

# A Survey of Deep Learning on Mobile Devices: Applications, Optimizations, Challenges, and Research Opportunities

Tianming Zhao\*, Yucheng Xie<sup>†</sup>, Yan Wang\*, Jerry Cheng<sup>‡</sup>, Xiaonan Guo<sup>†</sup>, Bin Hu<sup>§</sup> and Yingying Chen<sup>§</sup>

\*Temple University, Philadelphia, PA, USA 19122

<sup>†</sup>Indiana University–Purdue University Indianapolis, Indianapolis, IN, USA 46202

<sup>‡</sup>New York Institute of Technology, New York, NY, USA 10023

<sup>§</sup>Rutgers University, New Brunswick, NJ, USA 08901

Email: tum94362@temple.edu, yx11@iupui.edu, y.wang@temple.edu, jcheng18@nyit.edu, xg6@iupui.edu, bh439@scarletmail.rutgers.edu, yingche@scarletmail.rutgers.edu

**Abstract**—Deep learning (DL) has demonstrated great performance in various applications on powerful computers and servers. Recently, with the advancement of more powerful mobile devices (e.g., smartphones and touch pads), researchers are seeking DL solutions that could be deployed on mobile devices. Compared to traditional DL solutions using cloud servers, deploying DL on mobile devices have unique advantages in data privacy, communication overhead, and system cost. This paper provides a comprehensive survey for the current studies of adopting and deploying DL on mobile devices. Specifically, we summarize and compare the state-of-the-art DL techniques on mobile devices in various application domains involving vision, speech/speaker recognition, human activity recognition, transportation mode detection, and security. We generalize an optimization pipeline for bringing DL to mobile devices, including model-oriented optimization mechanisms (e.g., pruning, quantization, etc.) and non-model oriented optimization mechanisms (e.g., software accelerator and hardware design). Moreover, we summarize popular DL libraries regarding their support to state-of-the-art models (software) and processors (hardware). Based on our summarization, we further provide insights into potential research opportunities for developing DL for mobile devices.

**Index Terms**—Deep Learning, Mobile Sensing, Mobile Security, Optimization, Hardware and Software Accelerator Design

## I. INTRODUCTION

Recent developments of deep learning (DL) have made huge leaps in various mobile applications, such as image recognition [1]–[7], speech recognition [8]–[11], human activity recognition [12]–[17], etc. In general, DL could provide much higher inference accuracy than traditional machine-learning methods. However, DL usually involves complicated neural network models and possibly millions of parameters that require high computational power and large memories, which are not available on mobile devices. As a result, many mobile applications choose to use cloud computing technology to offload all or partial computational tasks of DL from resource-constraint mobile devices. While the cloud-based approaches are common and have initial success in many applications, people have privacy concerns because sensitive data need to be uploaded to cloud servers. Besides, the latency of cloud-based approaches is highly affected by network conditions and

the interfaces (e.g., LTE, WiFi, etc.) [18]. In addition, due to high maintenance costs, cloud computing services are primarily provided by a limited number of large companies that have the resources to support large systems [19]. Therefore, cloud-based approaches can no longer satisfy the emerging artificial intelligent (AI) and augmented reality (AR) mobile applications, new technologies aiming at enabling DL on mobile hardware are much in demanded.

Many research work has been conducted to develop DL-based approaches that are optimized for mobile devices. Topics of them span over a broad range, including miniature DL model design, DL model optimization, software accelerator, and hardware accelerator. A couple of earlier studies of DL have conducted for mobile device and applications. Nan *et al.* [57] explored typical DL compression methods and compared their practical performances using various metrics. Based on their performances, they discussed the advantages and bottlenecks of those compression methods. Solutions to reduce computational cost of DL with respect to algorithmic design, computational optimization, and hardware revolution have been explored in a more recent survey [58]. Moreover, some survey papers [59], [60] compare various attributes (e.g., latency, energy efficiency, costs, and scalability) of different mobile devices using DL and propose several evaluation criteria to choose appropriate development frameworks. In addition, Zhang *et al.* [61] present a survey of mobile and wireless networking research based on deep learning. They review recent deep learning applications to mobile and wireless networking and discuss how to tailor deep learning models for those problems. Chen *et al.* [62] has summarized DL in edge computing with optimization methods for fast inference on resource-constraint edge devices. Deng *et al.* [63] provide more optimal solutions to key problems in edge computing with the help of deep learning technologies and study how to carry out the entire process of building deep learning models on the edge. Deng [62] provides an accessible introduction to various aspects on the recent advancements of deep learning on mobile devices including mobile hardware architectures, optimization approaches, existing resources, mobile sensing

TABLE I  
IMAGE CLASSIFICATION AND SPEECH RECOGNITION APPLICATIONS ON MOBILE DEVICES.

Application	Image Classification/Object Detection		Speech/Speaker Recognition	
	Facial expression recognition, Road damage/Crack detection, etc. [1]–[7], [20]–[45]		Speech enhancement, Voice activity/Noise detection, etc. [8]–[11], [46]–[56]	
Adopted Model	Types	Complexity (architecture/layers)	Types	Complexity (architecture/layers)
	CNN (Image Classification)	GoogLeNet [20], [24], [34]; ResNet50 [22], [32], [36]; MobileNet V1, V2 [36]; DarkNet-19 [37]; VGG 16,19 [32]; AlexNet [20], [38]; InceptionV3 [32], [36]; CNN 5 [26], 7 [25], 8 [30], 9 [27], 12 [28] layers	CNN	SqueezeNet V1.1 [56], TrimNet [8]; CNN 2 [56], 3 [56], 5 [46], 8 [47], 15 [52], 21 [51] layers
	CNN (Object Detection)	YOLO v2 [23], [29], [37]; SSD [29], [33]; NIN [24]; Fast-RCNN [21], [29]; AlexNet [24]; InceptionV3 [22]; VGG 19 [22]	DNN	3 [55], 5 [49], 6 [55], 8 [9]
			RNN(LSTM)	3 [54], 5 [53], 8 [10], 17 [11]
Adopted DL Library	Caffe [3], [4], [20], [24], [30], [31], [38], [40] Tensorflow [2], [33], [35] Keras [32] Theano [25], [39] Torch [41]		Tensorflow [46], [47], [50] Keras [51] Theano [49]	
Adopted Optimization	Compression (pruning) [1], [3], [39], [40]		Compression(lower rank matrices) [10], [53], [54]	
	Quantization [4], [40], [43]		Quantization [9], [11]	
	Model scaling [5], [7], [44]			
	Programming strategies [4]		Programming strategies [9]	

applications, and challenges. While these works have looked into various aspects in DL optimized for edge and mobile devices, they either mainly discuss the DL optimization (e.g., model parallelism) on edge devices for applications such as mobile and wireless network or provide the introduction-level survey that is dedicated for quickly grasping the concepts. However, a comprehensive study of the technologies that could enable DL on various mobile hardware is still missing. We propose an optimization pipeline that is designed to adaptively support DL deployment on various kinds of mobile device for different application domains. This will provide an in-depth guidance for the researchers to customize their DL-based applications deployment.

This paper extensively studies the state-of-the-art literature that focuses on DL optimization for mobile applications and devices. We summarize DL-based approaches in different mobile applications in terms of evaluation platforms, adopted models, and optimization methods. In addition, we present the challenges of applying DL on mobile devices and introduce an optimization pipeline for guiding the optimization of DL targeting mobile devices and applications. Specifically, we divide the optimization pipeline into two parts (i.e., model-oriented and non-model-oriented) which cover typical optimization approaches with respect to both hardware and software. Moreover, we summarize software and libraries used in DL in terms of software and hardware optimization support, which could help us choose appropriate software and hardware combinations for their application development. We also point out the current trend and potential direction of research on the DL optimization for various mobile applications and resource-constrained mobile devices. It is worth noting that the mobile devices mentioned in this survey include a wide range of devices from resource-constrained IoT devices to high-end smartphones.

The remaining of this paper is organized as follows:

Section II summarizes DL in different mobile application. Section III discusses the challenges of applying DL on mobile devices and summarizes the optimization approaches. Section IV compares the state-of-art DL libraries. Section V discusses the on-device optimization in federated learning. Section VI points out the potential research in different mobile applications. Section VII concludes the paper.

## II. DL-ENABLED MOBILE APPLICATIONS

Many mobile applications adopt DL-based approaches due to their better performance than traditional machine learning-based approaches. In this section, we review the DL approaches developed for several popular application domains, including image classification and object detection, speech and speaker recognition, activity recognition, and mobile security. Our summary mainly focuses on the following aspects: the adopted evaluation platform, the adopted DL model, and the adopted performance improvement method.

### A. DL in Image Classification and Object Detection

In the past few years, decent performance has been achieved in image classification and object detection using DL approaches on mobile platforms (e.g., smartphones [1]–[7], embedded systems including NVIDIA Jetson TK1, Intel Edison platform, Raspberry pi [6], [39]–[41], [43], [67]–[69], and the mobile development kit such as Snapdragon 820, Snapdragon 845 [43]–[45]). We summarize the popular DL models for image classification and object detection in Table I. In particular, VGG, MobileNet V1 and V2, AlexNet, LeNet, SqueezeNet, GoogleNet, ShuffleNet are usually trained with popular database (e.g., ImageNet 2012 [1], [3]–[6], [40], [44], [70], [71] and MNIST [39], [40], [70]). For object detection, popular DL models such as SSD300, SSD512, YOLOv2, YOLOv3, YOLOv3-Tiny, SSD+MobileNet, SSDlite+MobileNet v2 are trained with the databases such as COCO [5], [6], [45],

