Procedure 1, in a while loop, splits the model, cuts the connection, and measures the approximate footprints of each branch. The $\text{Division}_{\text{Factor}}$, a hyperparameter, defines the granularity of
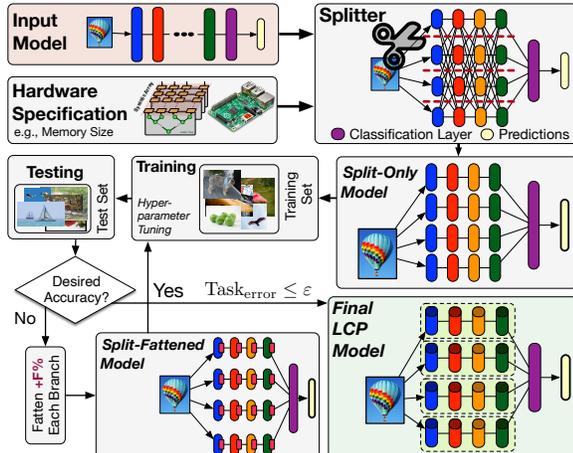


**Figure 7: Design Procedure of LCP models.**

---

**Procedure 1:** LCP Splitter (in Figure 7)

**Input** : DNN: Layer configurations $[0 : n]$
$\quad\quad\quad$ $\text{DNN}_{\text{Mem}}$, $\text{DNN}_{\text{MAC}}$: DNN memory and computational footprints
$\quad\quad\quad$ $\text{Division}_{\text{factor}}$: Division Factor for splitting
$\quad\quad\quad$ $\text{Dev}_{\text{Mem}}$, $\text{Dev}_{\text{MAC}}$: Hardware specification
**Output** : DNN: Layer configurations $[1 : n]$

1  Split(DNN, $\text{DNN}_{\text{Mem}}$, $\text{DNN}_{\text{MAC}}$, $\text{Division}_{\text{factor}}$, $\text{Dev}_{\text{Mem}}$, $\text{Dev}_{\text{MAC}}$)
2  $\quad$ $\text{Mem}_{\text{fit}} \leftarrow 0$; $\text{MAC}_{\text{Mac}} \leftarrow 0$;
3  $\quad$ **while** *not* $\text{Mem}_{\text{fit}}$ **and** *not* $\text{MAC}_{\text{Mac}}$ **do**
4  $\quad\quad$ $\text{Mem}_{\text{fit}} \leftarrow \text{DNN}_{\text{Mem}} < \text{Dev}_{\text{Mem}}$
5  $\quad\quad$ $\text{MAC}_{\text{Mac}} \leftarrow \text{DNN}_{\text{MAC}} < \text{Dev}_{\text{MAC}}$
6  $\quad\quad$ **for** layer $[0..n-1]$ *in* DNN **do**
7  $\quad\quad\quad$ layer.width $\leftarrow$ layer.width/ $\text{Division}_{\text{factor}}$
8  $\quad\quad$ RemoveNonBranchConnections(DNN)
9  $\quad$ **return** <DNN>

---

division/splitting. Here, we assume the $\text{Division}_{\text{Factor}}$ of two, but any number is viable. The loop exits when a single branch is fitted on a device (both memory and computation wise). If the number of devices is fewer than the number of branches, the execution is still possible, but will be inefficient. Then, we remove non-branch connections in a simple operation that keeps only one connection per layer. The derived model from the splitter is the *split-only model*. By training the split-only model and testing it, we measure its accuracy. The split-only models have fewer parameters and MAC operations than the original models (see Table 2) in total. Hence, after distribution, each branch has less computation and memory footprint than its model-parallelism version.

As a result of fewer number of parameters and removing several connections, a slight accuracy drop in split-only LCP models is expected. Depending on the accuracy requirement of the task, we either fatten each branch by $F\%$, a hyperparameter, or output the resulted model. We assumed a maximum of 3% bound for $\text{Task}_{\text{error}}$. Fattening each branch by $F\%$ is done by increasing the number of channels and output features of convolution and fully connected layers of the split-only model, respectively. Note that theses new *split-fattened models* are fattened within each branch. Thus, even with a high fattening percentage, still they have fewer parameters and MAC operations than the original model (see Table 3). When the accuracy is in the acceptable error range for our task, $\text{Task}_{\text{error}}$, we output the model architecture and its weights. It is expected that with similar number of parameters after fattening, LCP models achieve the same level of accuracy [58]. We showcase LCP models in §4.1 covering MLPerf [41].

**Key Features of LCP Models:** LCP models are designed by considering their underlying computation domain and have the following key features to address the challenges discussed in §2: **(1)** LCP models only communicate for input and pre-final activation. Therefore, they significantly reduce communication overhead in a distributed system. Additionally, the low communication load per inference helps with the straggler problem. This is in contrast to model parallelism, which highly depends on communication among all the intermediate layers; **(2)** LCP models split the size of a layer, so the total parameter size and computation complexity of the model are reduced. Therefore, they require fewer parameter sizes, less computa-

tion complexity, and no communication between the nodes for intermediate layers. These lower memory and computation footprints allow edge devices to efficiently operate within their limited resources (*e.g.*, no swap space activities due to limited memory); **(3)** LCP models replace the original wide model with several narrow and independent branches. Since the computations of branches are not dependent, in contrast to the single-chain of dependency in the original model, the scope of parallelism is not limited with each layer anymore. In other words, LCP models go beyond intra-layer parallelism.

## 3.2. Tailoring Hardware

Last section showed how we enable fast inference under resource constraints and at costly communication, by proposing a low-communication parallelization method that results in several narrow models. To further achieve the goal of fast inference and recognize the potential, the hardware can also be tailored. Recently, several popular tailored hardware designs for DNNs [47, 59–63] including TPU [47] use systolic arrays [64] that offer a high degree of concurrent processing through a dataflow compute arrays hence providing high *throughput*. In the edge applications, however, the main goal is *reducing single-batch inference latency*, rather than high throughput solely. This section introduces our microarchitecture (Figure 8a), an example of tailoring and simplifying the architecture of TPU to be implemented on small FPGAs or be fabricated as tiny (i.e., 0.107 mm$^2$ as shown in Figure 8b) low-power chips to be integrated with edge devices.

Figure 8a illustrates our tailored microarchitecture that similar to TPU, comprises a weight-stationary systolic array [64] for implementing matrix-matrix multiplication. The systolic array cells are organized in a 32×64 array ❶. To reduce the number of connections, only the first row of the systolic array is connected to the memory ❶. Moreover, each cell of the first row is only connected to one data stream line ❷. Based

on the type of an operand (*t*), streaming data is used for either initialization or for processing. Since the width of the systolic array is 32, a heuristic algorithm partitions the streaming (i.e., non-stationary operand) into blocks of 32 width and arbitrary length, and splits the stationary operand, into 32 × 64 blocks. To assist the smooth streaming of data from memory to the systolic array, we map these blocks along with their indices (*i*), type (stationary/non-stationary), and length to sequential memory addresses. We implement our 32×64 systolic array connected to *LPDDR2* memory with the data rate of 933Mb/s/pin @466 MHz [65], which gives a bandwidth of 3.7 GB/s. Other packaging options with higher memory bandwidths are also feasible. However, seeking a fair comparison with RPi3s, we choose this memory technology. The maximum data reuse rate of our design is 64 OPs/Byte, which leads to a peak throughput of 217.6 GOPs/s. The following explains three main modifications we made to this systolic architecture, to achieve our goal of *reducing single-batch latency*.

**(1) Adder Trees:** Instead of MAC-based systolic arrays, we separate adders from multiplications by integrating adder trees, the well established components for DNN accelerators [66–68], into systolic arrays architecture. Each cell of our systolic array is a *multiplier* with two integer operands, one stationary and one streaming (R1). Each row of the multiplier array is connected to an adder tree ❸, pipelined in five ($log_2 32$) stages. Adder trees reduce the result of multiplications into a single integer, which then contributes to creating an output element. The structure of the multiplier array connected to the adder trees reduces latency from O(n) to O(log(n)).

