# Tailor: Altering Skip Connections for Resource-Efficient Inference

OLIVIA WENG and GABRIEL MARCANO, University of California San Diego, USA

VLADIMIR LONCAR, Massachusetts Institute of Technology, USA

ALIREZA KHODAMORADI, AMD, USA

ABARAJITHAN G and NOJAN SHEYBANI, University of California San Diego, USA

ANDRES MEZA and FARINAZ KOUSHANFAR, University of California San Diego, USA

KRISTOF DENOLF, AMD, USA

JAVIER MAURICIO DUARTE and RYAN KASTNER, University of California San Diego, USA

Deep neural networks use skip connections to improve training convergence. However, these skip connections are costly in hardware, requiring extra buffers and increasing on- and off-chip memory utilization and bandwidth requirements. In this paper, we show that skip connections can be optimized for hardware when tackled with a hardware-software codesign approach. We argue that while a network's skip connections are needed for the network to learn, they can later be removed or shortened to provide a more hardware efficient implementation with minimal to no accuracy loss. We introduce Tailor, a codesign tool whose hardware-aware training algorithm gradually removes or shortens a fully trained network's skip connections to lower their hardware cost. Tailor improves resource utilization by up to 34% for BRAMs, 13% for FFs, and 16% for LUTs for on-chip, dataflow-style architectures. Tailor increases performance by 30% and reduces memory bandwidth by 45% for a 2D processing element array architecture.

CCS Concepts: • **Hardware → Hardware-software codesign**; • **Computer systems organization → Neural networks**.

Additional Key Words and Phrases: Hardware-software co-design, neural networks

## 1 INTRODUCTION

Convolutional neural networks (NNs) often rely on skip connections—identity functions that combine the outputs of different layers—to improve training convergence [17, 45]. Skip connections help mitigate the vanishing gradient problem [4, 15] that occurs when training large CNNs, which helps increase the network's accuracy. Skip connections allow NNs to have fewer filters/weights than architectures that lack skip connections [17], such as VGG [43].

However, skip connections are generally detrimental to hardware efficiency. They have an irregular design that is ill-suited for hardware acceleration. This is due to their long lifetimes, which span several NN layers, increasing memory utilization and bandwidth requirements. This is particularly true in ResNets [17], which introduced skip connections that spanned across five layers: two convolutions, two batch normalizations (BNs), and a ReLU activation [16, 34] (see Fig. 1a). The skip connection involves minimal computation—it is either the identity or a 1×1 convolutional layer

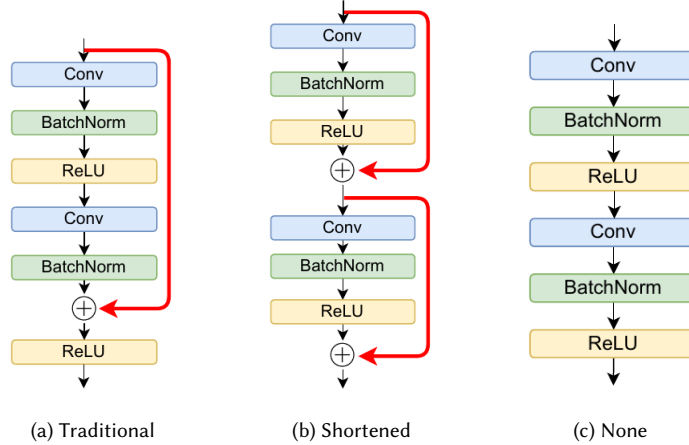|                   |                  |                 |
| :---------------: | :--------------: | :-------------: |
| (a) Traditional   | (b) Shortened    | (c) None        |

Fig. 1. Neural networks with traditional skip connections, like ResNet (a), have inefficient hardware implementations because the skip connection data must be preserved in memory during five layers of computation. This irregular topology increases memory resources and bandwidth. A more regular topology with reduced skip connection lifetimes would use fewer resources. TAILOR achieves this by shortening skip connections (b) or by eliminating them completely (c). Skip connections are in red.

for scaling—but it extends the necessary lifespan of the input data. Thus, we must store skip connection data for the duration of time needed to compute the five NN layers. In total, a model's skip connection data accounts for ~10% of its memory bandwidth either on or off chip. Buffering skip connections on chip increases on-chip memory utilization, whereas moving them off chip not only increases off-chip memory bandwidth but also requires extra control logic for scheduling [29, 30].

Optimizing skip connections requires careful *hardware-software codesign*. Skip connections are crucial for model convergence; naively removing them to reduce hardware resources leads to low accuracy [32, 50]. Instead, we must codesign how the model is (1) trained and (2) implemented in hardware to achieve a model that is both accurate and resource-efficient.

We develop TAILOR, a codesign method that gradually alters a NN's skip connections during training to produce a highly accurate and resource-efficient NN. Our results in Sec. 4 show that TAILOR can remove or shorten skip connections to achieve topologically regular NNs (Fig. 1b and 1c) that substantially reduce hardware resources, reduce memory bandwidth, and increase performance with minimal to no accuracy loss.

TAILOR takes an *existing* pre-trained model and reduces the hardware complexity of its skip connections with minimal to no accuracy loss. Moreover, TAILOR exploits the flexiblity of the FPGA architecture to customize the skip connection memories, which is not possible on a GPU or CPU. TAILOR accomplishes this dynamically during retraining in one of two ways: (1) SKIPREMOVER removes the skip connections altogether (Fig. 1c) to eliminate all associated hardware complications or (2) SKIPSHORTENER shortens each skip connection by splitting it into multiple shorter ones (Fig. 1b).

We evaluate TAILOR's applicability and benefit on ResNets [17, 18] and QuartzNets [23]—two important classes of NNs that contain skip connections of varying lengths. We also study implementing skip connections with an on-chip, dataflow-style FPGA architecture using hls4ml [2, 12] and a 2D array of multiply-accumulate processing elements. TAILOR reduces resource utilization of hls4ml architectures by up to 34% for BRAMs, 13% for FFs, and 16% for LUTs. TAILOR increases the performance of 2D array architecture by 30% and reduces memory bandwidth by 45%..

Tailor's hardware-software codesign approach reduces hardware complexity and resources by altering skip connections dynamically during retraining. Our contributions are:

- the Tailor software methodology of removing or shortening skip connections from existing NNs with minimal to no loss in accuracy,
- the Tailor hardware designs that exploit FPGA-specific architecture optimizations, which are not possible on GPU/CPU, to produce less resource-intensive skip connection implementations,
- experiments demonstrating that SkipShortener and SkipRemover models are implemented more efficiently with better performance and resource utilization than their traditional skip connection counterparts,
- and public release of the Tailor hardware-software codesign framework [1].

In Sec. 2, we review related work. In Sec. 3, we explain how Tailor's NN alterations optimize the hardware architecture. We then describe Tailor's two training methods, SkipRemover and SkipShortener, that alter skip connections with little to no loss in accuracy. Sec. 4 provides training, quantization, and hardware results for SkipRemover and SkipShortener. Sec. 5 discusses the tradeoffs Tailor presents between accuracy and hardware resource reductions. Sec. 6 concludes the paper.

## 2 BACKGROUND

### 2.1 Removing Skip Connections

While skip connection removal has been studied before [8, 25, 32, 50, 51], prior work is lacking in several ways: (1) preliminary work [32, 50, 51] only studies shallow models (up to 34 layers); (2) Li et al. [25] do not remove all of the skip connections in the models they evaluate; (3) Ding et al. [8] and Li et al. [25] both have limited architectural evaluations (e.g., GPU & mobile) that do not consider the highly customized skip connections memories enabled by FPGAs; and (4) Ding et al. [8] require starting with an entirely new NN topology whose skip connections are removable.

Monti et al. [32] start with a standard ResNet and introduce a new training method. This method uses an objective function that penalizes the skip connections and phases them out by the end of the training. This technique has only been applied to smaller ResNets (18 to 34 layers) with a small decrease in accuracy between 0.5 and 3%.

Zagoruyko and Komodakis [50] also develop a method for removing skip connections in a NN. They replace skip connections with Dirac parameterization, creating a new NN called DiracNet. The Dirac parameterization is shown in Eq. 1,

$$\text{DiracNet [50]: } y = \sigma(x + Wx) \tag{1}$$

$$\text{ResNet [17]: } y = x + \sigma(Wx), \tag{2}$$

where $\sigma(\cdot)$ is the nonlinear activation function, $W$ is the layer weight matrix, $x$ is the layer input, and $y$ is the layer output. For ease of comparison with ResNets, Eq. 2 is simplified to show only one convolutional layer. In fact, skip connections in ResNets hop over more than one convolutional layer, while in DiracNets, the identity mapping is over one single convolutional layer. Therefore, the weights and the identity mapping of the input can be folded because $x + Wx = (I + W)x$. This change requires DiracNets to widen the NN layers in the ResNets that they started with. The authors showed that their technique could be used to create models with up to 34 layers. Although it works for shallower models, DiracNets show a decrease in accuracy between 0.5% and 1.5% compared to ResNets. In contrast, SkipRemover eliminates skip connections without widening the layers in the NN and does not need to make this accuracy tradeoff.

Li et al. [25] develop residual distillation (RD), which is a modified knowledge distillation framework. RD starts with a ResNet as the teacher and a plain CNN without skip connections as the student. Unlike standard knowledge distillation, RD passes the teacher's gradients to the student during training. This differs from TAILOR because RD starts with a student model without skip connections, whereas TAILOR *gradually* modifies a model's skip connections every few epochs during training without sharing gradients. Moreover, while RD removes all skip connections from models evaluated on simpler datasets like CIFAR-10 and CIFAR-100 [24], it fails to remove all skip connections in its ImageNet evaluation, leaving 18% of them in the network, which is a costly choice. In our ImageNet evaluation (see Sec. 4.1), our SKIPREMOVER method removes all skip connections with minimal accuracy loss.

Ding et al. [8] introduce a new model architecture RepVGG, which trains using 3×3 convolutional layers that are each skipped over by both a 1×1 convolution and an identity connection. At inference time, these connections can be re-parameterized into the parameters of the 3×3 convolutional layers. While RepVGG is more accurate than our SKIPREMOVER model, it requires starting from their specialized training model architecture. This is costly to developers who have already trained a model with skip connections on their dataset. Similarly, transferring a pre-trained RepVGG model to a new dataset via transfer learning can be time-consuming given the many different methods [36, 47, 52] to evaluate. As such, TAILOR is ideal for these developers because it modifies the skip connections of an *existing* pre-trained model to be more resource-efficient with minimal to no accuracy loss. Developers can leverage the training they have already done and need not start from scratch with a brand new RepVGG architecture.