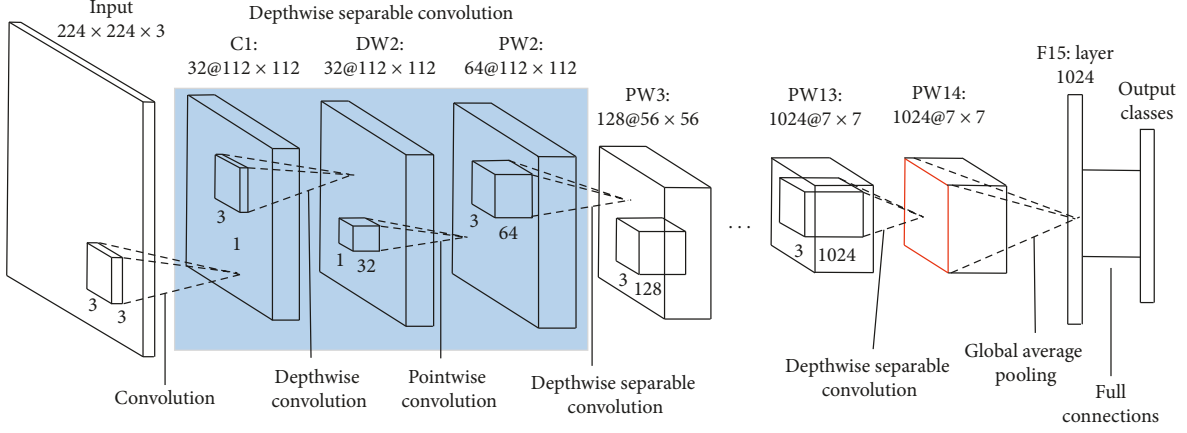


Fig. 1. The structure of MobileNet is based on depthwise separable filters for Image Classification [64]. MobileNet is a streamlined architecture that uses depthwise separable convolutions to construct lightweight deep convolutional neural networks and provides an efficient model for mobile and embedded vision applications. The model is composed of standard convolution layer (C), depthwise convolution layers (DW), pointwise convolution layers (PW), and full connection layers (F). The model also adopts a novel structure called depthwise separable convolution layer, which is a depthwise convolution layer followed by a pointwise convolution layer. Such a design could significantly reduce the model's parameters compared to the models composed of regular convolutions with the same depth, making MobileNet nearly as accurate as VGG16 [65] but 32 times smaller-sized and 27 times less compute-intensive. Moreover, MobileNet can achieve 4% better accuracy than AlexNet [66] with a 45 times smaller size and 9.4 times less computation.

VOC [6], [43], [45], Cifar [2], [39]. DL approaches have also been adopted in a range of special applications, such as CT images processing [72] and water ecosystem images captured by IoT devices [73]. Among all these models, some DL models (e.g., MobileNet in Figure 1 and SSD) are the most popular and usually adopted as the baseline for performance comparison with newly proposed models (e.g., PeeletNet, Enet, ThunderNet, etc) [2], [4]–[6], [41], [42], [44], [45], [68], [71].

Recently, new DL models are becoming faster, smaller, and more power efficient. For example, Mnasnet [5] adopted an automated mobile neural architecture search (MNAS) approach to achieve 1.8 times faster than MobileNetV2, PeeletNet [6] was 1.8 times faster than MobileNetV2 by using dynamic number of channels in bottleneck layers, Shufflenet [71] achieved 13 times faster than AlexNet by utilizing two new operations, pointwise group convolution and channel shuffle, to greatly reduce computation cost. Moreover, Mobiface [42] adopted two design strategies, bottleneck residual block with expansion layers and fast downsampling, to have its model size 10 times smaller than MobilenetV1. However, most existing DL models still have big sizes and complicated structure, making it hard to perform real-time image classification and object detection on resource-constrained mobile devices. To this end, compression and quantization methods have been adopted to reduce the complexity and improve the inference speed of DL models. In particular, AutoML [1] can reduce the model size of MobileNet by 34% from leveraging reinforcement learning to provide the model compression policy. DeepRebirth [3] achieved more than 3 times faster speed and 2.5 times run-time memory saving on GoogLeNet by the compression technique containing two operations (i.e., streamline slimming and branch slimming). DeepIoT [39] compressed neural network structures into smaller dense matrices by finding the minimum number of non-redundant hidden elements. It is able to reduce the model size of VGGNet by 97.6%, running time by 94.5%, and energy consumption by 95.6%. Deep compression [40]

pruned unimportant connections, quantized the network using weight sharing, and then applied Huffman coding on AlexNet, VGG-16, and LeNet to reduce their weight storage by 35, 49, 39 times, respectively. Along this line, model scaling methods are adopted by some work [5], [7], [44], which use the depth multiplier to change the width of channel of each layer and reduce the inference time. Moreover, advanced programming strategies (e.g., Fused Tile Partitioning (FTP)) were adopted by Zhao *et al.* [110] to reduce memory footprint of YOLOv2 by more than 68% without sacrificing accuracy.

Although these methods can improve DL performance on mobile devices, some work, for example, GoogLeNet and Thunder [45], still cannot provide real-time performance on mobile devices. The compressed GoogLeNet model has been shown good accuracy on low-end mobile devices (e.g., Moto E) [3], [4] but they can only process video streams at 10fps. Thunder with SNet535 targeted at achieving high accuracy but can only process 5.8 frames per second, which is not fast enough in real-time. Therefore, achieving both high accuracy and near real-time inference using DL models on the mobile devices is still a challenging research problem.

### B. DL in Speech and Speaker Recognition

In the speech/speaker recognition area, most existing work [51], [53]–[56], [111] claim that their proposed systems has the potential for mobile devices due to their less resource consumption. However, only few of them [8]–[11] have demonstrated their practicality in mobile applications. We summarize the models, software platform, and optimization methods of the existing work in speech and speaker recognition as shown in Table I.

In the speech/speaker recognition area, RNN architectures are commonly used because they allow extracting past and future dependencies at a given point of the audio. Some work [9], [11], [51], [52], [55] propose simple light-weight

TABLE II  
ACTIVITY RECOGNITION AND MOBILE SECURITY APPLICATIONS ON MOBILE DEVICES.

Application	Activity Recognition		Mobile Security	
	Human Activity Recogniton [12]–[15], [74]–[94] Transportation Mode Detection [95]–[100]		Malware Detection [101]–[107] Traffic Classification [108]	
Adopted Model	Types	Complexity (layers)	Types	Complexity (layers)
	CNN	4 [76], 5 [90], 6 [74], [79], [81], 7 [88], 8 [75], 10 [77], 13 [96], [97], 16 [85]	CNN	4 [102], 7 [103], [104], [108]
	DNN	5 [99], 10 [98]	DNN	3 [99]
	MLP	5 [83], 10 [84]	MLP	1 [108], 3 [101]
	SAE	4 [78]	SAE	5 [108]
	DBN	4 [76]	ANN	3 [106]
	RNN(LSTM)	6, 7 [80]	RNN(LSTM)	5 [108]
Adopted Library	DL			
	Tensorflow [75], [79], [89], [96], [97], [100] Keras [77], [84] Matlab tool [78], [85], [86] Torch [82] Pylearn [89] H2O (R) [92]		Tensorflow [102], [104], [108]	
Adopted Optimization	Compression [109] Quantization [13]		Quantization [104]	

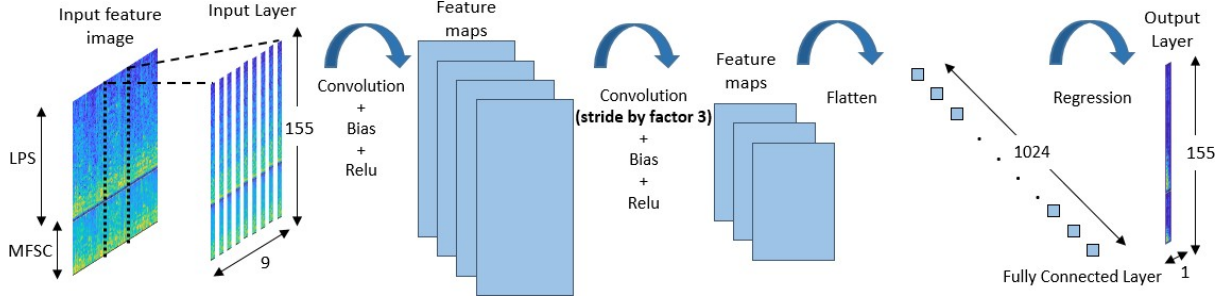


Fig. 2. CNN-based Speech/Speaker Recognition DL Model Example [47]. The architecture has 3 hidden layers, 2 convolutional layers, and 1 fully connected layer. The audio input is represented as STFT spectrograms and is fed as an input to the CNN for audio classification tasks.

neural networks (e.g., around 5-layer CNN, DNN, RNN) as shown in Figure 2, and compare with classic baseline models. For example, He *et al.* [11] proposed a compact E2E speech recognizer, which improved the word error rate by more than 20% over a conventional CTC-based baseline system; Deep KWS [55] achieved 45% relative improvement in false reject rate with respect to a competitive baseline (HMM). These initial results reveal advantages and potential of developing novel DL models for speech and speaker recognition. In addition, we observe that most of the existing studies in speech recognition are evaluated against baseline models, comprised of some classical machine learning models. In the future, DL models can be used as baseline standards.

Several methods have been adopted to improve performance in speech/speaker recognition, including compression (e.g., lower rank matrices) [10], [53], [54], quantization (e.g., float-point to 8-bit) [9], [11], light-weight neural network models (e.g., around 5-layer CNN, DNN, RNN) [9], [10], [52]–[56], and programming strategies (i.e., multi-threads, frame skipping, etc.) [9]. After adopting those methods, the proposed DL model size is reduced several times compared with the original model, and the inference time is also reduced significantly. For example, the model size reduces to 1/3 times [9], [53] or 1/4 times [11], [54] of the original size; the inference time is claimed to be many times faster than the original model [10], [53], [55], [56]; and the power consumption is reduced to

an acceptable level [111]. However, since most of these DL models have not been deployed in mobile platforms, and their performance improvements are only compared with the original DL models, it is hard to guarantee these models are suitable for mobile devices. Therefore, future research studies are needed to fully optimize the DL models to satisfy the hardware resource constraint and real-time application nature in speech or speaker recognition on mobile devices.

### C. DL in Activity Recognition

Activity recognition is another active research area in mobile computing, which could facilitate many emerging applications including AI and AR/VR on mobile devices. Existing researchers have adopted DL in activity recognition applications (e.g., human activity recognition and transportation mode detection) on smartphones [12]–[17], and embedded devices (e.g., Raspberry PI, MinnowBoard Turbot Dual-Core Board) [93], [94], [109]. We reviewed the recent literature on activity recognition and summarize them in Table II based on their complexity, software platform, and optimization methods.

Unlike the image classification area that always adopts complex DL models, the activity recognition domain usually adopts customized light-weight DL models, as shown in Figure 3. For example, some adopt CNN, DNN, LSTM, etc.