**(2) Simple Indexing Logic:** We use a data-driven execution model, in which data is pushed by the memory to the systolic array, triggered by the arrival of data. During execution, for each element, the indexing logic (❹) generates the appropriate row and column indices of the element using the index (*i*) of a block and its length to accompany the result. The row and column indices will later be used by the memory interface to write the results to physical locations in memory. By comparing the length and index (*i*), the end of the operations in the current layer is detected. The end of the current layer signals the start of activation and pooling functions (❺) for that layer.

**(3) Buffering Stationary Operands:** The stationary operands are often larger than the depth of the systolic array. In such cases, we have to partition a multiplication into several small operations that share a non-stationary operand, but have distinct stationary operands. To avoid multiple loads of stationary registers, we choose to integrate a buffer (❻) for stationary operands at each cell. As a result, the design serves requests with lower latency. Moreover, since each branch of the model has several layers, integrating these buffers allows fast context switching without the overhead of reloading the stationary operands. These buffers are connected in a column of cells, similar to streaming registers (R1)s. During the initialization, stationary operands are poured into these connected buffers to fill them by utilizing the connections between them.
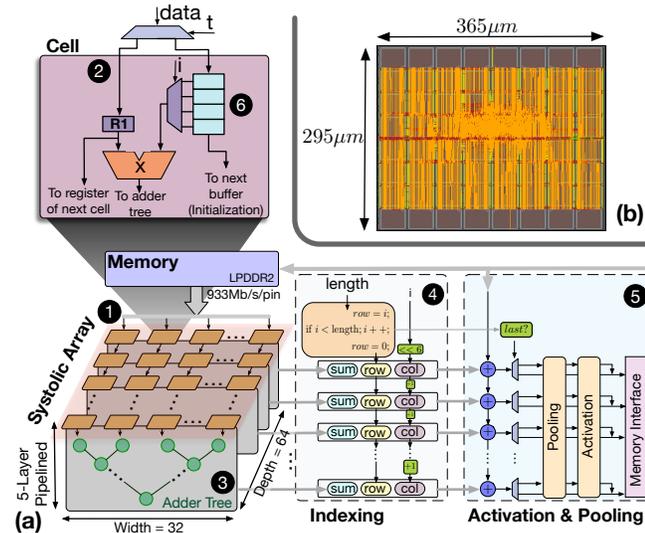


**Figure 8: Details of Tailored Hardware for Edge: (a) Microarchitecture overview, and (b) Layout of ASIC design at 7nm.**

# 4. Experimental Studies

This section shares our experimental results for generating LCP models and then their full-system implementation on RPi, TVM [69] on PYNQ boards, and AWS servers. Finally, we evaluated our hardware with edge FPGA implementation, and ASIC chip design. At the start of each subsection, the setup of related experiments is provided.

## 4.1. Generating LCP Models

**Training Specifications:** We train all the models, including the original model, from scratch to conduct a fair comparison. Normalization [70] layers are included. The training is done with an exponential learning rate with a decay factor of 0.94, initial learning rate $1e-2$, number of epoch per decay of two or 10, a dropout rate of 50%, and L2 regularization with weight decay of $5e-4$. We use ADAM optimizer [71] with $\beta_1 = 0.9$ and $\beta_2 = 0.99$. All biases are initialized to zeros and all weights are initialized with a normal distribution of mean 0 and a standard deviation of $4e-2$. All of our models are trained until the loss is flattened or least for 12 epochs. Test and accuracy measurements are done on at least 10% of datasets that have never been used in training to provide an unbiased evaluation of the model. For LCP, the $Division_{Factor}$, $F$, and $\varepsilon$, are 2%, 10%, and $\approx 3\%$, respectively.

**Datasets:** We use the following datasets: (1) MNIST [37], which contains 70k grayscale handwritten 28x28 images in 10 classes; (2) CIFAR10 [38], which contains 60k colored 32x32 images in 10 classes; (3) CIFAR100 [38], which contains 60k colored 32x32 images in 100 classes; (4) Flower102 [39], which contains 16,378 colored 224x224 images of flowers in 102 classes; and (5) ImageNet [40], which contains 1.33 M colored 224x224 images in 1000 classes.

**Models:** We use the representative model for each dataset, LeNet [72], LeNet-FC [72], VGG-S [73], CifarNet [38], VGG16 [73], AlexNetv2 [74], ResNet-18/50 [75], and MobileNet [76]. We cover all image-recognition models in MLPerf. In total, for brevity, we only report 53 instances of training results to show LCP extensibility using five datasets and nine models. Our additional results (not reported) with ResNet-34, DenseNet [77], and DarkNet19 [78] confirms extendibility. Simple sequential DNNs serve as a basis to confirm our method, while ResNets and MobileNet showcase LCP with modern models.

**Split-Only Models:** For split-only models, we use $Division_{Factor}$ of two, which results in models with two, four, and eight branches. Except the width, defined as output features in fully connected layers and the number of output channels (i.e., filters) in convolution layers, the rest of the parameters are similar to the original model as Splitter Procedure 1 only touches widths. Table 2 lists the training results. Figure 9a illustrates the accuracy difference of our models, shown in Table 2. As shown, the maximum accuracy drop is around 5% for CifarNet. Note that this accuracy drop occurs when we

**Table 2: Results of split-only LCP models.**

| Model Name | | Dataset | Layers[†] | Top-1 Accuracy | # Param | # MAC Opr. |
|---|---|---|---|---|---|---|
| LeNet-FC* | ‖ | MNIST | 3fc | 97.95 | 266.6k | 266.2k |
| LeNet | ‖ | MNIST | 2fc-3c-2p | 98.76 | 61.7k | 61.5k |
| LeNet-split2 | ‖ | MNIST | 3fc-6c-4p | 98.86 | 31.5k | 30.5k |
| LeNet-split4 | ‖ | MNIST | 5fc-12c-8p | 98.93 | 16.1k | 16.0k |
| LeNet-split8 | ‖ | MNIST | 9fc-24c-16p | 98.81 | 8.8k | 8.5k |
| CifarNet* | ‖ | Cifar10 | 2fc-2c-2p-2n-1d | 80.72 | 797.97k | 14.79M |
| CifarNet | ‖ | Cifar100 | 2fc-2c-2p-2n-1d | 52.87 | 815.34k | 14.81M |
| CifarNet-split2 | | Cifar100 | 5fc-4c-4p-4n-2d | 51.22 | 410.48k | 9.33M |
| CifarNet-split4 | | Cifar100 | 9fc-8c-8p-8n-4d | 48.48 | 208.05k | 6.59M |
| CifarNet-split8 | | Cifar100 | 17fc-16c-16p-16n-8d | 47.98 | 106.85k | 5.23M |
| VGG-S* | ‖ | Cifar100 | 3fc-5c-2p-1n-2d | 50.33 | 76.15M | 154.09M |
| VGG-S | ‖ | Flower102 | 3fc-5c-3p-1n-2d | 88.14 | 60.79M | 1.85G |
| VGG-S-split2 | | Flower102 | 5fc-10c-6p-2n-4d | 89.31 | 30.50M | 1.01G |
| VGG-S-split4 | | Flower102 | 9fc-20c-12p-4n-8d | 87.55 | 15.26M | 591.65M |
| VGG-S-split8 | | Flower102 | 17fc-40c-24p-8n-16d | 85.66 | 7.64M | 382.51M |
| ResNet-18 | ‖ | ImageNet | 18c-2p-17n | 70.68 | 11.69M | 1.82G |
| ResNet-18-split2 | | ImageNet | 35c-3p-34n | 69.85 | 6.11M | 0.98G |
| ResNet-18-split4 | | ImageNet | 69c-5p-68n | 68.07 | 3.32M | 0.55G |
| ResNet-18-split8 | | ImageNet | 137c-9p-136n | 66.76 | 1.93M | 0.34G |

[†] fc: fully-connected, c: convolution, p: pooling, n: normalization, and d: dropout.
* Detailed results are removed for brevity, refer to Figure 9. The results follows the same trend.

reduced the parameter size of our model extensively (around $1/8$). Figure 9b and c show reduction in the number of parameters and computation compared with the original DNN model; as seen, each split reduces both by about $split_{factor}$ times. This is because each convolution and fully connected layer in the split version create fewer outputs; therefore, the next layer requires fewer parameters. In the next section, we restore the accuracy of LCP models with split-fattened models.