## 2.2 Simplifying Skip Connection Hardware

ShuffleNet [28], DiracDeltaNet [48], and DenseNet [20] simplify skip connections by making them *concatenative*, i.e., they concatenate, rather than add, the skip connection data to the output of a layer. Concatenative skip connections take advantage of the fact that spatially consecutive memory accesses are typically faster than random accesses. This concatenation and off-chip data movement is possible using a simple controller (e.g., DMA engine).

TAILOR uses two techniques to simplify the skip connection hardware. SKIPREMOVER eliminates all logic and memory needed for a skip connection, making them less expensive than concatenative skip connections. Careful retraining allows skip connection removal in smaller networks with no degradation in accuracy. For larger networks, SKIPSHORTENER shortens the additive skip connections. By reducing their lifespans, the hardware implementation requires fewer resources. SKIPSHORTENER is not necessarily simpler than ShuffleNet [28] or DiracDeltaNet [48]. However, these concatenative skip connections have only been evaluated on image classification and object detection tasks. In our work, we demonstrate our SKIPREMOVER and SKIPSHORTENER methods on multifarious NNs and classification tasks, namely image-classifying ResNets of varying depths, DNA-basecalling QuartzNet-5×5, and automatic-speech-recognizing QuartzNet-10×5. With respect to DenseNet [20], SKIPSHORTENER ResNets use much less memory and bandwidth because DenseNet relies on significantly more skip connections throughout its NN. Given a NN with $L$ layers, DenseNet needs the memory and bandwidth to execute $L(L+1)/2$ concatenative skip connections, compared with SKIPSHORTENER ResNets' mere $L$ skip connections. With so many more skip connections, DenseNet is more expensive for hardware than SKIPSHORTENER ResNets.

Finally, all these techniques simplify skip connection hardware from the outset, building their models with modified skip connections and then training them from scratch. TAILOR differs because its hardware-aware training method *dynamically* alters the skip connections every few epochs during training, taking advantage of what the NN has learned with skip connections. Thus TAILOR allows the NN to gradually adapt to shortened skip connections (SKIPSHORTENER) or none at all (SKIPREMOVER).

(a) Traditional skip connection architecture



(b) No skip connection architecture
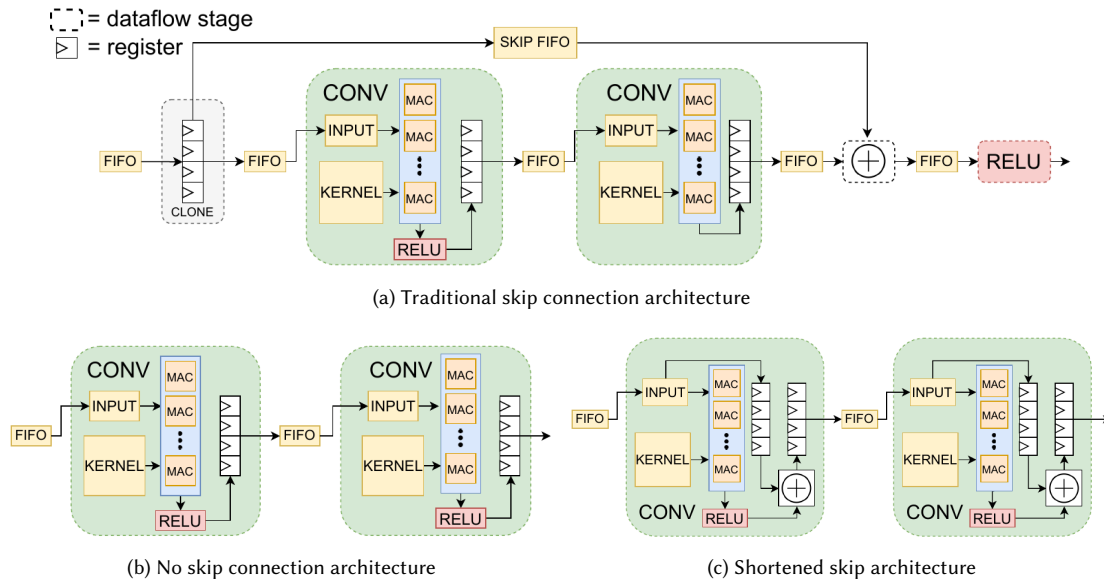
(c) Shortened skip architecture

Fig. 2. The hls4ml hardware architectures for traditional, shortened, and no skip connections. hls4ml pipelines each layer as is common for latency-critical tasks in resource-constrained environments [2, 12]. The three architectures correspond to a ResNet implemented with a traditional skip connection (a), shortened skip connections (b), and no skip connections (c). Note that we combine the batch normalization parameters with the kernel, as is commonly done [21].

## 3 TAILOR

Skip connections are important for training (to provide good accuracy), yet complicate implementation (requiring additional hardware resources and reducing performance). Tailor modifies skip connections to make their hardware implementation more efficient. Tailor uses a retraining method that gradually alters the network, resulting in little to no loss in accuracy.

### 3.1 Hardware Design

Fig. 2 shows three hardware implementations for NNs with traditional, shortened, and no-skip connections. The implementations correspond to accelerators produced by hls4ml—a tool that translates Python models into high-level synthesis code [11]. hls4ml creates a separate datapath for each layer and performs task-level pipelining across the layers. The layers communicate using FIFOs (AXI streams). Everything encapsulated by a dashed line resides in one pipeline stage. The inputs are fed into the architecture using a stream, and the results are given as an output stream. The weights are all stored on-chip, and all the internal results are stored on-chip. We evaluate each of these designs on FPGA later in Sec. 4.2 along with another style of architecture using a 2D array of processing elements. Tailor allows us to trade off between accuracy, performance, and resource usage through co-design of the neural network using hardware-aware training.

Fig. 2a shows the hardware needed to implement a single ResNet's skip connection. Note that in all of the designs shown in Fig. 2, we fuse the batch normalization parameters with the kernel, as is commonly done [21]. To be low latency and high throughput, the design uses task-level pipelining (i.e., the HLS dataflow pragma) for each NN layer, or

a small grouping of layers, and streams the data between each dataflow stage using first-in first-out buffers (FIFOs). Since FIFOs can only be read from once, skip connections complicate the design. We must spend a dataflow stage on cloning the skip connection data from its input FIFO into two other FIFOs so that it can be read twice for its two datapaths. The first path goes through a collection of convolutional and ReLU layers, and the second stores the data in a FIFO exclusive to skip connections (skip FIFO). Once the data has gone through the first path, we read from the skip FIFO to perform the addition to complete the skip connection's identity function. As such, implementing a skip connection on chip requires several extra FIFOs for handling the skip connection data, and this in turn increases on-chip memory resource utilization.

Ideally, we would eliminate the skip connections. As seen in Fig. 2b, without skip connections, we cut the number of dataflow stages in half (no more Clone, Add, or ReLU stages) and use less than half of the requisite FIFOs compared with Fig. 2a. All we need to do is pass the data through the convolutional and ReLU layers. This reduces resource utilization by up to 16% (see Sec. 4.2).

It may not be possible to remove the skip connections because they are essential for training convergence. In these cases, shortening the skip connections can simplify their hardware implementation. Fig. 2c shows a modified network with shortened skip connections such that *each skip connection's lifespan resides within a single dataflow stage*. We do not need additional dataflow stages to clone skip connection data. The shorter lifespans allow the shortened skip connections to be stored in *shift registers*, which can be implemented using the more abundant FFs as opposed to BRAMs, which is used in the traditional skip connection's hardware design. In this way, we exploit the short skip connections' lifetimes and use simpler, more efficient hardware memories to implement them (see Sec. 4.2). As such, we achieve a similar architecture to the version without skip connections (Fig. 2b), and similarly reduce resources spent on additional dataflow stages and FIFOs in Fig. 2a. SKIPSHORTENER is thus more resource-efficient than the traditional skip connection design. In fact, SKIPSHORTENER provides a tradeoff between the SKIPREMOVER and traditional designs because it uses more resources than SKIPREMOVER but less than the traditional one (see Sec. 4.2). But as we later show in Sec. 4.1, SKIPSHORTENER maintains accuracy in cases where SKIPREMOVER accuracy drops off. Thus, SKIPSHORTENER allows for design space exploration to balance accuracy and resource usage.

When used with hls4ml, TAILOR reduces resource consumption without changing the performance. This is a consequence of hls4ml's dataflow design; the resources we remove are not on the critical path—they are operating in parallel to the critical path. A dataflow design uses task-level pipelining, so reducing the resources spent on stages not on the critical path does not help or hurt overall throughput. Based on our Vivado co-simulation results, the clone stage executes in microseconds while the convolutional layer executes in milliseconds, an order of magnitude difference. Therefore, removing the clone buffer (Fig. 2b) or implementing it more efficiently (Fig. 2c) will not affect the overall dataflow latency because its latency is an order of magnitude less than the convolution's latency. This means TAILOR's resource reductions do not increase or decrease latency or throughput for this architecture style, as later shown in Tab. 7.

Another prevalent style of FPGA CNN architectures instantiates a 2D processing element (PE) array and iteratively programs the convolutions and other operations onto that PE array. We call this style of computation a Reconfigurable DNN Architecture. Fig. 3 provides an example architecture used in our experiments. We build this architecture using DeepSoCFlow [1]. Following the taxonomy described in [22], the reconfigurable DNN architecture is a 2D array of processing elements that optimally perform standard convolution and matrix multiplication with high data reuse.

---

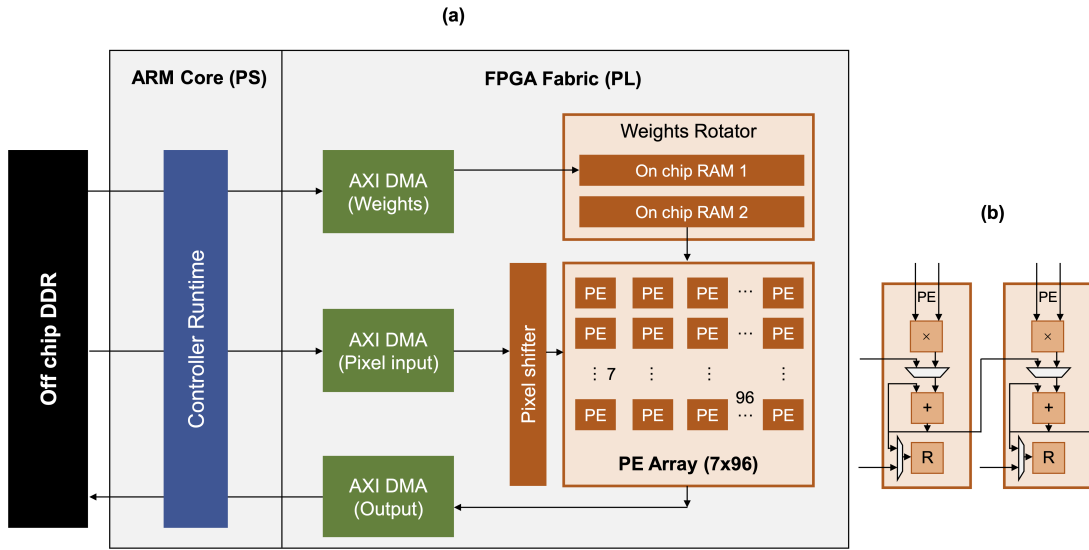[1]https://github.com/abarajithan11/deepsocflow

**(a)**



Fig. 3. (a) A Reconfigurable DNN Architecture synthesized on a ZCU102 FPGA development board. The architecture has a 2D array of processing elements that are iteratively programmed to compute layer operations. The controller runtime programs the DMA engines to load off-chip inputs and weights and store the intermediate and final results off-chip. (b) The processing elements (PE) are a multiply-accumulate datapath.