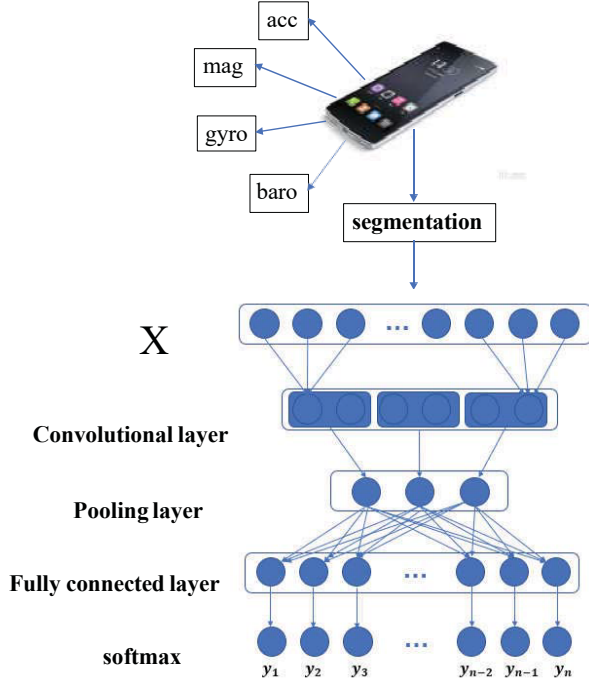


Fig. 3. CNN-based Activity Recognition DL Model Example [90]. The network consists of three types of layers including convolutional layer, pooling layer and fully connected layer. The input data is a vector containing data collected from accelerometer (acc), magnetometer (mag), gyroscope (gyro) and barometer (baro).

with less than 10 layers, as shown in Table II. Although these models achieve good performance, their model structures and databases are mostly not publicly available, thus it is hard to compare the performance among these approaches. Therefore, open-sourced DL models and datasets are highly demanded.

We also observe that the majority of DL-based activity recognition systems are developed and evaluated based on personal computers without considering the limited resources in mobile devices. Only few researchers have explored to improve the performance of DL models on mobile devices. As shown in Table II, Zebin *et al.* [13] used quantization in their method, and Stardust [109] compressed the structure of DL models to fit them into mobile devices. Some work [12]–[14], [93] adopted light-weight DL models (e.g., 3,4,8-layer CNN, 2-layer LSTM) to satisfy the requirements of their proposed applications. Since some light-weight DL models already achieved good performances (e.g., over 95% accuracy in five-class activity recognition [13]), further optimization is not vital to them. However, for complex activity recognition applications, DL models usually have complicated and huge-sized structures that cannot be miniaturized by using a simple method. Thus, new research on how to use the combination of multiple optimization methods (e.g., factorization and pruning) to reduce the size of such complex DL models while achieving high accuracy is essential for deploying complex DL models on mobile devices.

#### D. DL in Mobile Security

With the ever-increasing use of mobile and IoT devices, security is becoming more important these days. Therefore,

DL models for mobile security are also evolving quickly with recent development [101], [102], [104] for different resource constraints. Some of them are simulated using FPGA emulator (e.g., Xilinx Zynq-7000) [106] and Android emulator [105]. The existing work on DL models for mobile security are summarized in Table II. The light-weight DL models are often deployed on mobile devices as shown in Figure 4. For example, Yuan *et al.* [101] designed a lightweight on-device Android malware detector which is a three-layer multi-layer perceptron (MLP) neural network that mainly uses one-shot computation for model training. Hence it can be fully or incrementally trained directly on mobile devices (e.g., Samsung Galaxy S9). MalDozer [102] detected Android malware by using a four-layer CNN that is not only deployed on servers but also on IoT devices (e.g., Raspberry Pi 2). Different from aforementioned existing work that need the hand-engineered malware features, McLaughlin *et al.* [103] developed a five-layer CNN-based Android malware detection system which is capable of simultaneously learning to perform feature extraction and malware classification. Even though those models provide decent accuracy on mobile devices, they are not generally applicable to various mobile security applications since they are only trained and tested on a few malware datasets (e.g., VirusShare and Contagio, Malgenome and Drebin, Android malware repository from McAfee Labs are adopted by the work [101], [102], [103], respectively). Also, the proposed models are mostly not available. Considering the varying nature of malware in real scenarios, public malware datasets and open-sourced DL models are desirable in mobile security.

In addition, DL is applied as a viable strategy for network traffic monitoring and analysis applications, e.g., traffic classification and prediction [112]–[116], because a large portion of malware is spreading through the network traffic and hiding activities using various encryption protocols. However, the current state of research on mobile traffic classification has not reached a similar level of DL in other fields [108]. Researches about mobile networking and traffic monitoring have mostly been studied independently. Recently, only a few cross-over studies [107], [108] between the two research areas have emerged. Aceto *et al.* [108] presented a review for mobile traffic classification works that use deep learning techniques. They reproduced several DL classifiers, e.g., MLP, LSTM, CNN, and SAE, from the traffic classification literature and also made a detailed evaluation to compare the accuracy of these classifiers. Moreover, they found that the key issues of DL in traffic classification are the lack of advanced hybrid DL architectures and the lack of up-to-date human-generated public datasets. This problem is further worsened in mobile scenarios because the possibility of sharing large and up-to-date datasets is restricted by both higher privacy concerns and the fast-paced evolution of mobile traffic mix.

Because a timely response to malicious attacks could help prevent users' loss, mobile security applications are susceptible to inference time. Usually, it requires to reduce the model size and computational cost of DL models to achieve real-time detection. However, we find that majority of existing work mainly focused on improving inference accuracy using offline computing. While high accuracy has been achieved, not much



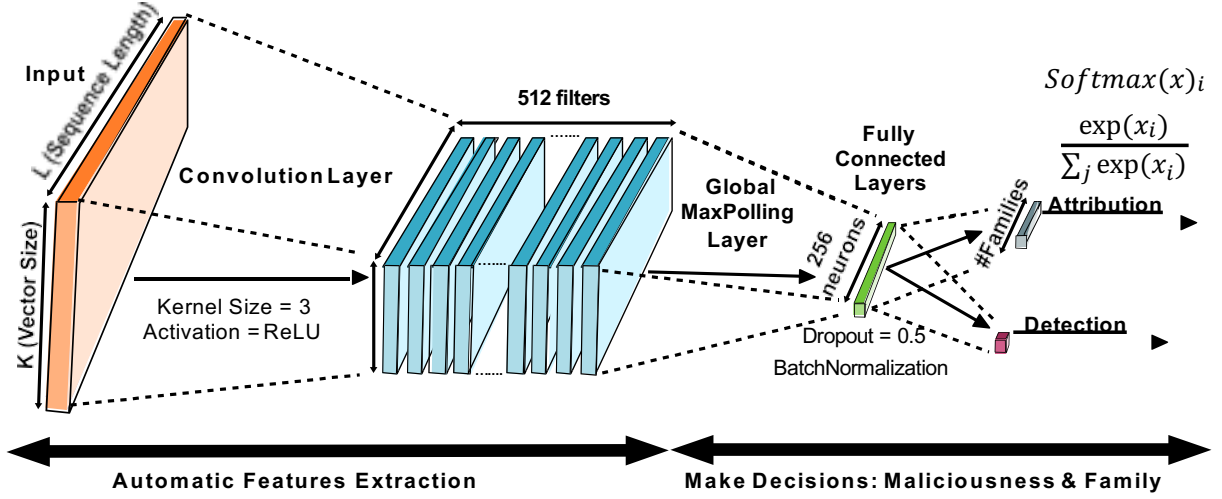


Fig. 4. Security DL Model Example [102]. The model contains a convolution layer followed by a global max pool layer, then connect to a fully-connected layer. In addition to Dropout used to prevent overfitting, they also utilize batch normalization to improve the results.

work has been done to reduce the model size and computational cost of DL models in mobile security. To achieve the real-time nature of the security applications on various mobile devices with different resource constraints, researchers will need to adaptively consider some typical performance improvement methods (e.g., quantization and pruning) in their research. Especially, when more generally applicable models with deeper architecture are adopted, more optimization methods would be needed to achieve real-time on the mobile devices.

### III. TOWARD ENABLING DL ON MOBILE DEVICES

#### A. Current Status and Issues

Currently, most of popular DL-based mobile applications collect sensor data from mobile devices and perform inferences using DL on powerful cloud servers. This is mainly because state-of-the-art DL models usually have complex and bulky structures, which contain hundreds, even thousands of nodes and potentially millions of parameters that cannot be deployed on resource-limited mobile devices. We have reviewed a broad range of literature about enabling DL on mobile devices and summarize current status and issues in this area as follows:

**How to Select/Design a Suitable Model for Heterogeneous Mobile Devices and Applications is Not Systematically Studied?** We have seen many application domains adopt various DL models [1]–[17] on server or desktop, however, there are no general models designed specifically for different domains on mobile devices. As the first step toward enabling DL on mobile devices, the current mobile computing domain needs the capability to select and design DL models suitable for various mobile devices and applications.

**What Kind of DL Optimization Approaches that Could Enable Existing DL Models on Various Mobile Devices is Still Lacking Comprehensive Analysis and Research?** While many researchers have been working on DL optimization in recent years, most of them focus on reducing the model size on a desktop or server that usually have sufficient hardware resources to run various compressed models smoothly.

Differently, due to various hardware resource limitations [117]–[120] on a range of mobile devices, deploying DL models on mobile devices need to consider more aspects (e.g., storage space, inferencing speed, power consumption, and real-time memory usage). No single optimization approach (e.g., weight pruning [121]) can easily meet the requirements for all those aspects. It is critical to find an optimization pipeline that adaptively adopt various exiting optimization approaches for DL deployment on mobile devices considering different application domains.

**How to Utilize the Full Power of Existing Mobile Device Hardware and Design New Dedicated Hardware for DL are Still to be Explored?** Current off-the-shelf mobile devices usually already have the strong computing capability at a low-end laptop. In addition, most of them have integrated various processors that are designed for different purposes, such as CPU, GPU, and DSP. While incorporating dedicated hardware for DL acceleration depends largely on hardware manufacture, it is more practical to have software accelerators on mobile devices to take full advantage of available hardware resources by using hardware and software co-design.

Another way to enable DL on mobile devices is to employ new hardware (e.g., mobile GPU [122], FPGA [123], [124], and ASIC [125]–[128]) dedicated to DL on mobile devices. So far, those designs require the special manufacturing processes. And not many new hardware designs have been done for efficient DL-based applications.