**Split-Fattened Models** Accuracy is a defining factor in several applications. Thus, we provide a remedy to restore the accuracy of split-only models. By fattening (i.e., adding more parameters) each branch, we aim to create larger layers in the split-only models. To do so, for each layer (excluding classification layer) in every branch, we increase the width by a fraction. So, fattening by 20% means the size of the output in each layer is increased 1.2x. We fatten every branch in 10% steps as Procedure 1 shows. Our experiments focus on split8, which have the highest accuracy drops. Figure 10 shows a summary of these models. As seen, 40% split-fattened models have higher accuracy than the original model while having fewer parameters and MAC operations. On average (for 30% and
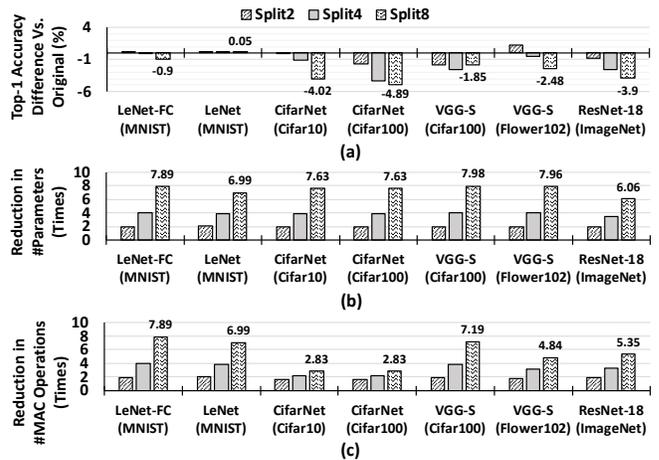


**Figure 9: Split-Only Models: (a) Accuracy, (b) reduction in the number of parameters, and (c) reduction in the number of MAC operations in comparison with the original model.**
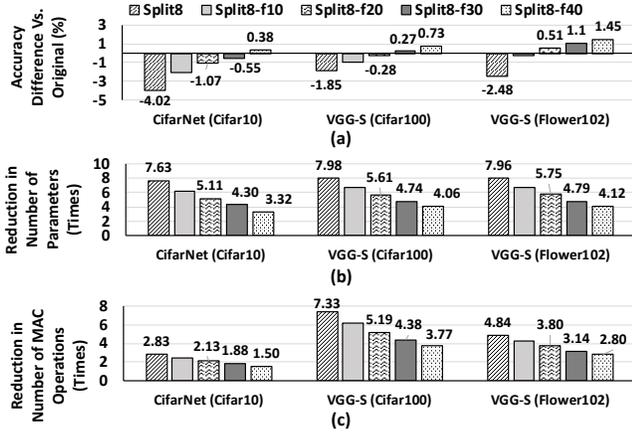
Figure 10: Split-Fattened Models – Common visual models (a) Accuracy difference, (b) reduction in the number of parameters, and (c) reduction in the number of MAC operations in comparison with the original one (Table 2).

Table 3: Results of ImageNet LCP models.

| Model Name | Dataset | Top-1 Acc. | Top-5 Acc. | # Param. | # MAC MAC Opr. |
|---|---|---|---|---|---|
| **AlexNet** | **ImageNet** | **57.02** | **80.32** | **50.3M** | **678.97M** |
| AlexNet-split8 | ImageNet | 49.03 | 73.10 | 6.32M | 145.37M |
| AlexNet-split8-f40 | ImageNet | 54.68 | 77.06 | 12.11M | 244M |
| **VGG16** | **ImageNet** | **70.48** | **90.02** | **138.36M** | **15.47G** |
| VGG16-split8 | ImageNet | 58.67 | 81.54 | 7.64M | 2.01G |
| VGG16-split8-f40 | ImageNet | 67.24 | 89.23 | 33.78M | 3.87G |
| **ResNet-50** | **ImageNet** | **75.4** | **93.1** | **22.80M** | **4.87G** |
| ResNet-split8 | ImageNet | 61.79 | 81.22 | 5.42M | 0.88G |
| ResNet-split8-f40 | ImageNet | 72.12 | 92.19 | 8.60M | 1.18G |
| **MobileNet** | **ImageNet** | **71.7** | **90** | **4.24M** | **4.86G** |
| MobileNet-split8 | ImageNet | 59.68 | 83.23 | 1.12M | 0.93G |
| MobileNet-split8-f40 | ImageNet | 68.05 | 89.12 | 2.12M | 1.34G |

For `[model_name]-f[number]`, number represent the percentage of fattening.

40% models), with 4.61x–3.81x fewer parameters and 2.95x–2.5x fewer MAC operations, split-fattened models achieve accuracy within our error bound of 3%, Task$_{error}$, while they jointly optimize memory, computation, and communication for edge.

**ImageNet Models:** Table 3 illustrates the results of ImageNet models. For the sake of brevity, we only show split8 and one fattened model. As shown, `f40` models restore the accuracy within 3% of the original model. The tradeoff for 3% accuracy loss is about 4x fewer parameters, 4x fewer computations, and 8x less communication load (vs. model parallelism). Figure 11 presents a comparative analysis for the communication load between distributed original models with model parallelism and distributed LCP models. Since LCP models avoid communication between their branches, the communication load is reduced significantly. In short, although split models are more complex than the original models in terms of the number of layers and connections, they achieve more parallelism with less communication load.

## 4.2. Exploring Performance on RPis, PYNQs, and AWS

**RPi Experiments Setup:** To study the benefits of LCP models versus only model-parallelism methods, we deploy several models on a distributed system of Raspberry Pi 3s (RPis), the specifications in Table 4. On each RPi, with the Ubuntu 16.04 operating system, we use TensorFlow [79] and Apache Avro [80], a remote procedure call (RPC) and data serialization framework, for communication between RPis. We measure
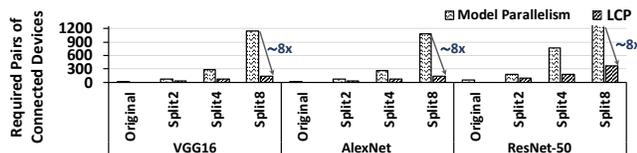


Figure 11: Communication reduction with LCP models compared to model parallelism (required pairs of connections).

power using a USB digital multimeter [81]. A local WiFi network with the measured bandwidth of 62.24 Mbps and a measured client-to-client latency of 8.83 ms for 64 B is used. All the real-world experiments are full-system measurements with all overheads included without any simulations/estimations.

**RPi Performance & Energy:** Figure 12 presents latency of inference per image on RPis. On a single device, AlexNet has 2.8 seconds latency, while VGG16 achieves 9.4 seconds latency. By deploying model-parallelism variants of the models on four and eight RPis, we achieve a maximum of 0.42s latency, a 6.6x increase, for AlexNet. But, for VGG16, on four RPis, we observe a slowdown, which is caused by high communication latency. LCP variants of split4 and split8 can reach up to 115 ms and 400 ms latency per image for AlexNet and VGG16, respectively. This is because LCP models are lightweight and parallelizable and have low communication. Figure 13 shows measured energy per inference for RPi implementations. To compare with previous related work, SplitNet [57], Figure 12 presents the performance of SplitNet models for AlexNet with different configurations. As seen, the performance is worse than LCP models. This is because SplitNet creates more merging/synchronization points with its tree-structured model design. The resulting model exponentially introduces more merging/synchronization with increased depth, which also does not equally split all the layers (causing load balancing issues). Finally, SplitNet performs parallelization based on dataset semantics, which means every dataset and model needs to be manually split. §2 provided more reasons on this performance difference.

**TVM Experiments on PYNQ Boards:** As a real-world ex-

Table 4: Specification of RPi, PYNQ FPGA, and AWS.