The dataflow is primarily output stationary while prioritizing maximal weight reuse and also reusing inputs to an extent. The engine performs fixed-point computations, where the input, weight, and output bit widths are adjustable as synthesis parameters, along with the number of rows and columns of processing elements. The weights rotator prefetches the weights of the next iteration into one block of on-chip memory while the other bank delivers weights, rotating them hundreds of times for maximal data reuse. The Pixel Shifter shifts perform vertical convolution. Partial sums are shifted to the PE on the right to compute horizontal convolution. The results are streamed out through the output DMA to the off-chip memory. The runtime controller would perform the residual addition, quantization, and activation on the processing system side while the engine computes the next iteration. Our implementation uses the ARM processor available in the Zynq chip. FPGAs without processors could instantiate a softcore processor to perform the controller runtime operations.

The TAILOR optimizations have different effects on the Reconfigurable DNN architecture as compared to hls4ml architecture. The Reconfigurable DNN architecture computes skip connections by loading input data from off-chip memory and performing the required operations upon it (addition, convoluation) Thus, unlike in hls4ml, removing a skip connection does not change the architecture; instead it changes how computations are mapped to that architecture. Skip connection removal eliminates the need to fetch the skip connection data and perform the associated convolution and addition operations. This increases the overall performance as we describe in Sec. 4.
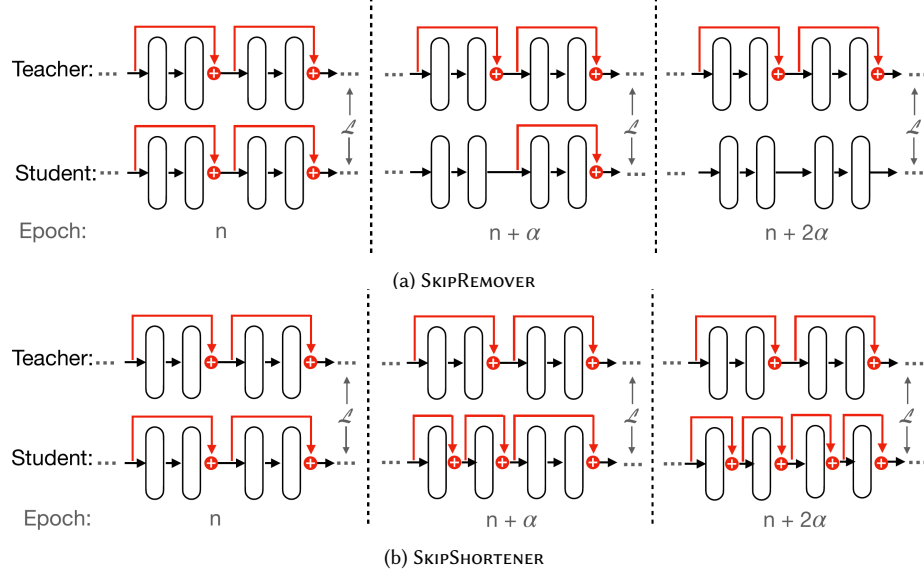
(a) SKIPREMOVER



(b) SKIPSHORTENER

Fig. 4. Three iterations in the SKIPREMOVER and SKIPSHORTENER algorithms as applied to a ResNet. In this example, skip connections are altered every $\alpha$ epochs and $\alpha | n$. Each pill block represents a set of convolutional, BN, and ReLU layers, and the skip connections are in red. $\mathcal{L}$ is the KD loss function defined in Eq. 3. Only the student model is used for inference.

## 3.2 Hardware-aware Training

It is difficult to modify a NN's skip connections without reducing accuracy. Naively removing all skip connections before or after training a NN is detrimental to its accuracy. Instead, TAILOR consists of two training algorithms, *SKIPREMOVER* and *SKIPSHORTENER*, that gradually alter a NN's skip connections on the fly—removing or shortening them every few epochs—in order to make them resource-efficient. Gradually altering the model during training tempers the performance drop of removing or shortening the skip connections, yielding minimal to no loss in accuracy as well as significant advantages in the hardware implementation, as described above.

TAILOR's iterative learning approach finetunes the altered NNs using a compression method known as *knowledge distillation (KD)* [19]. KD distills the knowledge of a larger, more complex NN (the teacher) into a smaller, simpler NN (the student). While the student model is training, it compares its output to the teacher model's output and thus learns from the teacher to perform better. KD provides impressive results for compressing NNs for various applications [31, 41, 44]. In traditional KD, the teacher model is already trained, and the student model is trained to match the teacher's behavior by replicating its output. The student achieves this by training with a loss function

$$\mathcal{L} = (1 - \beta)\mathcal{G}(\ell, s) + \beta\mathcal{H}(t, s) \tag{3}$$

where $\mathcal{G}$ and $\mathcal{H}$ are distance functions, $s$ and $t$ are student and teacher output vectors respectively, $\ell$ is the correct label vector, and $\beta$ is a tunable parameter [19].

With this idea in mind, both SKIPREMOVER and SKIPSHORTENER start with two identical pre-trained NNs with traditional skip connections, where one serves as the teacher and the other serves as the student. During the retraining stage, SKIPREMOVER removes a given skip connection every few epochs. SKIPSHORTENER takes a similar iterative

approach and, every few epochs, splits a given skip connection into multiple shorter ones. The skip connections are removed or shortened starting from the first skip connection encountered in the NN (from the input) to the last.

Fig. 4 visualizes both SKIPREMOVER's (Fig. 4a) and SKIPSHORTENER's (Fig. 4b) training algorithms for a ResNet-style NN. During training, we remove (SKIPREMOVER) or shorten (SKIPSHORTENER) one of the student's skip connections every $\alpha$ epochs. If $n$ is divisible by $\alpha$ (as in Fig. 4), then at epoch $n$, the student has had $n/\alpha$ skip connections altered, and we are viewing the next two skip connections to be modified in the student model: the $(n/\alpha) + 1$st and $(n/\alpha) + 2$nd. At epoch $n + \alpha$, the $(n/\alpha) + 1$st skip connection is altered (removed under SKIPREMOVER or split into two shorter skip connections under SKIPSHORTENER). The NN then trains for $\alpha$ epochs so that the student model can improve its weights given the latest model topology. Afterwards, at epoch $n + 2\alpha$, the $(n/\alpha) + 2$nd skip connection is similarly altered. During the entire skip modification retraining process, the student uses the KD loss function $\mathcal{L}$ defined in Eq. 3 to learn from the teacher and the true labels. The teacher's model topology and weights remain fixed during training. Once all skip connections have been altered, the student model continues training under KD for the remaining number of training epochs as defined by the user. Only the student model is used for inference.

TAILOR is novel because it *dynamically* transforms skip connections every few epochs during training. This is an instance of *hardware-aware training* because the skip connection are slowly altered specifically to reduce hardware resources, as previously discussed in Sec. 3.1. The gradual skip connection alterations allow the NN to take advantage of what it has learned with skip connections, so that it can dynamically adapt to shortened skip connections (SKIPSHORTENER) or none at all (SKIPREMOVER). Alg. 1 describes TAILOR's hardware-aware training process.

## 4 RESULTS

We evaluate TAILOR on two popular kinds of NNs that rely on skip connections: ResNets [17] and QuartzNets [23]. We study the effects of TAILOR on model accuracy, quantization, and hardware resource utilization.

### 4.1 Training results

To evaluate how TAILOR affects a NN's accuracy, we train ResNets and QuartzNets of varying depths using our SKIPREMOVER and SKIPSHORTENER algorithms in PyTorch [39]. The ResNets range from 20 to 110 layers and are trained on the CIFAR-10 [24], CIFAR-100 [24], and SVHN [35] datasets. We also evaluate ResNet50, which has a different skip connection topology than standard ResNets, on the ImageNet dataset [7]. The QuartzNets span between 29 and 54 layers. Their structure is determined by the number and lifetimes of their skip connections. For instance, a QuartzNet-10×5 has 10 skip connections that each have a lifetime of 5 sets of layers. We train a QuartzNet-5×5 on the Oxford Nanopore Reads dataset [42], a DNA basecalling task. We also train a QuartzNet-10×5 on the LibriSpeech dataset [37], an automatic speech recognition (ASR) task, which converts speech audio to text. ASR tasks are assessed using word error rate (WER), which measures the percent of words that the model predicted incorrectly. In all of our ResNet and QuartzNet-10×5 training experiments, we set $\alpha = 3$ in Alg. 1, so skip connections are removed or shortened every three epochs. For QuartzNet-5×5, we set $\alpha = 1$ instead because it trains better this way. For the ResNets, we set $\mathcal{G}$ and $\mathcal{H}$ in Eq. 3 to *categorical cross entropy* and *mean-squared error*, respectively, and set $\beta = 0.35$. For the QuartzNets, we set Eq. 3's parameters similarly, except for $\mathcal{G}$, which we set to *connectionist temporal classification loss*, which is used to train difficult tasks involving sequence alignment (like DNA basecalling and ASR). Note that in our training results, "Baseline" refers to the unmodified NN counterpart with conventional skip connections.