To facilitate the analysis of DL optimization for mobile devices and applications, we propose a DL optimization pipeline as shown in Figure 5. In particular, we identify four major mobile applications that have the most extensive use of DL as the focus of this survey, namely Image Classification/Object Detection, Speech/Speaker Recognition, Activity Recognition, Mobile Security. While DL has already been adopted in a variety of applications to improve their performance, these applications cover state-of-the-art DL design and optimization research. Overall, we categorize research on DL optimization for mobile devices and applications into two types: *Model-*

*Oriented Optimization Mechanisms* and *Non-Model-Oriented Optimization Mechanisms*. For most mobile applications, the first type can reduce DL models' size and inference time. These mechanisms include the *Optimum DL Model Selection*, which explores DL models that are most suitable for a particular type of mobile application and the *DL Model Optimization*, which develops compressing or pruning technologies to reduce the complexity of DL models for mobile deployment. In addition to the first type of optimization, the second type could further boost the performance of DL on mobile devices. Given a particular hardware environment, researchers have developed *Mobile DL Software Accelerators* to decompose a DL model and deploy the decomposed components on the most suitable mobile processing hardware (e.g., CPU, GPU, DSP, etc.). Except software-oriented optimization approaches, there are research work focusing on *Mobile Hardware Designs for DL*, which propose special computer architectures enabling highly efficient DL training and inference with low energy consumption. Next, we introduce our findings in each component of the pipeline.

## B. Model-Oriented Optimization Mechanisms for Mobile Devices and Applications

### 1) Optimum DL Model Selection

With the emerging use of mobile devices in AI or VR applications, research in the computer vision domain have led the trend of using DL on mobile devices. DL models that are specially designed for mobile devices have shown the initial success in image classification and object detection whose performance is comparable to that of the state-of-the-art DL models (e.g., VGG [65] and AlexNet [66]) on personal computers.

**DL Model Selection for Image Classification.** Specifically, MobileNets [64] and SqueezeNet [129] have already shown their effectiveness on mobile devices with power constraints for image classification. In MobileNets, by using the width and resolution multipliers, researchers can trade off a reasonable amount of accuracy and model size to build smaller and faster versions. Particularly, MobileNet is nearly as accurate as VGG16 but 32 times smaller and 27 times less compute intensive. Moreover, it is 4% better than AlexNet, yet 45 times smaller and 9.4 times less compute than AlexNet. While SqueezeNet achieves a model size less than 0.5MB with fifty times fewer parameters than AlexNet, but with same level accuracy. It can be entirely deployed in the FPGA and eliminate the need for off-chip memory accesses to load model parameters.

**DL Model Selection for Object Detection.** As for object detection, a fast single-shot object detector, SSD [130], has been built on MobileNet which can be deployed on the mobile and edge devices [131]. In particular, the prototype SSD could significantly outperform the state-of-the-art object detector counterparts (e.g., Faster R-CNN [132] and YOLOv1 [133]) in terms of both accuracy and speed. For example, the SSD512 (i.e., SSD with 512 x 512 input) model has 10% higher accuracy and three time faster than the Faster R-CNN. The SSD300 (i.e., SSD with 300 x 300 input) model

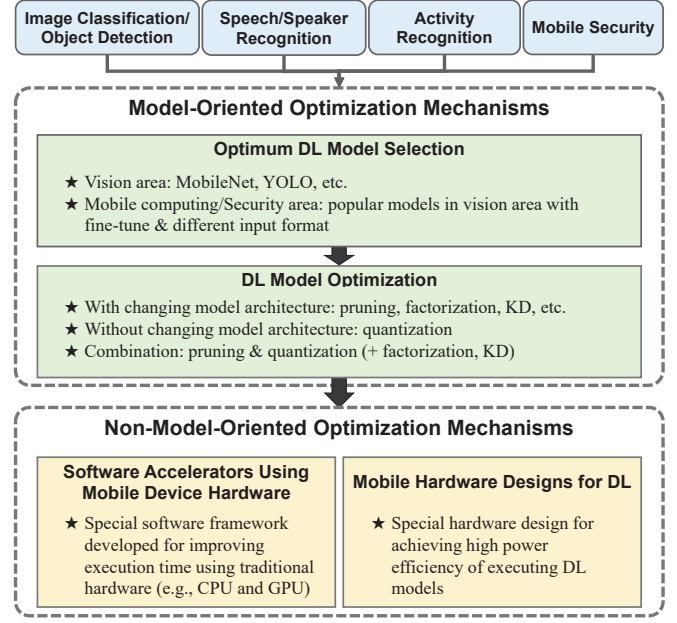


Fig. 5. DL optimization pipeline for mobile applications. We identify four major mobile applications that have the most extensive use of DL as the focus of this survey. Overall, we categorize research on DL optimization for mobile devices and applications into two types. The first one is Model-Oriented Optimization Mechanisms, which includes the Optimum DL Model Selection and DL Model Optimization. The second one is Non-Model-Oriented Optimization Mechanisms, which further boosts the performance of DL via Software Accelerators Using Mobile Device Hardware and Mobile Hardware Designs for DL.

could perform object detection at 59fps, which is about three times faster than the current real-time YOLOv1 model and producing superior detection accuracy of 74.3% on Pascal VOC2007 dataset. Moreover, YOLOv2 [134] was developed as an enhanced version of YOLOv1 whose inference speed could be twice of that of SSD. Particularly, YOLOv2 added the batch normalization on all convolutional layers and adopted anchor boxes as SSD and Faster R-CNN do. And it could also be run at a variety of image sizes to provide a smooth tradeoff between speed and accuracy which has been widely tested on the mobile devices for achieving real-time processing.

**DL Model Selection for Activity Recognition and Security.** While the DL models mentioned above are specially designed for computer vision-based applications, they could be adopted in other domains, such as activity recognition and security. There are two basic steps for adoption. First, according to different data type, it is necessary to change the input format for certain datasets (e.g., the activity data collected from accelerometer and gyroscope sensors in UCI HAR [135] and WISDM dataset [136], network traffic data collected from UNSW-NB15 [137], and CIDDs-001 [138] datasets in security area). Second, base model's hyperparameters (e.g., number of neurons, convolution kernel width, and learning rate) should be tuned based on both accuracy and inference time requirement of different applications and the constraint of adopted mobile hardware. Particularly, there are several benefits of using those state-of-art DL models in vision field. First, those models have been publicly available for usage. Therefore they only need to slightly change the input format. Second, adopting those

models could provide a general evaluation framework among research work. Third, those models should be able to support real-time applications on mobile devices.

The next optimization option is to apply the General DL Optimization Mechanisms. In particular, DL Model Optimization aims at reducing the model size so as to fit into resource-constrained mobile devices. Such methods can be roughly divided into two categories: *With-Changing-Model-Architecture* and *Without-Changing-Model-Architecture*, as is shown in Table III. The former reduces the model size by shrinking the original model architecture, while the latter shrinks the model weight size while remaining the original model architecture.

### 2) DL Model Optimization With Changing Model Architecture.

Several DL optimization methods in *With-Changing-Model-Architecture* have been proposed to change model architectures by either introducing new layers by factorization [139]–[141], or pruning connections between neurons [1], [3], [39], [40] or even generating totally new architecture based on original models [144]–[146].

**Factorization.** Factorization [149] changes the underlying architecture by adding new layers. Specifically, it compresses the weight matrix  $W$  of an original DL model into two lower-rank weight matrix  $W_a$  and  $W_b$  via singular value decomposition (SVD). As long as the summed parameters in two factorized matrix  $W_a$  and  $W_b$  is less than original matrix  $W$ , the model size will be compressed. Furthermore, this technique can uniformly be applied to convolutional layers and fully connection layers. For example, DXTK [139] showed that the DNN model after factorization can runs 9.7 times faster than its unmodified version. The latest work [150] showed that factorization could reduce the model size up to approximately 75% without decreasing the accuracy. However, factorization is computationally expensive and not suitable for large deep models since the computation cost increases exponentially as the layer increases.

**Pruning.** In addition to factorization, pruning could compress DL models by removing small-weight connections below a threshold. For example, Shi *et al.* [151] showed that DNN model after pruning could achieve up to 25.6 times reduction on transmission workload, 6 times acceleration on total computation and 4.81 times reduction on end-to-end latency as compared to the original DNN model without pruning. DeepIoT [39] presented the performance of adopting pruning on DL structures for sensing applications, including fully-connected, convolutional, and recurrent neural networks, as well as their combinations. It reduced the size of deep neural networks by 90% to 98.9%, the execution time by 71.4% to 94.5%, and energy consumption by 72.2% to 95.7% without loss of accuracy. The results show the potential of pruning for deploying deep neural networks on resource-constrained embedded devices. Currently, Tensorflow and Pytorch support weight pruning for both convolution layers and full connected layers. However, pruning increases training time and also generates a sparse matrix which requires some dedicated hardware.

**Knowledge Distillation (KD).** Different from pruning and factorization, Knowledge Distillation method [144] trains a small model (i.e., student model) based on large original models or whole ensemble of models (i.e., teacher model). This is also called a teacher-student network. Knowledge is the distribution of class probabilities predicted by the teacher model based on softmax layer output. In particular, knowledge is transferred to the distilled model (i.e., compressed model) by training it to match each class probability in the output of softmax layer of teacher model. For example, Hinton *et al.* [144] showed that they could significantly improve the model of a heavily used commercial system by distilling the knowledge in an ensemble of models into a single model. Particularly, more than 80% of the improvement in classification accuracy achieved by using an ensemble of 10 DNN models in Automatic Speech Recognition (ASR) was transferred to the distilled model. Ashok *et al.* [146] showed that KD could be combined with pruning to improve compression rate about 2 times than using pruning alone while also maintaining accuracy of the teacher model. However, KD heavily relies on softmax layer and has the strict assumption about the size and structure of teacher-student network. As a result, this method is more suitable for small and middle size models instead of large models.

### 3) DL Model Optimization Without Changing Model Architecture.

The DL optimization methods listed above can reduce the model size and retain the decent accuracy with changes in model architecture. Instead, the model optimization of *Without-Changing-Model-Architecture* can reduce model size without changing the original model architecture and have little degradation in model accuracy.