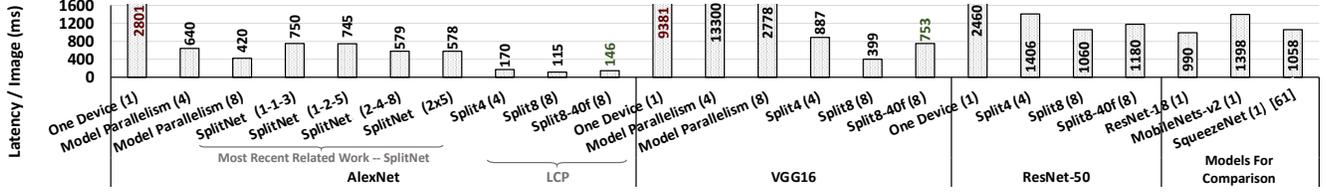| Raspberry Pi 3B+ | | | | |
|---|---|---|---|---|
| CPU | 1.2 GHz Quad Core ARM Cortex-A53 | | | |
| Memory | 1 GB LPDDR2 SDRAM @ **933Mb/s/pin** | | | |
| Die Size | ≈ 196$mm^2$ @ **28 nm** | | | |
| **Edge FPGA** (Zynq Artix 7 XC7Z020) | | | | |
| | | DSP48E | FF | LUT |
| Utilization | #Unit | 96 | 5427 | 2343 |
| | % | 44 | 5 | 4 |
| Static Power | 0.121 W | | | |
| Dynamic Power | Signals: 0.009 W | Logic: 0.003 W | | |
| **AWS** | | | | |
| AWS Instance Specification | T2.micro 1 vCPU, 1 GB Memory, 64 GB Storage | | | |

**Figure 12: Latency per image: Model-parallelism, SplitNet [57], and LCP models on RPi (number in parenthesis is #devices).**
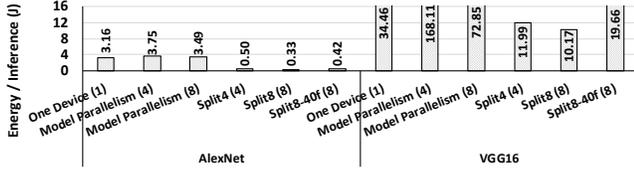


**Figure 13: All devices energy per inference: Model-parallelism, and LCP on RPi (number in parenthesis is #devices).**

ample for edge FPGA implementation, we use TVM [69] on the PYNQ [82] board. PYNQ is designed for embedded applications. We use the TVM VTA stack on the PYNQ as the architecture (RISC-style instructions) and only change the models (ResNet-18 vs. LCP ResNet-18 Split2 with <1 accuracy drop). In this way, we can measure the benefits of LCP models without relying on any special tailored hardware.

Our performance result shares the entire system pipeline performance, from a live camera feed to prediction output on two boards versus one board. Figure 14a shows a 2.7x speedup, including all communication and system overheads, network latency, and jitter because LCP models are parallelized on two devices and, in total, they have lower computation and memory footprints. The measured reduction in memory footprint is shown Figure 14b.
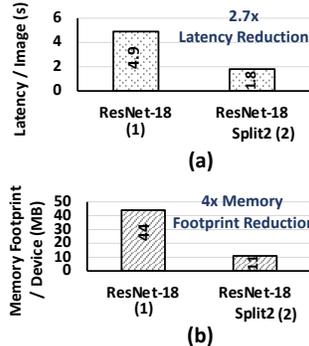


**Figure 14: TVM Experiments: (a) Latency per image, (b) memory footprint per device (number in parenthesis is #devices).**

**AWS Experiments:** To see the reduced communication and distributed execution benefits of LCP models further, we

deploy AlexNet, VGG16, and ResNet-50 models on AWS T2.micro instances with only one vCPU and 1 GB memory per instance. Figure 15 presents the derived statistics. In all cases, LCP models not only reduces the average latency but also significantly reduce maximum latency. Splits four and eight have lower speedup compared with our RPi experiments because all the 4/8 instances are not hosted on the same machine; thus, the communication cost is higher than the usual edge-specific cases that this paper targets.

## 4.3. Edge FPGA Experiments

**FPGA Experiments Setup:** We implement our tailored microarchitecture on a ZYNQ XC7Z020 FPGA targeting PYNQ-z1 boards [83]. We use Xilinx Vivado HLS for implementation and verify the functionality of our implementation using regression tests. We use relevant *#pragrma* as hints to describe our desired microarchitectures in C++. We synthesize and implement our design using Vivado and report post-implementation (*i.e.*, place & route) performance numbers and resource utilizations. Inputs and output of our design are transferred through the AXI stream interface. The clock frequency is set to 100 MHz. Communication for multiple devices is estimated with the network provided in §4.2.

**FPGA Performance:** Figure 16 shows the experiment results for our edge-tailored hardware. The latency per image is shown in Figure 16a, with improvement in communication
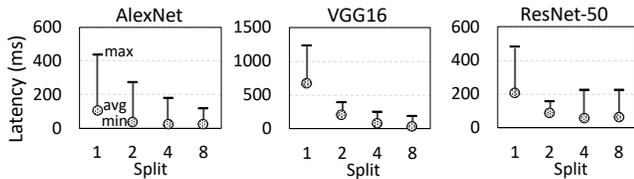


**Figure 15: Average, minimum, and maximum latencies of distributed LCP execution on AWS T2.micro instances with 1 vCPU and 1 GB memory per instance.**
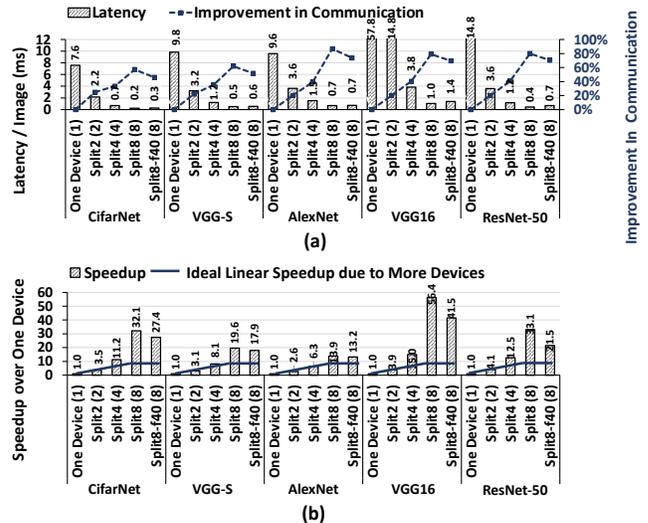


**Figure 16: Edge FPGA with tailored hardware latency and speedup: (a) Latency per image, (b) speedup over one device (number in parenthesis is #devices).**
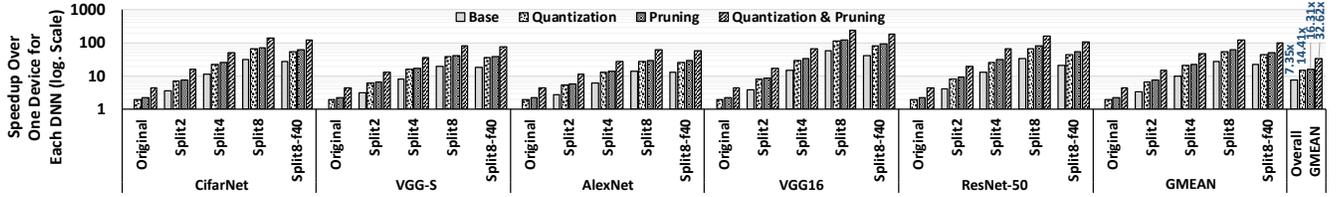
8

**Figure 17: Edge FPGA with tailored hardware speedup with quantization & pruning. Additional speedup is gained by applying lossless ($\leq 0.1\%$) quantization and structured pruning.**
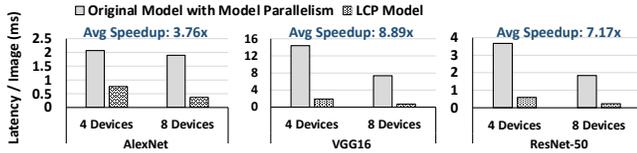


**Figure 18: Latency per image for edge FPGA with tailored hardware comparing LCP vs. model parallelism.**

overhead versus model-parallelism methods (86% and 60% for 8split and 4split). Depending on the model, the inference per latency on a single device is between 4–29ms; a 221–325x speedup compared to RPi results for AlexNet and VGG16. Our designed LCP models achieve acceptable performance for edge computing, which is 10s of inferences per second, around 1–10ms. As observed, the accuracy loss of our split-only models can be easily restored by fast split-fattened models of `f40` with a negligible performance overhead (maximum of 20 ms). Figure 16b illustrates the speedup numbers over one device. The ideal linear speedup shows the ideal scaling speedup with more available devices. As shown, we achieve superlinear speedups. An important parameter in scaling concerns how the *overheads* scale. The superlinear speedup stems from the dramatic reduction of communication overhead as parallelism increases. In traditional data and model parallelism, such overhead increases, which causes sublinear speedup. Figure 18 compares latency per image for LCP and model parallelism. On average, LCP models are 3.76x, 8.89x, and 7.17x faster than their model-parallelism counterparts for AlexNet, VGG16, and ResNet-50 (4 and 8 devices), respectively. LCP achieves a maximum and average speedups of 56x and 7x, compared to the originals (Figure 17, base bars).