Fig. 5 shows that SKIPREMOVER works well for ResNet-44 and smaller, at times even outperforming its baseline (traditional skip connection model). However, its accuracy drops as the number of layers increases. This indicates that

---

**Algorithm 1:** HARDWARE-AWARE TRAINING

---

1  set *alter* // `REMOVE or SHORTEN`
2  let $\alpha$ = how often to modify a skip connection
3  *teacher* = pre-trained model
4  *student* = pre-trained model
5  let current-skip = *student*'s first skip connection from the input side
6  let current-layers = all layers skipped by current-skip
7  **Function** `SkipRemover(`*current-skip*`)`:
       // see Fig. 4a
8      remove current-skip
9      **return** student model's next skip connection from the input side
10 **Function** `SkipShortener(`*current-skip, current-layers*`)`:
       // see Fig. 4b
11     Split current-skip into *len*(current-layers) skip connections
12     current-skip = student model's next skip connection from the input side
13     current-layers = student's next layers skipped by the new current-skip
14     **return** current-skip, current-layers
15 **for** *i in epochs* **do**
16     **if** *i* ≠ 0 *and i* mod $\alpha$ = 0 **then**
17         **if** *alter* == *REMOVE* **then**
18             current-skip = SkipRemover(*current-skip*)
19         **else if** *alter* == *SHORTEN* **then**
20             current-skip, current-layers = SkipShortener(*current-skip, current-layers*)
21     **end**
22     train *student* using Eq. 3
23 **end**
24 save the student model

---



(a) CIFAR-10                          (b) CIFAR-100                          (c) SVHN
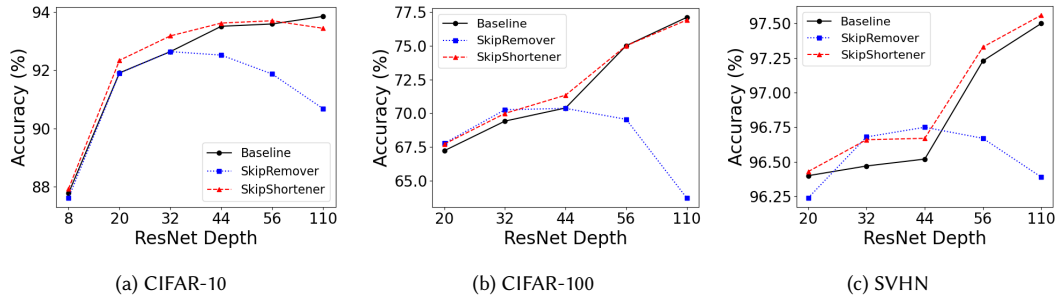
Fig. 5. Top-1 accuracy of SKIPREMOVER and SKIPSHORTENER ResNets of increasing depth on various datasets. "Baseline" refers to an unmodified ResNet with conventional skip connections.

shallower NNs do not need skip connections for these classification tasks, but they become more necessary for deeper networks. SKIPSHORTENER mostly outperforms the baseline on all three datasets, even on deep models.

*4.1.1 Ablation Studies.* We also perform *ablation studies* in which we remove key parts of TAILOR to understand why they are critical to minimizing accuracy loss. One key part of SKIPREMOVER/SKIPSHORTENER is the dynamic
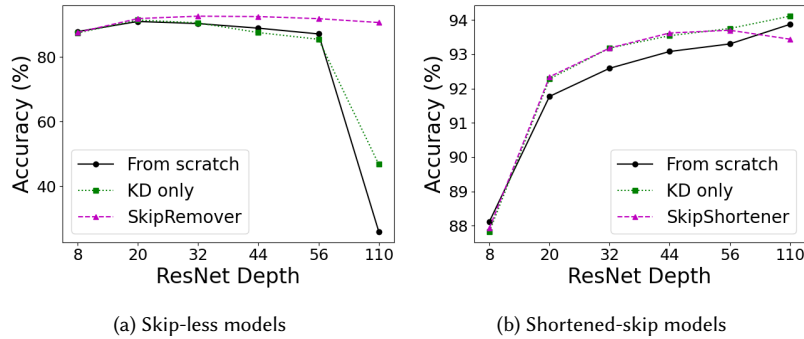
(a) Skip-less models　　　　　(b) Shortened-skip models

Fig. 6. Accuracy results for ResNets whose skip connections are all altered before training (apart from SKIPREMOVER and SKIPSHORTENER) on CIFAR-10. "From scratch" means training with randomly initialized weights without KD. "KD only" means training without dynamic skip alterations.

Table 1. Top-1 accuracy of ResNet-50 on the ImageNet dataset. *RD [25] only removes 82% of the skip connections.

| Model | Accuracy (%) |
|---|---|
| ResNet-50 | 75.85 |
| No skips (from scratch) | 58.36 |
| No skips (KD only) | 69.40 |
| Residual distillation (RD)* [25] | 76.08 |
| RepVGG-A2 [8] | 76.48 |
| **SKIPREMOVER** | **75.36** |

skip connection removal/shortening that occurs every few epochs during training under KD. We thus take away this dynamic model alteration by first altering the NNs to have either no skip connections or shortened skip connections. These pre-modified NNs are then trained under KD only. Another key part of SKIPREMOVER and SKIPSHORTENER is KD. We evaluate how skip-less and shortened-skip NNs perform without KD, training from randomly initialized weights (i.e., from scratch).

For ResNets trained on CIFAR-10, SKIPREMOVER and SKIPSHORTENER usually yield better results than either normal training or using KD-only on a statically pre-modified network on CIFAR-10 per Fig. 6a and Fig. 6b. The difference between all of the approaches in the figures is minimal for smaller models, but it becomes more apparent as NN depth increases. For instance, skip-less ResNet-110 under regular training yields an accuracy of 26.02% versus SKIPREMOVER, which achieves an accuracy of 90.68%, a 64.66% difference. SKIPREMOVER marginally outperforms regular training and KD-only on smaller skip-less models, but performs much better in comparison as the networks deepen. SKIPSHORTENER also generally performs better than the other two approaches for shortened skip models. Regular training mostly lags behind both KD and SKIPSHORTENER for shortened skip models.

For ResNet-50 on ImageNet, we only apply SKIPREMOVER because it uses an irregular skip connection architecture known as a "bottleneck block" to reduce the number of parameters [17]. This block has a skip connection spanning three layers: a 1×1 convolution, then a 3×3 convolution, then another 1×1 convolution (Fig. 7a). This irregular topology is not optimal for SKIPSHORTENER because it requires the majority of the shortened skip connections to pass through extra downsampling 1×1 convolutions to match the activation tensor shapes, significantly increasing the number of model parameters. As such, for ResNets with bottleneck blocks, like ResNet-50, we recommend SKIPREMOVER. As seen
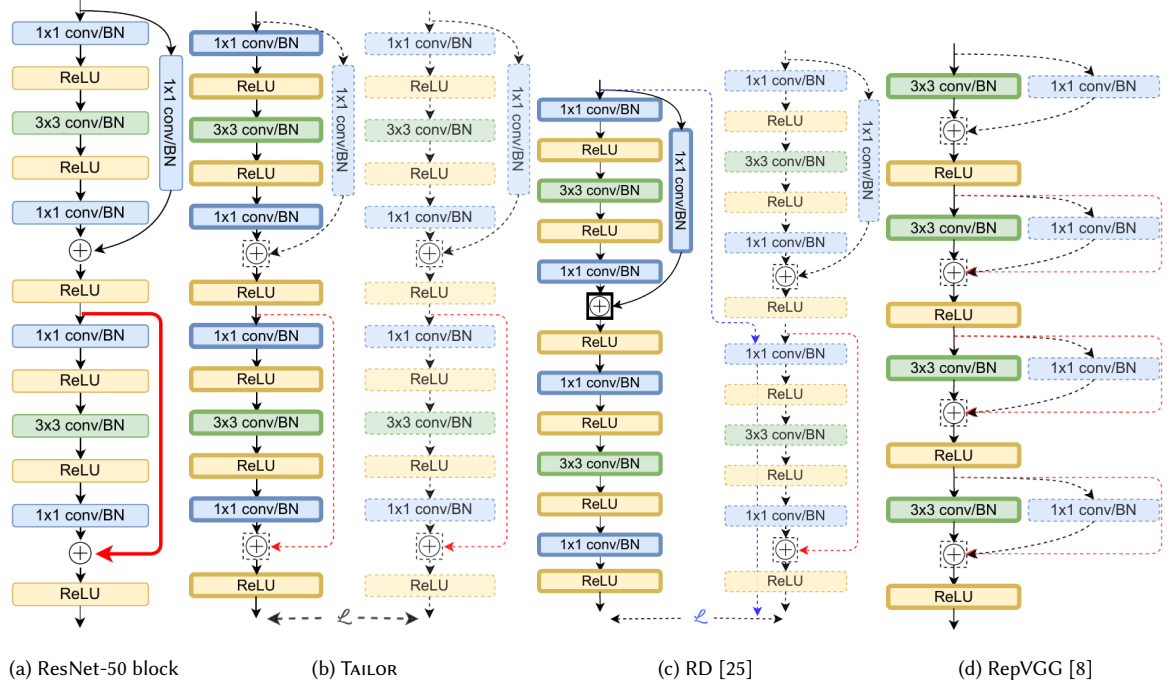
(a) ResNet-50 block     (b) Tailor     (c) RD [25]     (d) RepVGG [8]

Fig. 7. Comparing Tailor's ResNet-50 skip removal method with residual distillation (RD) [25] and RepVGG [8]. The dashed portions are only used during training and are later removed, leaving the final inference NNs, indicated by bolder lines. Note that Tailor (b) removes skip connections from a pretrained ResNet-50 (a). RD does the same but uses a modified KD method that does not remove the 1×1 convolution addition (c). RepVGG starts training from a different NN topology altogether (d).

in Tab. 1, SkipRemover incurs a 0.49% accuracy loss compared to the traditional ResNet-50. Compared to prior work such as RD [25] and RepVGG [8], SkipRemover has slightly lower accuracy (at most 1.12% accuracy difference)[2].

Nevertheless, SkipRemover has two advantages compared with these methods. First, SkipRemover removes all skip connections from ResNet-50, whereas RD only removes 82% of them. RD does not remove the 1×1 convolution addition used for downsampling (see Fig. 7c), which is particularly detrimental. In our experiments on hls4ml architectures, Vivado HLS estimates that ResNet-50's large 1×1 convolution skip connection consumes as many resources as the layers it skips over, effectively doubling resource consumption for that skip connection block. Although Vivado HLS has a tendency to overestimate the actual place-and-route (P&R) resource utilization, these estimates demonstrate that performing the 1×1 convolution is a nontrivial task that significantly affects resource consumption. Second, SkipRemover removes the skip connections from an *existing* pre-trained model, whereas RepVGG requires developers to adopt a new model topology (see Fig. 7d). If developers do not already have a model on hand, RepVGG is a better option. However, if developers already have a ResNet trained for their specific dataset, it is advantageous to use SkipRemover if they can afford a small accuracy loss. This prevents starting from scratch with RepVGG, which could require extensive hyperparameter tuning. Even finetuning a pre-trained RepVGG model to a new dataset using transfer learning is time consuming, as it is unclear which of the many methods [36, 47, 52] would work best. Instead, SkipRemover allows developers to take advantage of their existing work and achieve a more resource-efficient model.