**Quantization.** It can use quantization to optimize the model by replacing the full precision weight (i.e., 32 bit) into  $n$ -bit precision weight (e.g., binary quantization having 1-bit weight), and reducing the memory by a factor of  $32/n$ . Moreover, quantization can speed up inference since the costly multiplication and addition operations are replaced by cheap XNOR and bit-count operations of two quantized vectors with relatively lower energy requirement. For example, the system [147] showed that quantizing 16- or 32-bit weight to 4-10 bits could achieve energy reductions of up to 30 times without sacrificing algorithm performance. However, retraining is required after quantization and different layers may have different sensitivity to quantization resolution [148]. Additionally, computer hardware (e.g., CPU and GPU) does not support all the quantization resolution. For instance, Nvidia TensorRT only supports Integer quantization and Float16 quantization. Currently, popular DL libraries (e.g., Tensorflow and Pytorch) support only Integer quantization and Float16 quantization. A special kind of quantization model is Binarized Neural Networks (BNNs) [152]–[154] that perform the pure-logical computation by utilizing 1-bit weights and activations. Due to the fewer memory access, it significantly speeds up inference time and reduces energy consumption in resource-constrained devices. For example, Ding *et al.* [154] used 1-bit weights and activations to compute the distribution loss for regularizing the activation flow and develop a framework to



TABLE III  
CLASSIFICATION OF WITH-CHANGING-MODEL-ARCHITECTURE AND WITHOUT-CHANGING-MODEL-ARCHITECTURE DL MODEL OPTIMIZATION METHODS.

With-Changing-Model-Architecture	Without-Changing-Model-Architecture
Factorization [139]–[142]; Pruning [1], [3], [39], [40], [142], [143]; Generating new architecture [144]–[146]; Knowledge Distillation (KD) [144], [146];	Quantization [40], [142], [143], [147], [148];

systematically formulate the loss. In addition, the proposed approach is robust to select the hyperparameters (i.e., learning rate and optimizer) to train the state-of-the-art networks (i.e., AlexNet) on ImageNet dataset with high accuracy. The evaluation results demonstrate that the proposed BNN approach can realize efficient inference on resource-limited mobile devices with lower inference latency and energy consumption. Lin *et al.* [153] introduced two innovative optimization approaches to training BNNs at run-time. The approach used a linear combination of multiple binary weight bases to approximate full-precision weights and employ multiple binary activations to mitigate information loss. The experimental result shows the proposed approach significantly reduces the memory cost while maintaining high accuracy. Thus, the BNN approach is a potential option to enable the network to deploy on resource-constrained mobile devices.

#### 4) Combining DL Model Optimization With&Without Changing Model Architecture.

Although quantization technique can compress the network complexity by using limited numerical precision or clustering parameters, the compression ratio of quantization is usually less than the optimization methods with changing architecture such as pruning method. Since those two categories of methods are not conflicting with each other, some research [40], [142], [143] combined both of them for better performance. Deep compression [40] adopted both pruning and quantization to reduce the storage requirement of neural networks by 35 to 49 times without affecting their accuracy. It also achieved 3 to 4 times layer-wise speedup and 3 to 7 times better energy efficiency. This allowed fitting the model into on-chip SRAM cache rather than off-chip DRAM memory, which facilitated the use of complex neural networks in mobile applications with limited memory size and download bandwidth. CLIP-Q [143] combined network pruning and weight quantization in a single learning framework to perform pruning and quantization jointly. It obtained state-of-the-art compression rates of 51 times of AlexNet, 10 times of GoogLeNet, and 15 times of ResNet-50, respectively. Shi *et al.* [142] presented a compression approach based on the combination of low-rank matrix factorization and quantization training, to reduce complexity for neural network. For a three-layer long short-term memory network, the original model size could be reduced to 1% with negligible loss of accuracy. Moreover, factorization and KD can be further applied for reducing redundant computations and achieving faster inference.

In addition to the DL Model Optimization methods for reducing the model size, many existing works also focus on the generic optimization methods. These optimization methods

optimize the training algorithms, such as SGD, adaptive gradient methods, and distribution methods, to reduce the computational cost, energy consumption, and speed up the training process. They aim to solve the issues of gradient explosion, vanishing, and the undesirable spectrum during the training process meanwhile keeping high accuracy. For example, Le *et al.* [155] demonstrated that more sophisticated optimization methods (e.g., Limited memory BFGS (L-BFGS) and Conjugate Gradient (CG) with line search) could significantly simplify and speed up the process of training algorithms compared to the relatively simple optimization methods (e.e., SGD). MENNDL [156] developed a genetic algorithm to automate network selection on computational clusters through hyper-parameter optimization. It simplifies the process of applying deep networks to new applications and domains compared to the manual selection for the hyperparameter. To summary, such generic optimization methods reducing the training iterations result in computational cost, memory accesses and energy consumption reduction. Researchers could adopt these generic optimization methods combined with other optimization methods to enable the efficient training of the models on resource-constrained mobile devices.

#### C. Non-Model-Oriented Optimization Mechanisms

In addition to the model-oriented optimization mechanisms, there are some non-model-oriented ones that could further boost the performance of DL on mobile devices.

##### 1) Mobile DL Software Accelerators

It is a service to adopt pre-trained DL models on mobile devices. It decomposes a deep model across a mix of heterogeneous processors (e.g., CPU, GPU, DSP) to maximize energy-efficiency and execution time, with limited mobile resources, such as computation power and memory size. It has several advantages. Firstly, it removes the barriers preventing existing pre-trained DL models from being adopted by current mobile and wearable devices. Moreover, it allows complex DL models to run on commodity mobile devices with acceptable resource consumption levels and low latency. There are several software accelerators to optimize the DL for mobile devices.

DeepX [19] took an important step towards adopting DL models to mobile and wearable devices. This software accelerator significantly lowered the device resources (i.e., memory, computation, and energy) required by DL that currently acts as a severe bottleneck to mobile adoption. Specifically, a pair of resource control algorithms were proposed for the inference stage of DL. Firstly, DeepX decomposed monolithic deep model network architectures into various unit-blocks, which were then executed by heterogeneous local device processors (e.g., GPUs, CPUs); Secondly, it performed principled resource

TABLE IV  
THE SUMMARY OF SOFTWARE ACCELERATORS.

Accelerators	Design Principle	Advantages	Disadvantages
DeepX [19]	(1) Leverage a mix of heterogeneous processors (e.g., GPUs, LPUs); (2) Offer two resource control and optimization algorithms.	Decompose a deep model across available processors to maximize energy-efficiency and execution time.	DeepX is not applicable to those deep learning networks with temporal information.
CNNdroid [157]	Use parallel processing capabilities of mobile GPUs to accelerate deep CNN computations on current mobile devices.	(1) Compatible with CNN models trained by common desktop/server libraries; (2) Easy to configure and integrate into any Android app.	Do not consider the underlying hardware architecture to thoroughly make use of the hardware specifics.
RSTensorFlow [15]	Modify the kernels of TensorFlow operations to leverage the RenderScript heterogeneous computing framework on Android devices.	Utilize the power of available computation resources while running models trained with TensorFlow without requiring external tools.	Optimizing the energy costs of running on CPU vs GPU are potential future research directions.
DeepMon [18]	(1) Develop optimizations for processing convolutional layers; (2) Apply system-level optimizations to accelerate the matrix calculation.	Allow large DNNs to run on commodity mobile devices at a low latency.	Only extract features from video frames and need to be integrated with a complete cloud-enabled solution.
CADNN [36]	Leverage model compression to improve the DNN inference execution performance in the mobile environment.	Propose three major optimizations targeting modern mobile architectures.	Focus on the acceleration of existing DNN models designed for mobile applications.

TABLE V  
THE SUMMARY OF HARDWARE DESIGNS.

Accelerators	Categories	Advantages	Disadvantages
TrueNorth [125]	ASIC-based	Have associated software to aid in development.	The maximum matrix that one synaptic core can handle is $256 \times 256$ , and it supports only a simplified leaky-integrate-and-fire neuron model.
VPU [126]	ASIC-based	It can be deployed onboard a drone and used to perform inference on drone imagery, without significantly taxing system resources.	The execution time per inference using one chip is 4× slower compared to a reference CPU/GPU implementation.
EIE [127]	ASIC-based	Operating directly on compressed networks enables the large neural network models to fit in on-chip SRAM, which results in 120× better energy savings compared to accessing from external DRAM.	Not scalable due to all-to-all processing element broadcasts and a BW link of one element per cycle.
DianNao [128]	ASIC-based	The accelerator achieves high throughput in a small area, power and energy footprint.	The scalability and efficiency are severely limited by the bandwidth constraints of the memory system.
Zhang et al. [123]	FPGA-based	Successfully deploy the deeper VGGNet into an embedded FPGA platform, with several optimization techniques.	1) This approach is limited to VGG models. 2) Not flexible enough to match each layer's distinct features and result in underutilization of hardware resources.
FPGA15 [124]	FPGA-based	Under the constraints of computation resource and memory bandwidth, it explores all possible solutions in the design space using a roofline model.	The accelerator design was only applied to several CONV layers rather than the full CNN.

scaling that adjusted the architecture of deep models to reduce the overhead each unit-blocks introduces. Experiments showed that across various model and processor combinations, the mean energy benefit of DeepX was 7.12 times (Snapdragon) and 26.7 times (Tegra) under the time requirement.

Another accelerator CNNdroid [157] was a GPU-accelerated library for execution of trained deep CNNs on Android-based mobile devices. It was compatible with CNN models trained by common libraries, such as Caffe, Torch and Theano. Empirical evaluations showed that CNNdroid achieved up to 60 times faster speed and 130 times energy saving on current mobile devices. Moreover, the source code, documentation and sample projects were publicly available (<https://github.com/ENCP/CNNdroid>).

Although DeepX and CNNdroid could use computing resources on mobile devices to accelerate DL models, there are some limitations. These frameworks require some dedicated hardware (e.g., GPU) which are not available for some mobile devices. In addition, they have not been integrated well with popular DL frameworks, and thus need external tools to run models trained with existing libraries such as TensorFlow. Different from above two methods, RenderScript [15] was proposed as an extension to TensorFlow. By integrating the

acceleration framework tightly into TensorFlow, researchers could easily take advantage of the heterogeneous computing resources (CPU and GPU) on mobile devices without the need for other external tools. DL models could run three times faster by using RenderScript.

Another approach, DeepMon [18], allowed large DNNs to run on commercial mobile devices at a low latency. Prior work, such as DeepX, has focused on smaller DNNs, and more powerful non-commodity mobile devices such as the Tegra K1. DeepMon could complement DeepX by supporting various mobile GPUs (e.g., Adreno 420, Adreno 430, and Mali T 880) on the market. Moreover, DeepMon proposed to reduce the latency of convolutional layers. In particular, by devising a suite of optimization techniques, including the convolutional layer caching, decomposition, and matrix multiplication optimizations, DeepMon is two times faster than DeepX.