**Quantization & Pruning:** As mentioned in §5 and §1, techniques that reduce the footprint of DNNs can be applied to each individual LCP branch. Basically, the target output for each LCP branch is now its pre-final activations during optimizations. We study the benefits of lossless quantization and structured pruning on top of our LCP models. Based on our experiment, with 3.13 (<integer.fraction>) quantization, our models do not lose accuracy. Similarly, applying structured pruning [84], for which systolic arrays gain benefits, reduces the size of parameters between 40%–50% per convolution layer without an accuracy drop. Other pruning algorithms increase the sparsity of the data, which is not necessarily beneficial for systolic arrays. Figure 17 presents the speedup gained from these techniques normalized to the baseline implemen-

tation for each model, the execution performance of which shown in Figure 16a. Quantization and pruning themselves, improve the performance of the original models by 1.96x and 2.2x, respectively, and 4.31x when applied together. When quantization and pruning are combined with LCP, the overall performance speedup becomes 14.41x and 16.31x, respectively. Compared to the original models, LCP + quantization and pruning achieves up to 244x speedup (VGG16-split8), and an average of 33x (across all models and variants).

### 4.4. ASIC Implementation

We implement the ASIC design of LCP using an Arizona State Predictive PDK (ASAP) 7nm technology node [48]. Our tool chain includes the Synopsys design compiler (DC) for synthesis, Cadence Innovus for place and route, and Cadence Tempus for timing and power analysis. As an input to our ASIC design, we use our same Verilog code generated by Vivado HLS. Figure 8b show the layout of our chip of size 0.107 mm$^2$ (i.e., $295\mu m \times 365\mu m$). The memory cells shown in the figure represent the FIFO buffers, used for pipelining. Figure 19 shows the power consumption of our ASIC design. The breakdown of power consummation leading to a total 16.1 mW is listed in Figure 19a. As a comparison point, Eyeriss [60] and EIE [85] consume $\approx$250 mW and $\approx$590 mW, respectively. Besides, as Figure 19b shows, power distributes uniformly, which prevents hot spot creation.

## 5. Related Work

We review related techniques used to reduce the high demands of DNNs, distributing their computation, and current efforts on DNN hardware accelerators.

**Techniques Without Changing Model Architecture:** Several techniques have been developed to reduce the computation
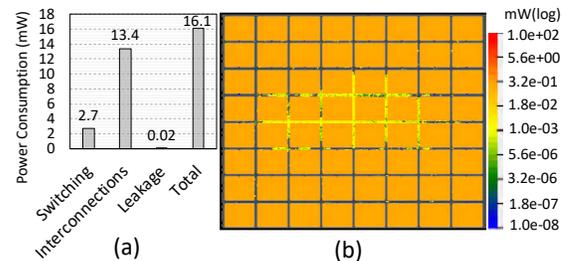


**Figure 19: Power Consumption for 7-nm ASIC Design @800MHz: (a) breakdown (b) distribution.**

and memory footprint of DNNs without changing the network architecture. For instance, pruning [32, 86–89] removes the close-to-zero weights and quantization or low-precision inference [36, 52, 53, 90, 91] change the representations of numbers, which results in simpler calculations. Other methods partition resources [92, 93] or binarize the weights [94–96]. Binarizing weights hurts accuracy. *The aforementioned techniques are orthogonal to our work and can be applied to each branch to further reduce the computational and memory costs (§4.3).*

**Techniques That Change Model Architecture:** With the prevalence of IoT and edge devices, specific frameworks such as ELL library [50] (see Figure 2) by Microsoft and Tensorflow Lite [97] have been developed by industry. Other proposals developed mobile-specific models [54, 77, 98–100] by handcrafting more efficient operations or models to reduce the number of parameters [54], create efficient operation to minimize computation density [98], or use resource-efficient connections [100]. Unlike LCP models, all these models have a single chain of dependency [58] that prevents efficient parallelism. Moreover, several of the models trade off the state-of-the-art accuracy with efficiency [100]. SplitNet [57] is one the few papers that focuses on higher parallelizability of models (evaluation on ResNet and AlexNet), but relying on dataset semantics creates imbalanced branches and the method is invariant to number of devices, as discussed in §2. Recently, with the growing interest in automating the design process [58, 101–103], learning new networks for mobiles has also gained attention by integrating the constraints of mobile platforms (*i.e.*, latency). These attempts are still limited to single-device execution. *In summary, these studies (1) have a high design cost, (i.e., they target only one specific model and dataset without extendibility); (2) target single mobile platforms; and (3) do not consider inter-layer layer parallelism and communication challenges.*

**Distributing DNN Inference Computations:** With large DNN models, distributing a single model has gained the attention of researchers [5, 14, 35, 104, 105]. Usually, the distribution is done in a high-performance computing domain with different goals in mind. In the resource-constrained edge devices, Neurosurgeon [105] dynamically partitions a DNN model between a *single* edge device and the cloud. DDNN [104] partitions the model between edge devices and the cloud but uses data parallelism. Hadidi et al. [5, 14, 106–109] investigate the distribution in edge with model-parallelism methods, showing the effect of the communication barrier in distributing by the diminishing return in performance with a large number of devices. *LCP models go beyond model parallelism methods, which was not the focus of the above studied, and enable efficient distribution that is not examined in the above studies.*

## 6. Discussions

**Intuition Behind LCP:** We conjecture that LCP models provide good performance because (1) independent branches learn complex non-overlapping features independently within a small search space, whereas original models need to create the same complex features from a higher dimension feature search. We observe that each branch eventually learns an almost disjoint feature representation; (2) In split models, gradient descent updates are more efficient in reaching early layers compared to the original models due to fewer number of parameters in their route.

**Extension to New Models:** We studied ResNets and MobileNet, which are still widely used models. Other models represent sequential DNNs that serve as the basis to confirm our method and are still used in robotics. Newer models such as EfficientNet [110] and MobileNetv3 [111] that use novel blocks such as Bottleneck or Squeeze & Excitation can be represented with convolution, fully connected, and basic matrix multiplications. All of which can be parallelized by LCP.

**System-Level Choices:** LCP is in conjunction with other technologies available today. LCP does not replace these technologies, but rather enables exploitation of local edge devices to enable intelligence in the edge [112, 113]. In a few cases, relying on cloud-based offloading for accuracy-critical tasks is necessary (*e.g.*, finding a specific license plate), whereas, in several others (*e.g.*, counting the cars passing an intersection) the system must rely on cloud or high-performance systems.

**SqueezeNet:** SqueezeNet [54] achieves an accuracy similar to that of AlexNet with fewer parameters by using compute-heavy Fire modules. SqueezeNet trades off parameters with computations, and requires 860M MAC operations, whereas our distributed AlexNet requires only 240M MAC operations. We also observe a 12x increase in the number of activations from 12.58 M in SqueezeNet vs. 1.39 M in AlexNet.

**Skip/Residual Connections:** LCP procedure similarly applies to more complex models with residual and skip connections as shown for ResNets in §4. Simply put, each branch has similar connections but with smaller depth.

**Alleviating Large Memory Footprints:** Sometimes large memory footprints are necessary and access to the next levels of the storage system is enforced. In our design (§3.2), such accesses do not cause slowdown because data is stored in sequential addresses (*i.e.*, streaming), and we overlap data transfer and computations for independent elements.

**Memory Layout Preprocessing:** Our simple algorithm to change the storage format is in $O(N)$ (§3.2(4)). Therefore, the host preprocessing for reordering the data can be done during writing the data to the memory with a single pass.

## 7. Conclusions

We proposed low-communication parallelization (LCP) models, designed for efficient in-the-edge distribution. LCP models optimize communication while reducing memory and computation by utilizing several narrow independent branches. We presented our results on the accuracy of LCP models. We build a systolic architecture for edge computing both on FPGA and ASIC. Finally, our results on RPis, edge-based FPGAs, AWS instances confirms the benefits.