---

[2]Ding et al [8] introduce RepVGG models of varying depths. We compare against RepVGG-A2 because it is about the same size as ResNet-50.

Table 2.   Top-1 accuracy of QuartzNet-5×5 on the Oxford Nanopore Reads dataset.

| Model | Accuracy (%) |
|---|---|
| QuartzNet-5×5 | 95.107 |
| No skips (from scratch) | 94.475 |
| No skips (KD only) | 94.863 |
| SkipRemover | **95.086** |
| Shortened skips (from scratch) | 95.019 |
| Shortened skips (KD only) | 95.016 |
| SkipShortener | 94.902 |

Table 3.   Word error rate (WER) of QuartzNet-10×5 on LibriSpeech dataset. This includes clear ("dev-clean") and noisy ("dev-other") audio samples. "—" indicates the model failed to converge.

| Model | dev-clean WER (%) | dev-other WER (%) |
|---|---|---|
| QuartzNet-10×5 | 5.56 | 16.63 |
| No skips (from scratch) | — | — |
| No skips (KD only) | — | — |
| SkipRemover | — | — |
| Shortened skips (from scratch) | **6.40** | **17.68** |
| Shortened skips (KD only) | 7.14 | 19.95 |
| SkipShortener | 7.86 | 21.16 |

For QuartzNet-5×5, the SkipRemover model performs the best—only 0.021% from the baseline (Tab. 2). These results all have high accuracy likely because DNA basecalling is an easier sequence alignment task (only four classes) and the model is more than sufficient. For a harder ASR task like LibriSpeech, QuartzNet-10×5 fails to converge without skip connections. Since the model must translate audio samples to text, the audio samples can be noisy, making ASR harder. LibriSpeech, in fact, divides its test samples into "dev-clean" for clearly spoken samples and "dev-other" for noisy samples. With such a challenging task, it is not possible to remove the skip connections (like with DNA basecalling). Nonetheless, QuartzNet-10×5 performs well under SkipShortener, as it is within 2% of the baseline WER (Tab. 3). For both QuartzNet-5×5 and -10×5, the best performing shortened skip connection model was one whose skip connections were shortened first and then trained from scratch. While SkipShortener has minimal accuracy loss for both QuartzNets, we recommend training a model with shortened skip connections from scratch for this task.

Overall, SkipRemover and SkipShortener perform better than either training on randomly initialized weights or training with KD only. For harder tasks like ASR though, training a shortened-skip model from scratch is a better choice. Nevertheless, the success of SkipRemover and SkipShortener lies in augmenting KD with dynamic skip alterations.

## 4.2  Hardware Results

We first quantize ResNets ranging from 20 to 56 layers deep to see how Tailor's accuracy fares under reduced precision. We then evaluate Tailor's effects on hardware resources and latency by performing a case study on ResNet-20-style skip connections implemented using the hls4ml architecture, i.e., the designs illustrated in Fig. 2. We select this style of skip connection because it is the fundamental building block of ResNets that range from 20 to 110 layers. In our case study, we vary the bit precision and number of filters to see how Tailor scales up. Based on how Tailor's resource reductions scale, designers can understand how Tailor extrapolates to their own hardware designs. We report latency as well

as P&R resource results on the Alveo U200 FPGA accelerator card (part no. `xcu200-fsgd2104-2-e`). For end-to-end application results, we evaluate the benefits of Tailor on two different styles of CNN architectures. The first uses the hls4ml tool to generate architectures. The second is the Reconfigurable DNN Engine—a 2D array of processing elements. Both styles of architectures are described in Sec. 3.1.

*4.2.1  Quantization.* The parameters of a hardware-accelerated NN are typically quantized from floating-point to fixed-point precision [6, 33, 48].
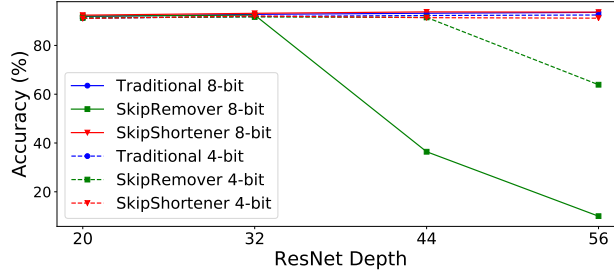


Fig. 8.  Quantized accuracy results for 8-bit and 4-bit fixed point using Brevitas.

Quantizing deep NNs with minimal accuracy loss is a largely manual and time-consuming task [14]. We use Brevitas [38] to quantize our SkipRemover and SkipShortener ResNets with depths of 20 to 56 from 32-bit floating-point (float32) to 8-bit and 4-bit fixed-point precision on the CIFAR-10 dataset. We modified Tailor's hardware-aware training algorithm where the teacher continues to use floating-point representation whereas the student is quantized. This results in the student undergoing quantization-aware training. In Fig. 8, we find that SkipShortener ResNets consistently outperform traditional ResNets under Brevitas quantization-aware training by 0.5%. SkipRemover ResNets start to suffer from the lack of bits as they get deeper, with accuracy dropping to random classification for ResNet-56. But, Brevitas is only one of dozens of ways to quantize neural networks [9, 10, 14, 33, 46], so it may be the case that a SkipRemover ResNet-56 requires a different method of quantization to achieve a quantized accuracy similar to its float32 counterpart.

*4.2.2  FPGA Evaluation.* Our first study looks solely at one ResNet block. The second study performs an end-to-end implementation of ResNet8 and ResNet50.

For our case study on a ResNet skip connection blocks (see designs in Fig. 2), we evaluate Tailor at `ap_fixed<8,3>` and `ap_fixed<16,6>` precisions using the hls4ml architecture. Under both bitwidths, we increase the number of filters for all designs from 16 to 32 to 64. This way, we can understand how Tailor scales with the number of filters. We use hls4ml [12] to translate these hardware designs into Vivado HLS, targeting the Alveo U200 FPGA accelerator card. hls4ml uses task-level pipelining (i.e., HLS dataflow) for each NN layer, or small group of layers and streams data between dataflow stages using FIFOs. hls4ml also exposes a knob known as *reuse factor*, which determines how often multipliers are reused in a design. To fairly compare our designs as the number of filters increases, we fix the reuse factor to 576. We then synthesize our designs to report P&R resource utilization as well as co-simulation latency results. Lastly, we run the designs on the U200 to verify correctness.

Under 8-bit precision, we find that both SkipRemover and SkipShortener reduce resources. Tab. 4 summarizes our P&R results. Since our model uses 8-bit precision, we see that all of our models exhibit low DSP usage and higher LUT

Table 4. Place-and-route resource utilization of a skip connection block as the number of filters increases for $\langle 8, 3 \rangle$ precision on an Alveo U200. SkipRemover reduces LUT and FF usage, whereas SkipShortener trades an increase in FFs for a decrease in LUTs. T = Traditional, R = SkipRemover, S = SkipShortener.

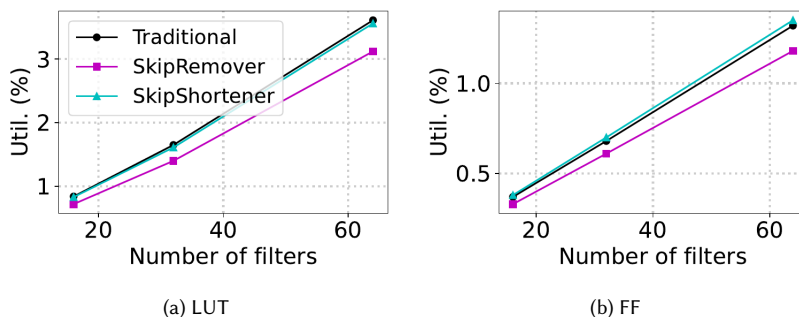| # filters | LUT | | | FF | | | DSP | BRAM |
|---|---|---|---|---|---|---|---|---|
| | T | R | S | T | R | S | T/R/S | T/R/S |
| 16 | 9,984 | 8,482 | 9,764 | 8,654 | 7,841 | 8,916 | 0 | 18.5 |
| 32 | 19,566 | 16,512 | 18,993 | 16,183 | 14,506 | 16,489 | 0 | 36.5 |
| 64 | 42,688 | 36,882 | 42,121 | 31,124 | 27,815 | 31,850 | 0 | 82 |



(a) LUT　　　　　　　　　　　　　　　　　　(b) FF

Fig. 9. Percent resource utilization of a $\langle 8, 3 \rangle$ skip connection block at various filter sizes on an Alveo U200. DSPs and BRAMs remain the same across the three designs, so they are not shown. SkipRemover and SkipShortener LUT and FF reductions scale linearly, as expected.
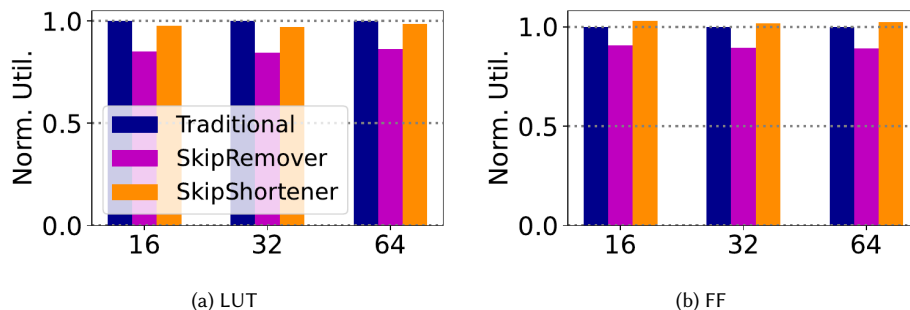


(a) LUT　　　　　　　　　　　　　　　　　　(b) FF

Fig. 10. Resource utilization normalized to the traditional design of a $\langle 8, 3 \rangle$ skip connection block at various filter sizes. DSPs and BRAMs remain the same across the three designs, so they are not shown. SkipRemover and SkipShortener LUT and FF reductions scale proportionally, as expected.

and FF utilization. This is because Vivado HLS maps multiplications on datatypes that are less than 10 bits to LUTs instead of DSPs, as noted by [2, 48]. It is possible to pack two 8-bit weights into a DSP [13], but this is out of scope and orthogonal to the effects Tailor has on hardware. Furthermore, all of the traditional and Tailor designs use the same amount of BRAMs with respect to the number of filters because here the BRAMs are used solely for on-chip weight storage, which does not differ across design. Nonetheless, SkipRemover decreases LUT usage by up to 16% and FF usage by up to 11% compared with the traditional design (Fig. 10). These resource savings represent the extra hardware needed to implement a skip connection and subsequently the resources saved. As previously mentioned in Sec. 3.1, the extra dataflow stages that carry out a skip connection are no longer necessary. More importantly, SkipRemover's