Furthermore, Wei *et al.* [36] found that most of the existing accelerators, such as DeepMon and DeepX did not explore possible optimization opportunities in computation and memory footprint reductions offered by model compression, including weight pruning and weight quantization. Therefore, a significant performance gap still remained between the peak performance

that could be potentially offered by mobile devices and what the existing systems actually achieved. Based on that, Wei *et al.* [36] proposed CADNN, a programming framework, to efficiently execute DNN inference on mobile devices with the compression (i.e., pruning) and a set of architecture-aware optimizations (e.g., computation pattern transformation, redundant memory load elimination, smart selection of memory, and computation optimization parameters, etc.). Based on evaluation studies, CADNN is about 8.8 times and 6.4 times faster than TensorFlow Lite and TVM respectively, two popular and highly optimized dense DNN execution frameworks.

As shown in Table IV, software accelerators can benefit from the deployments of pre-trained DL models on mobile devices based on various design principles. If the researchers would like to deploy their DL models on some resource-constraint mobile devices (e.g., IoT) and improve the energy efficiency and lower the memory and computation requirements of the DL, some common accelerator with resource control strategies and optimization approaches (e.g., DeepX, CADNN) might be considered. If the pre-trained DL model is developed with some specific libraries (e.g., TensorFlow) or is deployed on some specific devices (e.g., Android devices), some accelerators such as RendorScript and CNNdroid would benefit the researchers a lot, since they aimed at specific platforms or frameworks. If the researchers would like to adopt some large-scale DNN models, DeepMon is a better choice since it helps adopt large DNN models with low latency on various mobile devices in the market.

## 2) Mobile Hardware Designs for DL

Advances in DL are mainly due to powerful computer systems. Compared to such systems, mobile devices have limited computational power and battery lives, making them not suitable for DL deployment. In existing research, the most common computing processors used for DL on mobile devices are CPUs, GPUs, Field Programmable Gate Array (FPGA), and Application Specific Integrated Circuit (ASIC). We compare the power efficiency, time/cost budget, and compatibility of these processors in Figure 6. CPU and mobile GPU are widely used in the embedded platform and mobile devices due to their high availability and compatibility. They can be supported by multiple DL frameworks and various programming languages. Furthermore, some GPU-based execution frameworks [158], [159] can be utilized to provide low latency, high throughput, and efficient on-chip resource utilization. It is a good choice for the researchers to start with CPUs and GPUs to explore the feasibility of deploying DL models on mobile devices. However, as shown in Figure 6, CPUs and GPUs have the lowest power efficiency, which is not really mobile-friendly. To enable DL on power-constrained mobile devices, ASIC is the most power-efficient architecture suitable for processing DL on mobile devices. Compared to ASIC, FPGA is a balanced choice considering the power efficiency, compatibility, and developing budget. FPGA is also more flexible in development since it can be (re)programmed by firmware and supported by open-sourced libraries (e.g., Caffe).

Recent studies have shown that highly efficient DL training and inference with low energy consumption could be achieved using processors with specially designed architecture.

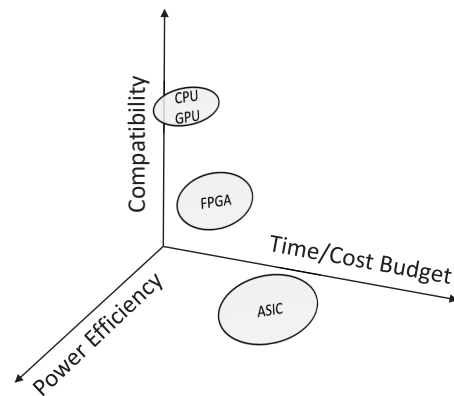


Fig. 6. Metrics of different computation components. There are three metrics we distinguish on various hardware devices, including Power Efficiency, Compatibility and Time/Cost Budget.

ASIC [125]–[128], [160], [161] and FPGA [123], [124], [162]) have offered a good starting point for developing DL processors for mobile devices. We summarized these works in Table V. For instance, TrueNorth [125] proposed a real-time neurosynaptic processor for mobile devices based on ASIC, which achieved an extremely low typical power consumption of 65mW. Besides, TrueNorth was fully configurable in terms of connectivity and neural parameters to allow custom configurations for a wide range of applications. The authors also introduced an asynchronous-synchronous tool that may facilitate general ASIC design. An ASIC-based VPU [126] has been used to enable low-power DL computation. This paper showed that VPU, as a highly-parallel vector co-processor with low power consumption, could reduce the thermal design power (TDP) up to 8 times compared to their CPU and GPU implementations. Han *et al.* [127] proposed an optimized energy-efficient engine, EIE, to operate on compressed deep neural networks. By leveraging sparsity in both the activations and weights, and integrating weight sharing and quantization techniques, this engine outperformed CPU, GPU, and mobile GPU by factors of 189 $\times$ , 13 $\times$  and 307 $\times$ , and consumed 24,000 $\times$ , 3,400 $\times$  and 2,700 $\times$  less energy than CPU, GPU and mobile GPU, respectively. Chen *et al.* [128] proposed an ASIC accelerator, DianNao, for the fast and low-energy execution of the inference of large CNNs and DNNs in a small form factor. This work had a special emphasis on the impact of memory on accelerator design, performance, and energy. This accelerator had high throughput, capable of performing 452GOPS in a small footprint of 3.02mm and 485mW. It achieved a speedup of 117.87 $\times$  and an energy reduction of 21.08 $\times$  over a 128-bit 2GHz SIMD core with a normal cache hierarchy.

Jiantao *et al.* [123] went with an FPGA-based approach for accelerating the large-scale image classification using a CNN. A state-of-the-art CNN, VGG16-SVD, was implemented on an FPGA platform. Experimental results showed that GPU consumed 26 times more power consumption than the proposed FPGA hardware. Zhang *et al.* [124] proposed a roofline-model-based FPGA accelerator. In this method they first optimized CNN's computation and memory access. Then, they modeled all possible designs in roofline model to find the best design for each layer. They realized an implementation on Xilinx VC707

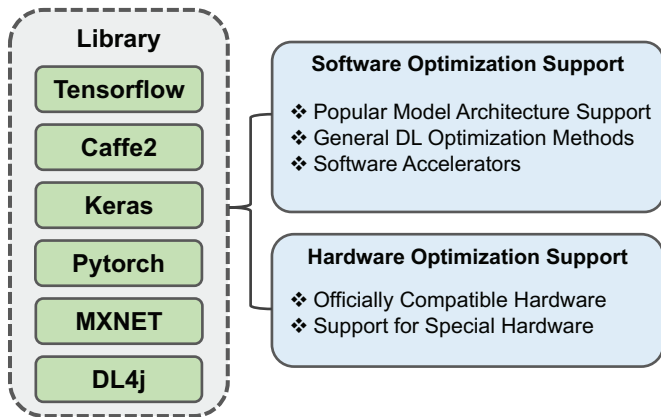


Fig. 7. Library Optimization Support. We summarize software libraries regarding their support for software optimization and hardware optimization respectively. In software optimization, we discuss popular model architectures and general DL optimization methods as well as software accelerators. In hardware optimization, we summarize the popular libraries for various types of hardware. Besides, we also consider special designs such as ASIC and FPGA.

board and achieved a 4.8x speedup over CPU.

Based on the proposed DL optimization pipeline for mobile applications, when deploying DL models on common mobile and embedded devices, researchers could apply different optimization mechanisms to improve the performance of mobile applications. The model-oriented optimization approaches, such as pruning, could be leveraged to reduce the memory and storage requirements and speed up the inference time. The non-model-oriented optimization mechanisms, such as software accelerators, could further improve the efficiency of executing DL operations without modifying models. While the improvement of software accelerators may be limited by various software environments such as operating systems and drivers, hardware accelerators based on special computer architecture (e.g., ASIC or FPGA) could be designed to improve the efficiency of DL operations on mobile devices significantly when time and money budget is sufficient.

#### IV. SOFTWARE LIBRARIES

This section summarizes software libraries regarding their support for software optimization and hardware optimization respectively, as shown in Figure 7. For software optimization, we discuss popular model architectures and general DL optimization approaches (e.g., quantization and pruning) as well as software accelerators. For hardware optimization, we summarize the popular libraries for various types of hardware such as mobile and powerful embedded devices and microcontroller (i.e., IoT device). Besides, we also consider the related designs such as ASIC and FPGA.

**Software Optimization Support.** Popular deep learning libraries, such as Tensorflow, Pytorch, Keras, Caffe2, Deeplearning4j (DL4j) support most of the state-of-art model architectures in various application domains as shown in Table I, II. In particular, some pretrained models for image, text, video applications, such as Mobilenet, Resnet, Inception, Convnet, and VGG are available [163]–[165] and finetuning in user-specific applications for downloading.

Besides supporting model architectures, those libraries also provide general optimization procedures. As shown in Table VI, Tensorflow, Pytorch, and Keras currently support common optimization algorithms via both quantization and pruning [166]. In particular, Tensorflow officially provides full integer quantization, float16 quantization and magnitude-based weight pruning. For example, it can reduce the model size fourfold and increase the computational speed threefold on CPU, Edge TPU, and Microcontrollers via full integer quantization [167]. Moreover, its supported float16 quantization can reduce the model size twofold on CPU with potential GPU acceleration [167]. In addition, Tensorflow achieves model sparsity by supporting magnitude-based weight pruning [168]. Since Keras can adopt Tensorflow as backend, it can carry those optimization functionalities provided by Tensorflow. In contrast, PyTorch supports only integer quantization allowing a four times reduction in the model size and memory bandwidth requirements [169]. And it provides pruning capability via the third part package (e.g., Intel AI Lab proposes Distiller [170], an open-source Python package for neural network compression research). Caffe2 and MXNet currently support only integer quantization via the official API [171] and Apache/MXNet toolkit, respectively. Optimization techniques such as factorization and knowledge distillation can also be adopted when using those libraries. In addition, some existing software accelerators (e.g., DeepX, RenderScript, DeepMon, CNNdroid, etc.) can also be adopted to further enhance computational performance by leveraging various mobile hardware components such as GPU and LPU. Table VI summarizes the accelerators and related libraries.

**Hardware Optimization Support.** Different hardware can support various deep learning frameworks and libraries which made it convenient for researchers to deploy DL on their devices. Mobile phones and powerful embedded devices are the most common authorized and supported mobile hardware to deploy DL. Specifically, Android and iOS mobile phones are supported by most of the popular libraries as shown in Table VII. In addition, TensorFlow, Caffe2, PyTorch, MXNet, etc are available in the embedded device such as Raspberry Pi, NVIDIA Tegra and Nvidia Jetson TX2. Some libraries have been adopted for resource-constraint mobile devices such as microcontrollers (IoT). For example, TensorFlow Lite, a light-weight version of TensorFlow (<https://www.tensorflow.org/lite/microcontrollers>), supports the processors based on the Arm Cortex-M Series and ESP32 architecture (e.g., Arduino Nano 33 BLE Sense, SparkFun Edge, etc.). Caffe2 facilitates certain microcontroller (e.g., Arm Cortex-M processor cores) by using the CMSIS-NN optimized libraries. Furthermore, researchers also design customized hardware (e.g., ASIC and FPGA) to implement libraries such as TensorFlow, Python, and MXNet.