# References

[1] Alessandro Giusti, Jérôme Guzzi, Dan C Cireşan, Fang-Lin He, Juan P Rodríguez, Flavio Fontana, Matthias Faessler, Christian Forster, Jürgen Schmidhuber, Gianni Di Caro, et al. A machine learning approach to visual perception of forest trails for mobile robots. *IEEE Robotics and Automation Letters*, 1(2):661–667, 2016.

[2] Mark Pfeiffer, Michael Schaeuble, Juan Nieto, Roland Siegwart, and Cesar Cadena. From perception to decision: A data-driven approach to end-to-end motion planning for autonomous ground robots. In *2017 ieee international conference on robotics and automation (icra)*, pages 1527–1533. IEEE, 2017.

[3] Peter Corcoran and Soumya Kanti Datta. Mobile-edge computing and the internet of things for consumers: Extending cloud computing and services to the edge of the network. *IEEE Consumer Electronics Magazine*, 5(4):73–74, 2016.

[4] Manuela Veloso, Joydeep Biswas, Brian Coltin, Stephanie Rosenthal, Tom Kollar, Cetin Mericli, Mehdi Samadi, Susana Brandao, and Rodrigo Ventura. Cobots: Collaborative robots servicing multi-floor buildings. In *2012 IEEE/RSJ international conference on intelligent robots and systems*, pages 5446–5447. IEEE, 2012.

[5] Ramyad Hadidi, Jiashen Cao, Matthew Woodward, Michael S Ryoo, and Hyesoon Kim. Distributed perception by collaborative robots. *IEEE Robotics and Automation Letters*, 3(4):3709–3716, 2018.

[6] Arti Singh, Baskar Ganapathysubramanian, Asheesh Kumar Singh, and Soumik Sarkar. Machine learning for high-throughput stress phenotyping in plants. *Trends in plant science*, 21(2):110–124, 2016.

[7] Huimin Lu, Yujie Li, Shenglin Mu, Dong Wang, Hyoungseop Kim, and Seiichi Serikawa. Motor anomaly detection for unmanned aerial vehicles using reinforcement learning. *IEEE internet of things journal*, 5(4):2315–2322, 2018.

[8] Zhangjie Fu, Yuanhang Mao, Daojing He, Jingnan Yu, and Guowu Xie. Secure multi-uav collaborative task allocation. *IEEE Access*, 7:35579–35587, 2019.

[9] Nader Mohamed, Jameela Al-Jaroodi, Imad Jawhar, and Sanja Lazarova-Molnar. A service-oriented middleware for building collaborative uavs. *Journal of Intelligent & Robotic Systems*, 74(1-2):309–321, 2014.

[10] Shuochao Yao, Yiran Zhao, Aston Zhang, Shaohan Hu, Huajie Shao, Chao Zhang, Lu Su, and Tarek Abdelzaher. Deep learning for the internet of things. *Computer*, 51(5):32–41, 2018.

[11] Omer Berat Sezer, Erdogan Dogdu, and Ahmet Murat Ozbayoglu. Context-aware computing, learning, and big data in internet of things: a survey. *IEEE Internet of Things Journal*, 5(1):1–27, 2018.

[12] He Li, Kaoru Ota, and Mianxiong Dong. Learning iot in edge: Deep learning for the internet of things with edge computing. *IEEE Network*, 32(1):96–101, 2018.

[13] Tuyen X Tran, Abolfazl Hajisami, Parul Pandey, and Dario Pompili. Collaborative mobile edge computing in 5g networks: New paradigms, scenarios, and challenges. *IEEE Communications Magazine*, 55(4):54–61, 2017.

[14] Ramyad Hadidi, Jiashen Cao, Michael S Ryoo, and Hyesoon Kim. Towards collaborative inferencing of deep neural networks on internet of things devices. *IEEE Internet of Things Journal*, 2020.

[15] Luigi Alfredo Grieco, Alessandro Rizzo, Simona Colucci, Sabrina Sicari, Giuseppe Piro, Donato Di Paola, and Gennaro Boggia. Iot-aided robotics applications: Technological implications, target domains and open issues. *Computer Communications*, 54:32–47, 2014.

[16] Milan Erdelj, Michał Król, and Enrico Natalizio. Wireless sensor networks and multi-uav systems for natural disaster management. *Computer Networks*, 124:72–86, 2017.

[17] Markus Quaritsch, Emil Stojanovski, Christian Bettstetter, Gerhard Friedrich, Hermann Hellwagner, Bernhard Rinner, Michael Hofbaur, and Mubarak Shah. Collaborative microdrones: applications and research challenges. In *Proceedings of the 2nd International Conference on Autonomic Computing and Communication Systems*, pages 1–7, 2008.

[18] Nathan Michael, Shaojie Shen, Kartik Mohta, Vijay Kumar, Keiji Nagatani, Yoshito Okada, Seiga Kiribayashi, Kazuki Otake, Kazuya Yoshida, Kazunori Ohno, et al. Collaborative mapping of an earthquake damaged building via ground and aerial robots. In *Field and service robotics*, pages 33–47. Springer, 2014.

[19] Avital Bechar and Clément Vigneault. Agricultural robots for field operations: Concepts and components. *Biosystems Engineering*, 149:94–111, 2016.

[20] H Anil, KS Nikhil, V Chaitra, and BS Guru Sharan. Revolutionizing farming using swarm robotics. In *2015 6th International Conference on Intelligent Systems, Modelling and Simulation*, pages 141–147. IEEE, 2015.

[21] Y Baudoin and Maki K Habib. *Using robots in hazardous environments: Landmine detection, de-mining and other applications*. Elsevier, 2010.

[22] Marcos Dias de Assuncao, Alexandre da Silva Veith, and Rajkumar Buyya. Distributed data stream processing and edge computing: A survey on resource elasticity and future directions. *Journal of Network and Computer Applications*, 103:1–17, 2018.

[23] Stuart Golodetz, Tommaso Cavallari, Nicholas A Lord, Victor A Prisacariu, David W Murray, and Philip HS Torr. Collaborative large-scale dense 3d reconstruction with online inter-agent pose optimisation. *IEEE transactions on visualization and computer graphics*, 24(11):2895–2905, 2018.

[24] Binita Gupta. Discovering cloud-based services for iot devices in an iot network associated with a user, June 4 2015. US Patent App. 14/550,595.

[25] Hui Li and Xiaojiang Xing. Internet of things service architecture and method for realizing internet of things service, March 17 2015. US Patent 8,984,113.

[26] Inc. Gartner. Gartner says 6.4 billion connected "things" will be in use in 2016, up 30 percent from 2015. https://www.gartner.com/newsroom/id/3165317, 2015. [Online; accessed 04/10/20].

[27] Jayavardhana Gubbi, Rajkumar Buyya, Slaven Marusic, and Marimuthu Palaniswami. Internet of things (iot): A vision, architectural elements, and future directions. *Future generation computer systems*, 29(7):1645–1660, 2013.

[28] Ben Zhang, Nitesh Mor, John Kolb, Douglas S Chan, Ken Lutz, Eric Allman, John Wawrzynek, Edward Lee, and John Kubiatowicz. The cloud is not enough: Saving iot from the cloud. In *7th {USENIX} Workshop on Hot Topics in Cloud Computing (HotCloud 15)*, 2015.

[29] Shancang Li, Li Da Xu, and Shanshan Zhao. The internet of things: a survey. *Information Systems Frontiers*, 17(2):243–259, 2015.

[30] F Biscotti, J Skorupa, R Contu, et al. The impact of the internet of things on data centers. *Gartner Research*, 18, 2014.

[31] In Lee and Kyoochun Lee. The internet of things (iot): Applications, investments, and challenges for enterprises. *Business Horizons*, 58(4):431–440, 2015.

[32] Song Han, Huizi Mao, and William J Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and huffman coding. In *4th International Conference on Learning Representations*. ACM, 2016.

[33] Ramyad Hadidi, Jiashen Cao, Yilun Xie, Bahar Asgari, Tushar Krishna, and Hyesoon Kim. Characterizing the deployment of deep neural networks on commercial edge devices. In *Proceedings of IEEE International Symposium on Workload Characterization*, 2019.