Table 5. Place-and-route resource utilization of a skip connection block as the number of filters increases for $\langle 16, 6 \rangle$ precision on an Alveo U200. SKIPREMOVER reduces resources across the board, whereas SKIPSHORTENER trades an increase in LUTs for a decrease in FFs and BRAMs. T = Traditional, R = SKIPREMOVER, S = SKIPSHORTENER.

| # filters | LUT | | | FF | | | DSP | BRAM | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | T | R | S | T | R | S | T/R/S | T | R | S |
| 16 | 14,733 | 13,320 | 14,933 | 17,044 | 14,935 | 16,438 | 12 | 60.5 | 52.5 | 42.5 |
| 32 | 28,498 | 25,330 | 28,184 | 32,923 | 28,747 | 31,764 | 48 | 124 | 108 | 84.5 |
| 64 | 55,699 | 50,074 | 55,720 | 64,564 | 56,263 | 62,252 | 192 | 267.5 | 235.5 | 203.5 |

savings scale linearly as the number of filters increases from 16 to 64 (Fig. 9). SKIPSHORTENER's resource reductions present a tradeoff, increasing FFs by 2% in exchange for decreasing LUTs by 3% (Fig. 10). SKIPSHORTENER lowers LUT utilization because the lifespan of each skip connection lasts only one dataflow stage instead of the traditional two. This means we need not spend extra logic on the dataflow stages needed to copy the skip connections to buffers that last longer than one stage. However, since the shortened skip connection now fully resides in a single dataflow stage (previously described in Fig. 2c), this requires some extra FFs. This represents the tradeoff SKIPSHORTENER provides at 8-bit precision: some extra FFs for fewer LUTs. These resource tradeoffs also scale linearly as the number of filters scales up, as seen in Fig. 9.

We find more dramatic resource reductions when we look at our 16-bit designs. Tab. 5 summarizes our P&R results. In contrast with our 8-bit designs, at higher precision, our designs rely more on DSPs and BRAMs. This time the BRAMs are used not only to store weights on chip but also to implement the FIFOs that connect the dataflow stages. Therefore, as we tailor the dataflow stages according to each design (e.g., SKIPREMOVER or SKIPSHORTENER), the BRAMs now also reflect these changes. At its best, SKIPREMOVER lowers LUTs by 11%, FFs by 13%, and BRAMs by 13%. Without a skip connection to implement, SKIPREMOVER uses fewer resources than the traditional design. The DSPs remains unchanged because they are used solely for the convolutional layers' multiplications and not the skip connection, which is also the case for SKIPSHORTENER.

Similar to the 8-bit designs, SKIPSHORTENER presents a resource tradeoff—this time trading a small increase in LUTs (at most 1%) for decreases in FFs and BRAMs. In the best case, SKIPSHORTENER reduces LUTs by 1%, FFs by 4%, and BRAMs by 34%. While SKIPSHORTENER uses fewer LUTs than the traditional case for 32 filters, SKIPSHORTENER pays about a 1% increase in LUTs for 16 and 64 filters in exchange for decreases in FFs and BRAMs. This small disparity is likely an artifact of the heuristics Vivado P&R uses to allocate resources. Again, these resource tradeoffs and savings are possible because the shortened skip connections can be implemented within a single dataflow stage due to its reduced lifetime. Tab. 6 shows that the lifetime of each shortened skip connection is a little less than half the lifetime of the traditional one. With shorter lifetimes, we find that the SKIPSHORTENER's skip connections' FIFOs can now be implemented using shift registers instead of BRAMs, which is what the traditional design still uses (Tab. 6). Shift registers are much more efficient memories compared to BRAMs. As such, it is advantageous to hardware designers to consider how SKIPSHORTENER provides opportunity to implement skip connections with a more efficient memory architecture like shift registers. This leads to 30–34% fewer BRAMs than the traditional design, even as the number of filters scales up. While in this case SKIPSHORTENER uses fewer BRAMs than SKIPREMOVER does, SKIPSHORTENER offsets this difference by using more FFs than SKIPREMOVER does. For both SKIPREMOVER and SKIPSHORTENER, resource utilization (and the associated reductions) scale linearly, as seen in Fig. 11.

TAILOR does not affect latency for hls4ml architectures. As seen in Tab. 7, for each number of filters, all designs exhibit the same latency, according to co-simulation on an Alveo U200. The slight decrease in latency as the number of filters
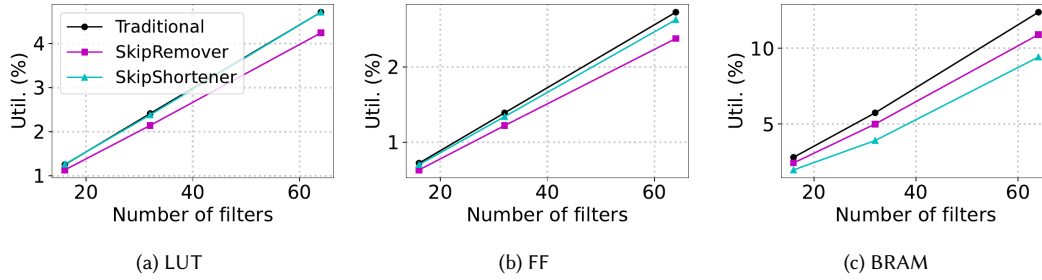
(a) LUT                          (b) FF                          (c) BRAM

Fig. 11. Percent resource utilization of a ⟨16, 6⟩ skip connection block at various filter sizes on an Alveo U200. SkipRemover and SkipShortener resource reductions scale linearly, as expected.
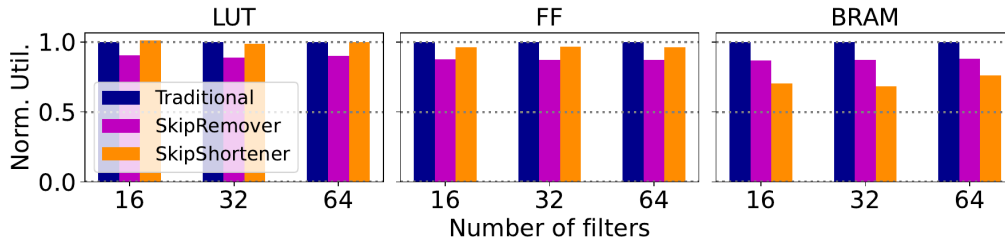


Fig. 12. Resource utilization normalized to the traditional design of a ⟨16, 6⟩ skip connection block at various filter sizes. The SkipRemover and SkipShortener resource savings scale proportionally as the number of filters scales up.

Table 6. FIFO depths of a single skip connection hardware design at 16-bit precision. SkipRemover has no skip connections, so it has no skip connection FIFOs.

| Hardware Design | FIFO Depth | FIFO Implementation |
|---|---|---|
| Traditional | 69 | BRAM |
| SkipRemover | 0 | — |
| SkipShortener 1st skip | 33 | Shift Register |
| SkipShortener 2nd skip | 34 | Shift Register |

Table 7. Latency co-simulation results of a skip connection block at ⟨8, 3⟩ and ⟨16, 6⟩ precision. The latency for the Traditional, SkipRemover, and SkipShortener designs are the same for each number of filters because they all rely on task-level pipelining that reuses multipliers at the same rate (576×).

| # filters | Latency (ms) Traditional/SkipRemover/SkipShortener |
|---|---|
| 16 | 23.38 |
| 32 | 23.05 |
| 64 | 22.39 |

scales is due to an increase in DSPs and a higher degree of parallelism. As discussed in Sec. 3.1, hls4ml designs pipeline their tasks. The convolutions' multiplication tasks dominate the overall dataflow latency. The tasks that SkipRemover eliminates and SkipShortener implements more efficiently, namely the skip connection cloning and addition stages, have significantly lower latency than the convolutions and are thus not on the critical path. The throughput thus remains the same.

By shortening skip connections, we reduce their lifespans, which provides an opportunity for simplifying their hardware implementation specifically for hls4ml architectures. However, shortening skip connections is not beneficial for all architectures. As seen in Tab. 8, shortening skip connections is worse for both GPU and CPU because doing so increases off-chip memory accesses. These extra accesses lower throughput by 5% on GPU and 2% on CPU. On FPGAs with hls4ml architectures, however, we can modify the architecture to take advantage of shortened skip connections, reducing resource consumption without negatively affecting throughput (Tab. 8).

Table 8. Normalized throughput of a ResNet20. The GPU and CPU both were run with batch size = 64, whereas FPGA was run with batch size = 1. Throughput is normalized column-wise to the top entry. GPU = 1080Ti. CPU = AMD Ryzen 9 5900X. FPGA = Alveo U200. SkipRemover increases GPU and CPU throughput because it decreases off-chip memory accesses. SkipShortener, however, decreases GPU and CPU throughput because it increases off-chip memory accesses. For a fully on-chip, dataflowed FPGA architecture, neither SkipRemover nor SkipShortener have any effect on throughput.

| Model | GPU | CPU | FPGA |
|---|---|---|---|
| Traditional skip connections | 1× | 1× | 1× |
| SkipRemover | 1.11× | 1.03× | 1× |
| SkipShortener | 0.95× | 0.98× | 1× |

We performed two studies to understand how Tailor performs for end-to-end implementations of ResNet models. The first is ResNet8 from MLPerf Tiny that was designed in hls4ml [3, 5]. The second is ResNet50 implemented on the Reconfigurable DNN architecture.

The ResNet8 model targets the Alveo U200. It uses 16-bit fixed-point representation with six integer bits. The reuse factor for the layers was hand-tuned to 72, which directly affects the resource usage and latency of the layers. The reuse factor is one of the more important knobs for design space exploration in hls4ml and is often hand-tuned to maximize resource usage of the platform while optimizing the overall network performance.

Table 9. MLPerf Tiny ResNet8 model implemented using hls4ml with skip connection, with shortened skip connections, and without skip connections.

| | With Skip Connections | Shortened Skip Connections | Without Skip Connections |
|---|---|---|---|
| Accuracy (%) | 87.39 | 87.93 | 87.62 |
| LUTs | 158609 | 165699 | 144206 |
| FFs | 196012 | 204914 | 181768 |
| DSP48s | 1083 | 1083 | 1043 |
| BRAMs | 173 | 158.5 | 156 |

Tab. 9 shows the resource usage results for the ResNet8 model with skip connections, with shortened skip connections, and without skip connections. Removing the skip connections has clear benefits across all the resources. Shortening the skip connections reduces BRAMs while increasing LUTs and FFs. Both the shortened skip connections and the removed skip connections models show improved accuracy over traditional skip connections. In all cases, the latency remains the same, requiring 304,697 cycles running at 100 MHz (approximately 3ms/inference).