**Model Conversion Tools.** There are many available libraries and each of them has its advantages and use cases. To make the best use of different features provided by these libraries and allow researchers to take full advantage of the libraries based on its situation, model conversion tools are developed. Then, users can smoothly convert DL models from one library to another library without reprogramming. For example, MMDnn [172] is

TABLE VI  
SOFTWARE OPTIMIZATION SUPPORT.

	General DL Optimization Methods				Software Accelerators
	pruning	quantization	factorization	KD	
<b>TensorFlow</b>	weight pruning*	Integer quantization*, float16 quantization*	✓	✓	DeepX, RenderScript, DeepMon
<b>Caffe2</b>	N/A	Integer quantization*	✓	✓	DeepX, CNNdroid
<b>Keras</b>	weight pruning*	Integer quantization*, float16 quantization*	✓	✓	DeepX
<b>Pytorch</b>	weight pruning	Integer quantization*	✓	✓	DeepX, CNNdroid
<b>MXNET</b>	N/A	Integer quantization*	✓	N/A	N/A
<b>DL4j</b>	N/A	N/A	N/A	N/A	DeepX

TABLE VII  
HARDWARE OPTIMIZATION SUPPORT.

	Official Supported Hardware		Special Hardware Design
	Mobile & powerful embedded device	Microcontroller support (IoT)	
<b>TensorFlow</b>	standard TensorFlow Lite (smartphone and Raspberry Pi)	TensorFlow Lite for Microcontrollers (processors based on the Arm Cortex-M Series and ESP32 architecture, e.g., Arduino Nano 33 BLE Sense, SparkFun Edge)	ASIC support
<b>Caffe2</b>	iOS, Android, Tegra, and Raspberry Pi platforms	CMSIS-NN optimized libraries for Arm Cortex-M processor cores	N/A
<b>Keras</b>	mobile and embedded devices	N/A	N/A
<b>PyTorch</b>	iOS, Android, and Raspberry Pi	N/A	ASIC support
<b>MXNet</b>	iOS, Android, Raspberry Pi, and NVIDIA Jetson TX	N/A	ASIC support
<b>DL4j</b>	Android, iOS	N/A	N/A

a comprehensive and cross-framework tool to convert, visualize, and diagnose deep learning models. It is a universal converter to convert DL models among a number of existing frameworks, making a trained model of one framework be easily deployed on another framework for inference and prediction. Open Neural Network Exchange (ONNX) [173] is an open standard format for representing machine learning models. It is facilitated by a community of partners for implementations in many frameworks, such as PyTorch, MXNet, Caffe2, TensorFlow, etc.

## V. FEDERATED LEARNING ON MOBILE DEVICES

Training deep learning network on mobile device is another challenging issue. Federated learning (FL) is a way of enabling deep learning training by leveraging many client devices (e.g., mobile devices and IoT devices) in a network. This solves the issue that a single resource-constraint device may never support training a deep learning model due to its limited hardware resources. Particularly, all mobile devices need to devote their storage and computing resources for their data training. However, this requirement is not always satisfied. When FL deploys identical neural network models to heterogeneous devices, the ones with weak computational capacities may significantly delay the synchronous parameter aggregation [174]. Besides, a large amount of memory and power are needed for the training of advanced DL models and

the storage of model parameters [175]. So, many aspects still need to be carefully considered including memory limitation, energy budget, communication, synchronization, resource distribution, privacy, etc. In this work, we focus on the on-device optimization in FL domain and mainly discuss how to save memory usage and energy consumption using the proposed DL optimization pipeline. Other aspects are out of the scope of discussion of this paper.

### A. Model-Oriented Optimization Mechanisms in Federated Learning

#### 1) DL Model Selection for FL

Since the FL runs on multiple heterogeneous devices, which vary in computing ability and memory size, so it is critical to first select the suitable models that are smoothly executed on various devices. Specifically, due to the limited memory and computing resources of some resource-constraint FL clients, the large size of CNN models (i.e., VGG and ResNet) are hardly adopted on them. Researchers have been exploring lightweight and sparse deep learning architecture for more general applicability. For instance, a novel tri-layer FL scheme [176] is proposed to consume less amount of resources to accomplish a target convergence. A variant of the Long Short-Term Memory (LSTM) recurrent neural network called the Coupled Input and Forget Gate (CIFG) [177] is adopted for mobile keyboard predication with federated learning. The CIFG



architecture is advantageous for the mobile device environment because the number of computations and the parameter set size are reduced with no impact on model performance. Moreover, a personalized federated learning framework (e.g., three layer CNN) [178] is advocated to cope with the heterogeneity issues in FL environments.

## 2) DL Model Optimization for FL

Since some FL clients may not have enough memory size and energy budget to meet the system requirements for training, the researchers have been exploring how to reduce the memory cost and save the energy cost in FL training and prolong battery-power life duration. To address these issues, some researchers have adopted the *With-Changing-Model-Architecture* model optimization approaches. For example, PruneFL [179] proposes an adaptive and distributed parameter pruning approach, which adapts the model size during FL to reduce both communication and computation overhead and minimize the overall training time while maintaining a similar accuracy as the original model. To reduce the memory requirement for training, Helios [174] devises soft-training method to dynamically compress the original training model into the expected size through a rotating neuron training approach. FedPARL [176] reduces the model size by performing sample-based pruning, which is demonstrated useful for resource-constraint IoT devices. Moreover, *Without-Changing-Model-Architecture* model optimization approaches have also been adopted. For instance, Wu *et al.* [180] propose a quantization framework that constrains all layers to low-bit width integers in both training and inference. This approach can reduce the energy about 5 times compared to 32-bit float point operations. In addition, dynamic computation approaches are explored to optimize the model during training and inference phase. The basic idea of those approach is to only activate the partial neural network to process the input. For example, dynamic channel pruning techniques [181], [182] are adopted to identify channels in the deep neural network that are considered as unimportant and skip their computations. Some researchers [182], [183] apply dynamic layer skipping techniques on ResNet models and RNN models with skip connections.

## B. Non-Model-Oriented Optimization Mechanisms in Federated Learning

### 1) Mobile DL Software Accelerators for FL

Several possible solutions should be considered to facilitate on-device DL learning for FL systems. One direction is to improve DL hardware usage on mobile devices via the mobile DL software accelerators [15], [18], [19], [36], [157] as we discussed in Section III. Since these DL software accelerators can decompose a deep model across a mix of heterogeneous processors (e.g., CPU, GPU, DSP) to maximize energy-efficiency and execution time, with limited mobile resources (e.g., computation power and memory size), it removes the barriers preventing existing pre-trained DL models from being adopted by current mobile devices in FL systems. For example, DeepX [19] significantly lowers the device resources (i.e., memory, computation, energy) required by

deep learning that currently acts as a severe bottleneck to mobile adoption through a pair of resource control algorithms, designed for inference stage of deep learning. Experiments demonstrate its superior performance in terms of low execution time and energy consumption in DL running, compared to cloud offloading-based approaches. This research is promising and likely to facilitate the development of sensor processing and mobile DL inference, which would enable on-device FL implementation at scale.

### 2) Mobile Hardware Designs for FL

Another direction is mobile hardware designs for FL. In recent years, few hardware designs are emerging to solve some specific problems in FL area. For example, the complicated operations and large operands of privacy-preserving mechanisms (e.g., homomorphic encryption) impose significant overhead on federated learning. Maintaining accuracy and security more efficiently has been a key problem of federated learning. Yang *et al.* [184] propose a hardware solution to accelerate the training phase in federated learning by designing an FPGA-based homomorphic encryption framework. Experiments show that the proposed accelerator achieves a near-optimal execution clock cycle, with a better digital signal processing efficiency and reduces the encryption time by up to 71% than existing designs. Another work [122] proposes HAFLO, a GPU-based solution to improve the performance of FL by reducing the significant computational overhead of homomorphic encryption. The core idea of HAFLO is to summarize a set of performance-critical homomorphic operators used by federated logistic regression and accelerate the execution of these operators through joint optimization of storage, IO, and computation. Experimental results show that the proposed method achieves a  $49.9\times$  speedup acceleration on a popular FL framework. Furthermore, as we discussed in Section III, some studies have demonstrated that highly efficient DL training and inference with low energy consumption could be achieved using specially designed processors. Some ASICs [125]–[128], [160], [161] and FPGAs [123], [124], [162] have offered a good starting point for developing DL processors to mitigate the constraints of hardware, memory, and power resources of mobile devices. These works are promising and likely to further enable on-device FL implementation at scale. Exploring how these work execute in FL systems might help tailor these solutions to solve existing issues in FL.

## C. Cooperative DNN Inference in Federated Learning

A fast enough DNN inference with acceptable resource consumption is also an essential consideration in federated learning due to the latency requirements of many applications and the resource constraints of a single device. The existing researchers have been taking advantage of the aggregated computational power of a cluster of mobile devices, coordinating them for faster DNN inference. For example, DeepThings [110] proposes a Fused Tile Partitioning (FTP) method for dividing convolutional layers into independently distributable tasks and parallelizing them on multiple devices. It provides scalable CNN inference speedups of  $1.7\text{--}3.5\times$  on 2–6 edge devices.

MoDNN [185] partitions trained DNN models on several mobile devices, assigning more workload to the more powerful devices for DNN inference acceleration. It can accelerate the DNN computation by 2.17-4.28 $\times$  when the number of devices increases from 2 to 4. Hadidi *et al.* [186] adopt both data parallelism and model parallelism to accelerate DNN inference using several robots. CoEdge [187] dynamically partitions the DNN inference workload based on the computing capabilities of devices and network conditions. It achieves up to 4.49-7.21 $\times$  latency speedup over the local approach and 25.5%-66.9% energy reduction for four widely-adopted CNN models. These works demonstrate that cooperative DNN inference in FL could significantly speed up model inference and reduce resource consumption on every mobile device.

## VI. POTENTIAL RESEARCH OPPORTUNITIES IN MOBILE APPLICATIONS

In this section, we provide some potential research opportunities for mobile applications in the areas of database and model, inference time, power consumption, hardware design, and optimization trade-off on different hardware.