[34] Alex Krizhevsky, Ilya Sutskever, and Geoffrey E Hinton. Imagenet classification with deep convolutional neural networks. In *26th Annual Conference on Neural Information Processing Systems (NIPS)*, pages 1097–1105. ACM, 2012.

[35] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Andrew Senior, Paul Tucker, Ke Yang, Quoc V Le, et al. Large scale distributed deep networks. In *NIPS'12*, pages 1223–1231. ACM, 2012.

[36] Yunchao Gong, Liu Liu, Ming Yang, and Lubomir Bourdev. Compressing deep convolutional networks using vector quantization. *arXiv preprint arXiv:1412.6115*, 2014.

[37] Yann LeCun. The mnist database of handwritten digits. *http://yann.lecun.com/exdb/mnist/*, 1998.

[38] Alex Krizhevsky and Geoffrey Hinton. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.

[39] M-E. Nilsback and A. Zisserman. Automated flower classification over a large number of classes. In *Proceedings of the Indian Conference on Computer Vision, Graphics and Image Processing*, Dec 2008.

[40] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *IJCV*, 115(3):211–252, 2015.

[41] Peter Mattson, Vijay Janapa Reddi, Christine Cheng, Cody Coleman, Greg Diamos, David Kanter, Paulius Micikevicius, David Patterson, Guenther Schmuelling, Hanlin Tang, et al. Mlperf: An industry standard benchmark suite for machine learning performance. *IEEE Micro*, 40(2):8–16, 2020.

[42] Vladimir Vujović and Mirjana Maksimović. Raspberry pi as a sensor web node for home automation. *Computers & Electrical Engineering*, 44:153–171, 2015.

[43] Richard Grimmett. *Raspberry Pi robotics projects*. Packt Publishing Ltd, 2015.

[44] Alan G Millard, Russell Joyce, James A Hilder, Cristian Fleşeriu, Leonard Newbrook, Wei Li, Liam J McDaid, and David M Halliday. The pi-puck extension board: a raspberry pi interface for the e-puck robot platform. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 741–748. IEEE, 2017.

[45] Isaiah Brand, Josh Roy, Aaron Ray, John Oberlin, and Stefanie Oberlix. Pidrone: An autonomous educational drone using raspberry pi and python. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 1–7. IEEE, 2018.

[46] Sean Wilson, Ruben Gameros, Michael Sheely, Matthew Lin, Kathryn Dover, Robert Gevorkyan, Matt Haberland, Andrea Bertozzi, and Spring Berman. Pheeno, a versatile swarm robotic research and education platform. *IEEE Robotics and Automation Letters*, 1(2):884–891, 2016.

[47] Norman P Jouppi, Cliff Young, Nishant Patil, David Patterson, Gaurav Agrawal, Raminder Bajwa, Sarah Bates, Suresh Bhatia, Nan Boden, Al Borchers, et al. In-datacenter performance analysis of a tensor processing unit. In *2017 ACM/IEEE 44th Annual International Symposium on Computer Architecture (ISCA)*, pages 1–12. IEEE, 2017.

[48] Lawrence T Clark, Vinay Vashishtha, Lucian Shifren, Aditya Gujja, Saurabh Sinha, Brian Cline, Chandarasekaran Ramamurthy, and Greg Yeric. Asap7: A 7-nm finfet predictive process design kit. *Microelectronics Journal*, 53:105–115, 2016.

[49] Swagath Venkataramani, Ashish Ranjan, Subarno Banerjee, Dipankar Das, Sasikanth Avancha, Ashok Jagannathan, Ajaya Durg, Dheemanth Nagaraj, Bharat Kaul, Pradeep Dubey, et al. Scaledeep: A scalable compute architecture for learning and evaluating deep networks. In *ISCA'17*, pages 13–26. ACM, 2017.

[50] Microsoft. Embedded learning library (ell). https://microsoft.github.io/ELL/, 2017. [Online; accessed 04/10/20].

[51] Ofer Dekel - Microsoft Research. Compiling ai for the edge. *SysML 2019 Keynote*, 2019.

[52] Vincent Vanhoucke, Andrew Senior, and Mark Z Mao. Improving the speed of neural networks on cpus. In *Proceeding Deep Learning and Unsupervised Feature Learning NIPS Workshop*, volume 1, page 4. ACM, 2011.

[53] Darryl Lin, Sachin Talathi, and Sreekanth Annapureddy. Fixed point quantization of deep convolutional networks. In *International Conference on Machine Learning*, pages 2849–2858, 2016.

[54] Forrest N Iandola, Song Han, Matthew W Moskewicz, Khalid Ashraf, William J Dally, and Kurt Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5 mb model size. *arXiv preprint arXiv:1602.07360*, 2016.

[55] Sharan Chetlur, Cliff Woolley, Philippe Vandermersch, Jonathan Cohen, John Tran, Bryan Catanzaro, and Evan Shelhamer. Cudnn: Efficient primitives for deep learning. *arXiv preprint arXiv:1410.0759*, 2014.

[56] Stefan Hadjis, Firas Abuzaid, Ce Zhang, and Christopher Ré. Caffe con troll: Shallow ideas to speed up deep learning. In *Proceedings of the Fourth Workshop on Data analytics in the Cloud*, page 2. ACM, 2015.

[57] Juyong Kim, Yookoon Park, Gunhee Kim, and Sung Ju Hwang. Splitnet: Learning to semantically split deep networks for parameter reduction and model parallelization. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1866–1874. JMLR. org, 2017.

[58] Saining Xie, Alexander Kirillov, Ross Girshick, and Kaiming He. Exploring randomly wired neural networks for image recognition. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1284–1293, 2019.

[59] Song Han, Junlong Kang, Huizi Mao, Yiming Hu, Xin Li, Yubin Li, Dongliang Xie, Hong Luo, Song Yao, Yu Wang, et al. Ese: Efficient speech recognition engine with sparse lstm on fpga. In *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 75–84. ACM, 2017.

[60] Yu-Hsin Chen, Tushar Krishna, Joel S Emer, and Vivienne Sze. Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks. *IEEE Journal of Solid-State Circuits*, 52(1):127–138, 2017.

[61] Zidong Du, Robert Fasthuber, Tianshi Chen, Paolo Ienne, Ling Li, Tao Luo, Xiaobing Feng, Yunji Chen, and Olivier Temam. Shidiannao: Shifting vision processing closer to the sensor. In *ACM SIGARCH Computer Architecture News*, volume 43, pages 92–104. ACM, 2015.

[62] Jeff Dean. Machine learning for systems and systems for machine learning, 2017.

[63] Bahar Asgari, Ramyad Hadidi, Hyesoon Kim, and Sudhakar Yalamanchili. Eridanus: Efficiently running inference of dnns using systolic arrays. *IEEE Micro*, 39(5):46–54, 2019.

[64] Hsiang-Tsung Kung. Why systolic architectures? *IEEE computer*, 15(1):37–46, 1982.

[65] JEDEC. Jedec standard: Low power double data rate 2 (lpddr2). https://www.jedec.org/sites/default/files/docs/JESD209-2B.pdf, 2019. [Online; accessed 04/10/20].

[66] Chen Zhang, Peng Li, Guangyu Sun, Yijin Guan, Bingjun Xiao, and Jason Cong. Optimizing fpga-based accelerator design for deep convolutional neural networks. In *Proceedings of the 2015 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, pages 161–170. ACM, 2015.

[67] Xiaofan Zhang, Junsong Wang, Chao Zhu, Yonghua Lin, Jinjun Xiong, Wen-mei Hwu, and Deming Chen. Dnnbuilder: an automated tool for building high-performance dnn hardware accelerators for fpgas. In *2018 IEEE/ACM International Conference on Computer-Aided Design (ICCAD)*, pages 1–8. IEEE, 2018.

[68] Yunxuan Yu, Chen Wu, Tiandong Zhao, Kun Wang, and Lei He. Opu: An fpga-based overlay processor for convolutional neural networks. *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, 2019.

[69] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Haichen Shen, Eddie Q Yan, Leyuan Wang, Yuwei Hu, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. Tvm: end-to-end optimization stack for deep learning. *arXiv preprint arXiv:1802.04799*, pages 1–15, 2018.

[70] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML'17*, pages 448–456. ACM, 2015.