Our second full model case study implemented a Reconfigurable DNN architecture on the ZCU102 development board which contains a Zynq UltraScale+ MPSoC. The Reconfigurable DNN array is configured to have 7 rows × 96 columns for a total of 672 of processing elements (PEs) that support 8-bit inputs and 8-bit weights. Each PE contains a multiplier and an accumulator implemented using DSPs on FPGA fabric. Input pixels and weights are streamed into the engine as AXI-Stream packets. Images are processed in batches of 7, to increase the reuse and reduce memory accesses.

The Reconfigurable DNN architecture was synthesized, placed & routed at a clock frequency of 250 MHz on a ZCU102. The architecture with $7 \times 96 = 672$ PEs used 49057 LUTs (18%), 81446 flip flops (15%), 114 BRAMs (13%), and 1344 DSPs (53%) on the FPGA fabric.

We implemented a ResNet50 model with and without skip connections on a 672-element Reconfigurable DNN architecture running on the ZCU102. Tab. 10 shows the performance of ResNet50. Removing the skip connections largely benefits the performance due to the removal of the $1 \times 1$ convolution blocks. Removing the skip connections also removes those layers, which no longer need to be scheduled on the PE array. The results are much better performance in terms of all metrics: approximately 30% increases in FPS and latency and approximately 45% decrease in memory accesses.

Table 10. ResNet50 performance with and without skip connections on the Reconfigurable DNN architecture. The architecture has 672 processing elements and runs on the ZCU102 development board at 250 MHz.

|  | With skip connections | Without skip connections |
|---|---|---|
| Accuracy (%) | 75.85 | 75.36 |
| Frames per second (FPS) | 28.69 | 37.47 |
| Time per image (s) | 0.035 | 0.027 |
| Latency (s) | 0.244 | 0.187 |
| Memory access per image (Mb) | 140.95 | 92.71 |

## 5  DISCUSSION

With these results in hand, designers can now consider which accuracy versus resource tradeoffs they are willing to make during the hardware-software codesign process.

SKIPREMOVER provides minimal accuracy loss while reducing resource consumption and increasing performance—a win-win scenario. As seen in Sec. 4.1, SKIPREMOVER ResNet-50 is only 0.49% less accurate than the baseline on ImageNet. But, SKIPREMOVER is less effective on deeper NNs (such as QuartzNet-10×5 and ResNet-110). In fact, QuartzNet-10×5 fails to converge when trained under SKIPREMOVER. For such deep NNs trained on difficult tasks like ASR, skip connections are instrumental in training convergence [17]. By removing skip connections, we expect and see a degradation in accuracy for deeper NNs. This degradation is not as drastic for other tasks. For instance, ResNet-110 still converges when trained using SKIPREMOVER, but it is 3.72% less accurate on CIFAR-10 and 9.61% less accurate on CIFAR-100, compared to the original baseline model. We propose this tradeoff between NN size and SKIPREMOVER performance as an additional consideration during design space exploration. In response, SKIPSHORTENER is more suitable for deeper NNs when SKIPREMOVER is less effective. SKIPSHORTENER maintains accuracy comparable to its original skip connection models and reduces resource requirements by up to 34% compared to the traditional skip connection model.

Based on our hls4ml evaluation, designers can extrapolate to their own designs because, as we have shown in Fig. 9 and Fig. 11, the resource usage and savings scale linearly as the number of filters grows. We have also shown that at the higher 16-bit precision, TAILOR provides significant resource reductions, so if designers need more precision, TAILOR's savings will follow. If they need lower 8-bit precision, SKIPREMOVER still manages to lower the 8-bit designs' LUTs by 16% and FFs by 11%. Even SKIPSHORTENER decreases LUTs by 3% despite a 2% increase in FFs, though these smaller resource savings are offset by its overall higher accuracy performance compared with SKIPREMOVER. As a result, it is up to the designer to consider how to best apply TAILOR's codesign methods given their accuracy and resource requirements.

## 5.1 Theoretical Understanding

Prior work investigated why skip connections are so helpful to ResNets. Veit et al [45] argue that ResNets behave like ensembles of smaller subnetworks that vary in depth and allow the NN to train and converge more easily. Li et al [26] and Yao et al [49] show that introducing skip connections make the NN loss landscape to be much smoother and have less nonconvexity. They show that naively removing these skip connections causes an explosion of nonconvexity in the loss landscape, which makes training significantly more difficult. We confirm these results in our ablation studies (Sec. 4.1), as accuracy indeed drops when skip connections are removed naively. With both KD and SkipRemover, we see an improvement in accuracy. Since the student is trying to mimic the teacher's outputs, it is possible that the teacher's outputs guide the student in such a way that prevents the loss landscape from becoming less smooth. Theoretical work from Lin et al [27] has proven that a ResNet with one-neuron hidden layers is a universal approximator. This work suggests that adding more neurons to the hidden layers creates an over-parameterized ResNet. Since stochastic gradient descent performs better in the presence of over-parameterization, having more neurons per hidden layer increases training efficiency, making it easier to converge. This work also argues that a ResNet is essentially a sparse version of a fully connected NN because the identity skip connections create simpler paths within the ResNet, which was similarly posited by Lin et al [27]. Given that CNNs and ResNets have both been proven to be universal approximators [27, 40], this implies that there exists a set of parameters for a CNN that can mimic a ResNet such that they equal the same function. It is mainly easier to find a well performing ResNet because Lin et al [27] showed that one-neuron hidden layers is sufficient for a ResNet to be a universal approximator.

## 5.2 Future Work

In our work, Tailor has taken removing and shortening skip connections to their extremes: it either fully removes or fully shortens all the skip connections in a NN. It would be worthwhile to understand the accuracy versus resource utilization tradeoff under less extreme cases, e.g., removing only half of the skip connections. It would also be interesting to mix SkipRemover and SkipShortener together to try and recover accuracy in the instances when SkipRemover fails. These approaches may help address SkipRemover's scalability issues and strike a balance between SkipShortener's high accuracy and SkipRemover's resource savings and performance improvements.

## 6 CONCLUSION

Tailor introduces two new methods, SkipRemover and SkipShortener, that alters NNs with skip connections dynamically during retraining to fit better on hardware, achieving resource-efficient inference with minimal to no loss in accuracy. With SkipRemover, NNs no longer need to rely on skip connections for high accuracy during inference. With SkipShortener, we retrain NNs to use shorter skip connections with minimal to no loss in accuracy. Shortening skip connections is beneficial for hardware architectures generated by the hls4ml tool as it reduces the skip connection lifetime. We demonstrate FPGA resource consumption reductions of up to 34% for BRAMs, 13% for FFs, and 16% for LUTs. We show that Tailor is also valuable for optimizing 2D PE array architectures. SkipRemover increases performance by 30% and decreases memory bandwidth by 45%. Designers can decide which accuracy versus resource tradeoffs offered by SkipRemover and SkipShortener are suitable to their design requirements. As a result, Tailor is another tool in the hardware-software codesign toolbox for designers to use when building customized accelerators.

## ACKNOWLEDGMENTS

## REFERENCES

[1] 2023. Tailor. https://github.com/oliviaweng/tailor.

[2] Thea Aarrestad et al. 2021. Fast convolutional neural networks on FPGAs with hls4ml. *Mach. Learn.: Sci. Technol.* 2, 4 (2021), 045015. https://doi.org/10.1088/2632-2153/ac0ea1 arXiv:2101.05108 [cs.LG]

[3] Colby Banbury, Vijay Janapa Reddi, Peter Torelli, Jeremy Holleman, Nat Jeffries, Csaba Kiraly, Pietro Montino, David Kanter, Sebastian Ahmed, Danilo Pau, et al. 2021. Mlperf tiny benchmark. *arXiv preprint arXiv:2106.07597* (2021).

[4] Y. Bengio, P. Simard, and P. Frasconi. 1994. Learning long-term dependencies with gradient descent is difficult. *IEEE Trans. Neural Netw.* 5, 2 (March 1994), 157. https://doi.org/10.1109/72.279181

[5] Hendrik Borras et al. 2022. Open-source FPGA-ML codesign for the MLPerf Tiny Benchmark. In *Workshop on Benchmarking Machine Learning Workloads on Emerging Hardware (MLBench)*. arXiv:2206.11791 [cs.LG]

[6] Sung-En Chang, Yanyu Li, Mengshu Sun, Runbin Shi, Hayden K-H So, Xuehai Qian, Yanzhi Wang, and Xue Lin. 2021. Mix and Match: A novel FPGA-centric deep neural network quantization framework. In *2021 IEEE International Symposium on High-Performance Computer Architecture (HPCA)*. IEEE, 208. https://doi.org/10.1109/HPCA51647.2021.00027 arXiv:2012.04240

[7] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. ImageNet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*. 248. https://doi.org/10.1109/CVPR.2009.5206848

[8] Xiaohan Ding, Xiangyu Zhang, Ningning Ma, Jungong Han, Guiguang Ding, and Jian Sun. 2021. RepVGG: Making VGG-style convnets great again. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 13733. https://doi.org/10.1109/CVPR46437.2021.01352 arXiv:2101.03697

[9] Zhen Dong, Zhewei Yao, Daiyaan Arfeen, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2020. HAWQ-V2: Hessian Aware trace-Weighted Quantization of Neural Networks. 33 (2020), 18518. arXiv:1911.03852 https://proceedings.neurips.cc/paper/2020/file/d77c703536718b95308130ff2e5cf9ee-Paper.pdf

[10] Zhen Dong, Zhewei Yao, Amir Gholami, Michael W Mahoney, and Kurt Keutzer. 2019. HAWQ: Hessian aware quantization of neural networks with mixed-precision. In *Proc. IEEE/CVF International Conference on Computer Vision*. 293. arXiv:1905.03696

[11] Javier Duarte et al. 2018. Fast inference of deep neural networks in FPGAs for particle physics. *J. Instrum.* 13, 07 (2018), P07027. arXiv:1804.06913

[12] Farah Fahim et al. 2021. hls4ml: An Open-Source Codesign Workflow to Empower Scientific Low-Power Machine Learning Devices. In *1st TinyML Research Symposium*. arXiv:2103.05579 [cs.LG]

[13] Yao Fu, Ephrem Wu, Ashish Sirasao, Sedny Attia, Kamran Khan, and Ralph Wittig. 2017. *Deep learning with INT8 optimization on Xilinx devices*. White Paper WP486. https://docs.xilinx.com/v/u/en-US/wp486-deep-learning-int8

[14] Amir Gholami, Sehoon Kim, Zhen Dong, Zhewei Yao, Michael W Mahoney, and Kurt Keutzer. 2021. A survey of quantization methods for efficient neural network inference. (2021). arXiv:2103.13630