[71] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*, 2014.

[72] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.

[73] Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. In *3rd International Conference on Learning Representations*. ACM, 2015.

[74] Alex Krizhevsky. One weird trick for parallelizing convolutional neural networks. *arXiv preprint arXiv:1404.5997*, 2014.

[75] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *CVPR'16*, pages 770–778. IEEE, 2016.

[76] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[77] Gao Huang, Shichen Liu, Laurens Van der Maaten, and Kilian Q Weinberger. Condensenet: An efficient densenet using learned group convolutions. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2752–2761, 2018.

[78] Joseph Redmon. Darknet: Open source neural networks in c. pjreddie.com/darknet, 2013–2016.

[79] Martín Abadi et al. TensorFlow: Large-scale machine learning on heterogeneous systems, 2015. Software available from tensorflow.org.

[80] The Apache Software Foundation. Apache avro. https://avro.apache.org, 2017. [Online; accessed 04/10/20].

[81] Makerhawk. Um25c usb power meter. makerhawk.com, 2019. [Online; accessed 04/10/20].

[82] Xilinx Inc. Pynq: Python productivity for zynq. pynq.io, 2019. [Online; accessed 04/10/20].

[83] Younmin Bae, Ramyad Hadidi, Bahar Asgari, Jiashen Cao, and Hyesoon Kim. Capella: Customizing perception for edge devices by efficiently allocating fpgas to dnns. In *2019 29th International Conference on Field Programmable Logic and Applications (FPL)*, pages 421–421. IEEE, 2019.

[84] Sajid Anwar, Kyuyeon Hwang, and Wonyong Sung. Structured pruning of deep convolutional neural networks. *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, 13(3):32, 2017.

[85] Song Han, Xingyu Liu, Huizi Mao, Jing Pu, Ardavan Pedram, Mark A Horowitz, and William J Dally. Eie: efficient inference engine on compressed deep neural network. In *2016 ACM/IEEE 43rd Annual International Symposium on Computer Architecture (ISCA)*, pages 243–254. IEEE, 2016.

[86] Jiecao Yu, Andrew Lukefahr, David Palframan, Ganesh Dasika, Reetuparna Das, and Scott Mahlke. Scalpel: Customizing dnn pruning to the underlying hardware parallelism. In *44th International Symposium on Computer Architecture (ISCA)*, pages 548–560. IEEE, 2017.

[87] Ji Lin, Yongming Rao, Jiwen Lu, and Jie Zhou. Runtime neural pruning. In *Advances in Neural Information Processing Systems (NIPS)*, pages 2181–2191, 2017.

[88] Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks. In *Advances in neural information processing systems*, pages 2074–2082, 2016.

[89] Bahar Asgari, Ramyad Hadidi, Hyesoon Kim, and Sudhakar Yalamanchili. Lodestar: Creating locally-dense cnns for efficient inference on systolic arrays. In *Proceedings of the 56th Annual Design Automation Conference 2019*, pages 1–2, 2019.

[90] Matthieu Courbariaux, Yoshua Bengio, and Jean-Pierre David. Training deep neural networks with low precision multiplication. *arXiv preprint arXiv:1412.7024*, 2014.

[91] Urs Köster, Tristan Webb, Xin Wang, Marcel Nassar, Arjun K Bansal, William Constable, Oguz Elibol, Scott Gray, Stewart Hall, Luke Hornof, et al. Flexpoint: An adaptive numerical format for efficient training of deep neural networks. In *Advances in Neural Information Processing Systems (NIPS)*, pages 1742–1752, 2017.

[92] Yongming Shen, Michael Ferdman, and Peter Milder. Maximizing CNN accelerator efficiency through resource partitioning. In *44th International Symposium on Computer Architecture (ISCA)*. IEEE, 2017.

[93] Jianxin Guo, Shouyi Yin, Peng Ouyang, Leibo Liu, and Shaojun Wei. Bit-width based resource partitioning for cnn acceleration on fpga. In *25th Annual IEEE International Symposium on. Field- Programmable Custom Computing Machines*, 2017.

[94] Fengfu Li, Bo Zhang, and Bin Liu. Ternary weight networks. *arXiv preprint arXiv:1605.04711*, 2016.

[95] Matthieu Courbariaux, Itay Hubara, Daniel Soudry, Ran El-Yaniv, and Yoshua Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or- 1. *arXiv preprint arXiv:1602.02830*, 2016.

[96] Mohammad Rastegari, Vicente Ordonez, Joseph Redmon, and Ali Farhadi. Xnor-net: Imagenet classification using binary convolutional neural networks. In *ECCV'16*, pages 525–542. Springer, 2016.

[97] Google. Introduction to tensorflow lite. https://www.tensorflow.org/mobile/tflite/, 2017. [Online; accessed 04/10/20].

[98] Andrew G Howard, Menglong Zhu, Bo Chen, Dmitry Kalenichenko, Weijun Wang, Tobias Weyand, Marco Andreetto, and Hartwig Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. *arXiv preprint arXiv:1704.04861*, 2017.

[99] Xiangyu Zhang, Xinyu Zhou, Mengxiao Lin, and Jian Sun. Shufflenet: An extremely efficient convolutional neural network for mobile devices. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 6848–6856, 2018.

[100] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4510–4520, 2018.

[101] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V Le. Learning transferable architectures for scalable image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 8697–8710, 2018.

[102] Barret Zoph and Quoc V Le. Neural architecture search with reinforcement learning. 2016.

[103] Bowen Baker, Otkrist Gupta, Nikhil Naik, and Ramesh Raskar. Designing neural network architectures using reinforcement learning. *arXiv preprint arXiv:1611.02167*, 2016.

[104] Surat Teerapittayanon, Bradley McDanel, and HT Kung. Distributed deep neural networks over the cloud, the edge and end devices. In *37th IEEE International Conference on Distributed Computing Systems (ICDCS)*, pages 328–339. IEEE, 2017.

[105] Yiping Kang, Johann Hauswald, Cao Gao, Austin Rovinski, Trevor Mudge, Jason Mars, and Lingjia Tang. Neurosurgeon: Collaborative intelligence between the cloud and mobile edge. In *22nd ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, pages 615–629. ACM, 2017.

[106] Ramyad Hadidi, Jiashen Cao, Matthew Woodward, Michael S Ryoo, and Hyesoon Kim. Musical chair: Efficient real-time recognition using collaborative iot devices. *arXiv preprint arXiv:1802.02138*, 2018.

[107] Jiashen Cao, Fei Wu, Ramyad Hadidi, Lixing Liu, Tushar Krishna, Micheal S Ryoo, and Hyesoon Kim. An edge-centric scalable intelligent framework to collaboratively execute dnn. In *Demo for SysML Conference, Palo Alto, CA*, 2019.

[108] Ramyad Hadidi, Jiashen Cao, Matthew Woodward, Michael S Ryoo, and Hyesoon Kim. Real-time image recognition using collaborative iot devices. In *Proceedings of the 1st on Reproducible Quality-Efficient Systems Tournament on Co-designing Pareto-efficient Deep Learning*, page 1. 2018.

[109] Ramyad Hadidi, Jiashen Cao, Micheal S Ryoo, and Hyesoon Kim. Collaborative execution of deep neural networks on internet of things devices. *arXiv preprint arXiv:1901.02537*, 2019.

[110] Mingxing Tan and Quoc Le. Efficientnet: Rethinking model scaling for convolutional neural networks. In *International Conference on Machine Learning*, pages 6105–6114, 2019.

[111] Andrew Howard, Mark Sandler, Grace Chu, Liang-Chieh Chen, Bo Chen, Mingxing Tan, Weijun Wang, Yukun Zhu, Ruoming Pang, Vijay Vasudevan, et al. Searching for mobilenetv3. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1314–1324, 2019.

[112] Zhi Zhou, Xu Chen, En Li, Liekang Zeng, Ke Luo, and Junshan Zhang. Edge intelligence: Paving the last mile of artificial intelligence with edge computing. *Proceedings of the IEEE*, 107(8):1738–1762, 2019.

[113] Ramyad Hadidi, Jiashen Cao, Michael S Ryoo, and Hyesoon Kim. Robustly executing dnns in iot systems using coded distributed computing. In *Proceedings of the 56th Annual Design Automation Conference 2019*, pages 1–2, 2019.