[15] Xavier Glorot and Yoshua Bengio. 2010. Understanding the difficulty of training deep feedforward neural networks. In *Proc. 13th International Conference on Artificial Intelligence and Statistics*, Yee Whye Teh and Mike Titterington (Eds.), Vol. 9. 249. https://proceedings.mlr.press/v9/glorot10a.html

[16] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. 2011. Deep Sparse Rectifier Neural Networks. In *Proc. 14th International Conference on Artificial Intelligence and Statistics*, Geoffrey Gordon, David Dunson, and Miroslav Dudík (Eds.), Vol. 15. 315. http://proceedings.mlr.press/v15/glorot11a.html

[17] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Deep Residual Learning for Image Recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 770. https://doi.org/10.1109/CVPR.2016.90 arXiv:1512.03385

[18] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. 2016. Identity Mappings in Deep Residual Networks. In *ECCV 2016*, Bastian Leibe, Jiri Matas, Nicu Sebe, and Max Welling (Eds.). 630. https://doi.org/10.1007/978-3-319-46493-0_38 arXiv:1603.05027

[19] Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. 2015. Distilling the Knowledge in a Neural Network. (2015). arXiv:1503.02531

[20] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. 2017. Densely connected convolutional networks. In *Proc. IEEE conference on computer vision and pattern recognition*. 4700. https://doi.org/10.1109/CVPR.2017.243 arXiv:1608.06993

[21] Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. 2018. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proc. IEEE conference on computer vision and pattern recognition*. 2704. https://doi.org/10.1109/CVPR.2018.00286

[22] Leonardo Rezende Juracy, Rafael Garibotti, Fernando Gehm Moraes, et al. 2023. From CNN to DNN Hardware Accelerators: A Survey on Design, Exploration, Simulation, and Frameworks. *Foundations and Trends® in Electronic Design Automation* 13, 4 (2023), 270–344.

[23] Samuel Kriman, Stanislav Beliaev, Boris Ginsburg, Jocelyn Huang, Oleksii Kuchaiev, Vitaly Lavrukhin, Ryan Leary, Jason Li, and Yang Zhang. 2020. QuartzNet: Deep automatic speech recognition with 1D time-channel separable convolutions. In *2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 6124. https://doi.org/10.1109/ICASSP40776.2020.9053889

[24] Alex Krizhevsky. 2009. Learning Multiple Layers of Features from Tiny Images. *Tech Report* (2009).

[25] Guilin Li, Junlei Zhang, Yunhe Wang, Chuanjian Liu, Matthias Tan, Yunfeng Lin, Wei Zhang, Jiashi Feng, and Tong Zhang. 2020. Residual distillation: Towards portable deep neural networks without shortcuts. *Advances in Neural Information Processing Systems* 33 (2020), 8935. https://proceedings.neurips.cc/paper/2020/file/657b96f0592803e25a4f07166fff289a-Paper.pdf

[26] Hao Li, Zheng Xu, Gavin Taylor, Christoph Studer, and Tom Goldstein. 2018. Visualizing the loss landscape of neural nets. *Advances in neural information processing systems* 31 (2018).

[27] Hongzhou Lin and Stefanie Jegelka. 2018. Resnet with one-neuron hidden layers is a universal approximator. *Advances in neural information processing systems* 31 (2018).

[28] Ningning Ma, Xiangyu Zhang, Hai-Tao Zheng, and Jian Sun. 2018. ShuffleNet v2: Practical guidelines for efficient CNN architecture design. In *Proc. European conference on computer vision (ECCV)*. 116. https://doi.org/10.1109/ASPCON49795.2020.9276669 arXiv:1807.11164

[29] Y. Ma, Y. Cao, S. Vrudhula, and J. Seo. 2018. Optimizing the Convolution Operation to Accelerate Deep Neural Networks on FPGA. *IEEE Trans Very Large Scale Integr. VLSI Syst.* 26, 7 (2018), 1354. https://doi.org/10.1109/TVLSI.2018.2815603

[30] Y. Ma, M. Kim, Y. Cao, S. Vrudhula, and J. Seo. 2017. End-to-end scalable FPGA accelerator for deep residual networks. In *2017 IEEE International Symposium on Circuits and Systems (ISCAS)*. 1. https://doi.org/10.1109/ISCAS.2017.8050344

[31] Seyed Iman Mirzadeh, Mehrdad Farajtabar, Ang Li, N. Levine, A. Matsukawa, and H. Ghasemzadeh. 2020. Improved Knowledge Distillation via Teacher Assistant. In *AAAI*, Vol. 34. 5191. https://doi.org/10.1609/aaai.v34i04.5963 arXiv:1902.03393

[32] Ricardo Pio Monti, Sina Tootoonian, and Robin Cao. 2018. Avoiding Degradation in Deep Feed-Forward Networks by Phasing Out Skip-Connections. *Artificial Neural Networks and Machine Learning (ICANN)* 11141 (2018). https://doi.org/10.1007/978-3-030-01424-7_44

[33] Bert Moons, Koen Goetschalckx, Nick Van Berckelaer, and Marian Verhelst. 2017. Minimum Energy Quantized Neural Networks. In *2017 51st Asilomar Conference on Signals, Systems, and Computers*. 1921. https://doi.org/10.1109/ACSSC.2017.8335699 arXiv:1711.00215 [cs.NE]

[34] Vinod Nair and Geoffrey E. Hinton. 2010. Rectified Linear Units Improve Restricted Boltzmann Machines. In *Proc. 27th International Conference on Machine Learning*. 807. https://icml.cc/Conferences/2010/papers/432.pdf

[35] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y. Ng. 2011. Reading Digits in Natural Images with Unsupervised Feature Learning. *NIPS Workshop on Deep Learning and Unsupervised Feature Learning* (2011).

[36] Sinno Jialin Pan and Qiang Yang. 2009. A survey on transfer learning. *IEEE Trans. Knowl. Data Eng.* 22, 10 (2009), 1345. https://doi.org/10.1109/TKDE.2009.191

[37] Vassil Panayotov, Guoguo Chen, Daniel Povey, and Sanjeev Khudanpur. 2015. LibriSpeech: An ASR corpus based on public domain audio books. In *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. 5206. https://doi.org/10.1109/ICASSP.2015.7178964

[38] Alessandro Pappalardo. 2022. *Xilinx/brevitas*. https://doi.org/10.5281/zenodo.3333552

[39] Adam Paszke, Sam Gross, Francisco Massa, Adam Lerer, James Bradbury, Gregory Chanan, Trevor Killeen, Zeming Lin, Natalia Gimelshein, Luca Antiga, Alban Desmaison, Andreas Kopf, Edward Yang, Zachary DeVito, Martin Raison, Alykhan Tejani, Sasank Chilamkurthy, Benoit Steiner, Lu Fang, Junjie Bai, and Soumith Chintala. 2019. PyTorch: An Imperative Style, High-Performance Deep Learning Library. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.). Vol. 32. 8024. arXiv:1912.01703 http://papers.neurips.cc/paper/9015-pytorch-an-imperative-style-high-performance-deep-learning-library.pdf

[40] Philipp Petersen and Felix Voigtlaender. 2020. Equivalence of approximation by convolutional neural networks and fully-connected networks. *Proc. Amer. Math. Soc.* 148, 4 (2020), 1567–1581.

[41] Bharat Bhusan Sau and Vineeth N. Balasubramanian. 2016. Deep Model Compression: Distilling Knowledge from Noisy Teachers. (2016). arXiv:1610.09650

[42] Jordi Silvestre-Ryan and Ian Holmes. 2021. Pair consensus decoding improves accuracy of neural network basecallers for nanopore sequencing. *Genome Biol.* 22 (2021), 38. https://doi.org/10.1186/s13059-020-02255-1

[43] Karen Simonyan and Andrew Zisserman. 2015. Very Deep Convolutional Networks for Large-Scale Image Recognition. In *3rd International Conference on Learning Representations, ICLR 2015, San Diego, CA, USA, May 7-9, 2015, Conference Track Proceedings*, Yoshua Bengio and Yann LeCun (Eds.). http://arxiv.org/abs/1409.1556

[44] Antti Tarvainen and Harri Valpola. 2017. Mean Teachers Are Better Role Models: Weight-Averaged Consistency Targets Improve Semi-Supervised Deep Learning Results. In *Advances in Neural Information Processing Systems*, I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Eds.), Vol. 30. 1195. arXiv:1703.01780 https://proceedings.neurips.cc/paper/2017/file/68053af2923e00204c3ca7c6a3150cf7-Paper.pdf

[45] Andreas Veit, Michael J Wilber, and Serge Belongie. 2016. Residual networks behave like ensembles of relatively shallow networks. *Advances in Neural Information Processing Systems* 29 (2016). arXiv:1605.06431 https://proceedings.neurips.cc/paper/2016/file/37bc2f75bf1bcfe8450a1a41c200364c-Paper.pdf

[46] Kuan Wang, Zhijian Liu, Yujun Lin, Ji Lin, and Song Han. 2019. HAQ: Hardware-aware automated quantization with mixed precision. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition.* 8612. https://doi.org/10.1109/CVPR.2019.00881 arXiv:1811.08886

[47] Karl Weiss, Taghi M Khoshgoftaar, and DingDing Wang. 2016. A survey of transfer learning. *J. Big Data* 3, 1 (2016), 1. https://doi.org/10.1186/s40537-016-0043-6

[48] Yifan Yang, Qijing Huang, Bichen Wu, Tianjun Zhang, Liang Ma, Giulio Gambardella, Michaela Blott, Luciano Lavagno, Kees Vissers, John Wawrzynek, et al. 2019. Synetgy: Algorithm-hardware co-design for convnet accelerators on embedded FPGAs. In *Proc. 2019 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays.* 23. https://doi.org/10.1145/3289602.3293902

[49] Zhewei Yao, Amir Gholami, Kurt Keutzer, and Michael W Mahoney. 2020. Pyhessian: Neural networks through the lens of the hessian. In *2020 IEEE international conference on big data (Big data).* IEEE, 581–590.

[50] Sergey Zagoruyko and Nikos Komodakis. 2017. DiracNets: Training Very Deep Neural Networks Without Skip-Connections. (2017). arXiv:1706.00388

[51] Sergey Zagoruyko and Nikos Komodakis. 2018. DiracNets: Training Very Deep Neural Networks Without Skip-Connections. (2018). arXiv:1706.00388

[52] Fuzhen Zhuang, Zhiyuan Qi, Keyu Duan, Dongbo Xi, Yongchun Zhu, Hengshu Zhu, Hui Xiong, and Qing He. 2020. A comprehensive survey on transfer learning. *Proc. IEEE* 109, 1 (2020), 43. https://doi.org/10.1109/JPROC.2020.3004555