**Database and Model.** As discussed in Section II, the existing DL work in non-vision fields usually adopts various lightweight DL models. And there is also a lack of the baseline models in non-vision field based on which to compare with each other's performance. Therefore, it is desirable to select general DL models as basis for performance evaluation. In addition, given the fact that various types of mobile devices could be used, a public repository of benchmark performance in various hardware could be of great benefit to the research community. Moreover, the non-vision fields still lack comprehensive, diverse, and high-quality datasets for DL model training and evaluating. There are two reasons for such inadequacy. First, service providers prefer to keep their data confidential and rarely publish their datasets. Second, due to limitations of hardware resource and network conditions, data collected from mobile devices usually come with loss, redundancy, mislabeling, and class imbalance. To promote adopting DL on mobile devices for various application domains in both academic and industry communities, researchers and companies are encouraged to collect and publish more high-quality datasets.

**Inference Time.** As shown in section II, the existing work using DL in non-vision fields (e.g., activity recognition and mobile security) focuses much more in accuracy than inference time. However, those applications require certain speed for inference in reality. For example, human daily activities recognition requires the inferencing time (about 0.3 to 3s) to achieve real-time nature. In security area, a detecting time with a delay of 2 seconds is usually considered as acceptable for detecting attacks [188], [189]. Therefore, researchers should also evaluate inference time for those non-vision applications in the future. Besides, since different mobile hardware has different computational power, researchers may consider an alternative way (e.g., FLOPs widely used in vision field) to fairly evaluate the potential inference speed of their systems. In addition, improving the inference speed is also critical for mobile devices. Particularly, developing and adopting various

compression techniques to reduce the inference time on mobile devices still needs exploration.

**Power Consumption.** With the aid of highly effective computation components (e.g., GPU, ASIC, etc.) in mobile devices, deep learning models can be easily deployed with decent performances [190]. However, since deep learning is computationally intensive, energy and capability constraints should be considered when we deploy DL models on mobile devices because of their limited battery capacities. As a result, reducing the energy consumption of mobile devices with respect to data collection (e.g., from motion sensors and GPS) in motion-sensor-based applications might be an open research area. Another interesting research direction is to optimize energy consumption when deploying DL models on mobile devices. Although designing some specific hardware chips (e.g., FPGAs and ASICs) can achieve this purpose in section III, the generality of these special hardware design, as well as their interoperability and compatibility with existing hardware platforms, remain as challenges.

**Hardware Design.** DL on mobile devices is facilitated by mobile GPUs but is usually limited within several layers, as shown in section II, due to the constraint of computation and memory resources. Thus, the performance of some mobile applications is usually not satisfactory. In future research, in order to execute more complex models with better performance, mobile GPU with more CUDA cores and large graphic memories need to be developed. In addition to powerful hardware components, some algorithms, such as the Toeplitz matrix, and Winograd and Strassen algorithms, can be explored to improve the computing performance by taking advantage of the low-latency temporary storage architecture of GPU [191].

Compared with powerful mobile devices, such as smartphones, the resource bottleneck in microcontrollers is more serious. Specifically, smartphones can now equip with several gigabytes of RAM, but microcontrollers, such as the ARM Cortex series, are limited to just hundreds of kilobytes. Some techniques, such as binary deep architectures, have the potential to fill this gap [192]. Such architectures could not only build extremely small models but also remove the requirements for expensive multiplication operations. Besides, the limited on-chip memory capacity often causes massive off-chip memory access and leads to very high energy consumption. Some frameworks such as retention-aware neural acceleration (RANA) [193], can be explored to save the energy consumption of microcontrollers.

Furthermore, FPGAs and ASICs are also promising for DL. However, they currently lack adequate software supports to fully achieve their potential in DL. Since FPGAs and ASICs are built with spatial architectures with low-energy on-chip memory, in future research, reusable dataflow algorithms can provide solutions to reduce data movements. Moreover, a promising field called neuromorphic computing, enables information processing at meager energy cost on electronic devices via emulating the electrical behaviors of biological neural networks in human brains. Based on this technique, some companies are trying to develop ASIC chips, such as TrueNorth from IBM. Since neuromorphic architectures are more suitable for brain-like computations and achieve decent

power efficiency, as shown in section III, they would be an attractive topic for future research.

**Optimization Trade-off on Different Hardware.** In DL, larger and deeper models usually produce higher accuracy. However, they may also lead to larger inference time and more energy consumption. Building DL models on mobile devices with an accurate scheme, fast inference time, and without mass power consumption and huge memory usage should be a good research direction. Besides, even with various existing optimization approaches for DL on mobile platforms, it is still challenging to adapt DL models to various device hardware, while conventional optimization methods always ignore different hardware architectures and optimize all the DL models in a uniform way. Recent work [194] has highlighted that pruning and quantization methodologies relied on formulations are hardware-unaware, and they do not necessarily result in optimal configurations in terms of hardware efficiency. In order to solve the above issues, some new techniques, such as design automation technique [126], hardware-aware modeling methodologies [125], and stochastic computing [160] could be explored to achieve optional optimization on different hardware. The insights of these design policies will inspire future research in optimization with various mobile hardware to achieve efficient deep learning computing.

## VII. CONCLUSION

This paper provides a comprehensive review of the recent advancements of deep learning on mobile devices. Applications in various areas are summarized and demonstrated at the intersection of deep learning and mobile computing. The challenges of bringing DL on mobile devices are discussed, and a thorough DL optimization pipeline is introduced. Compared with existing work, optimization approaches including general and optional optimization approaches are studied in detail. Moreover, the popular DL libraries are summarized with respect to software and hardware optimization support. These resources may serve as references and recommendations for researchers when they try to find appropriate optimization approaches and suitable DL libraries to deploy DL on mobile devices. Furthermore, the potential research opportunities of DL on mobile devices are also discussed with respect to database and model, inference time, power consumption, as well as hardware design, and optimization for trade-off on different hardware. With the fast advances of deep learning algorithms and sustained emergence of powerful mobile hardware, deep learning on mobile devices would remain as a hot topic, and more challenging and interesting research directions would spring up in the future.

## ACKNOWLEDGMENT

This work was partially supported by the National Science Foundation Grants CCF2000480, CCF2028858, CCF2028894, CCF2028873, CCF2028876, CCF1909963, CNS1815908, and CNS1717356.

## REFERENCES

- [1] Y. He, J. Lin, Z. Liu, H. Wang, L.-J. Li, and S. Han, "Amc: Automl for model compression and acceleration on mobile devices," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 784–800.
- [2] K. Yang, T. Xing, Y. Liu, Z. Li, X. Gong, X. Chen, and D. Fang, "cdeeparch: A compact deep neural network architecture for mobile sensing," *IEEE/ACM Transactions on Networking*, vol. 27, no. 5, pp. 2043–2055, 2019.
- [3] D. Li, X. Wang, and D. Kong, "Deeprebirth: Accelerating deep neural network execution on mobile devices," in *Thirty-second AAAI conference on artificial intelligence*, 2018.
- [4] K. Yanai, R. Tanno, and K. Okamoto, "Efficient mobile implementation of a cnn-based object recognition system," in *Proceedings of the 24th ACM international conference on Multimedia*, 2016, pp. 362–366.
- [5] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, "Mnasnet: Platform-aware neural architecture search for mobile," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 2820–2828.
- [6] R. J. Wang, X. Li, and C. X. Ling, "Peleee: A real-time object detection system on mobile devices," in *Advances in Neural Information Processing Systems*, 2018, pp. 1963–1972.
- [7] Y. He, Z. Pan, L. Li, Y. Shan, D. Cao, and L. Chen, "Real-time vehicle detection from short-range aerial image with compressed mobilenet," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 8339–8345.
- [8] H. Ai, W. Xia, and Q. Zhang, "Speaker recognition based on lightweight neural network for smart home solutions," in *International Symposium on Cyberspace Safety and Security*. Springer, 2019, pp. 421–431.
- [9] X. Lei, A. Senior, A. Gruenstein, and J. Sorensen, "Accurate and compact large vocabulary speech recognition on mobile devices," 2013.
- [10] I. McGraw, R. Prabhavalkar, R. Alvarez, M. G. Arenas, K. Rao, D. Rybach, O. Alsharif, H. Sak, A. Gruenstein, F. Beaufays *et al.*, "Personalized speech recognition on mobile devices," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 5955–5959.
- [11] Y. He, T. N. Sainath, R. Prabhavalkar, I. McGraw, R. Alvarez, D. Zhao, D. Rybach, A. Kannan, Y. Wu, R. Pang *et al.*, "Streaming end-to-end speech recognition for mobile devices," in *ICASSP 2019-2019 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2019, pp. 6381–6385.
- [12] F. Cruciani, I. Cleland, C. Nugent, P. McCullagh, K. Synnes, and J. Hallberg, "Automatic annotation for human activity recognition in free living using a smartphone," *Sensors*, vol. 18, no. 7, p. 2203, 2018.
- [13] T. Zebin, P. J. Scully, N. Peek, A. J. Casson, and K. B. Ozanyan, "Design and implementation of a convolutional neural network on an edge computing smartphone for human activity recognition," *IEEE Access*, vol. 7, pp. 133 509–133 520, 2019.
- [14] N. Mairitha, T. Mairitha, and S. Inoue, "On-device deep learning inference for efficient activity data collection," *Sensors*, vol. 19, no. 15, p. 3434, 2019.
- [15] M. Alzantot, Y. Wang, Z. Ren, and M. B. Srivastava, "Rstensorflow: Gpu enabled tensorflow for deep learning on commodity android devices," in *Proceedings of the 1st International Workshop on Deep Learning for Mobile Systems and Applications*, 2017, pp. 7–12.
- [16] S. Tuli, N. Basumatary, and R. Buyya, "Edgelens: Deep learning based object detection in integrated iot, fog and cloud computing environments," *arXiv preprint arXiv:1906.11056*, 2019.
- [17] M. Mohammadi, A. Al-Fuqaha, M. Guizani, and J.-S. Oh, "Semisupervised deep reinforcement learning in support of iot and smart city services," *IEEE Internet of Things Journal*, vol. 5, no. 2, pp. 624–635, 2017.
- [18] L. N. Huynh, Y. Lee, and R. K. Balan, "Deepmon: Mobile gpu-based deep learning framework for continuous vision applications," in *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*, 2017, pp. 82–95.
- [19] N. D. Lane, S. Bhattacharya, P. Georgiev, C. Forlivesi, L. Jiao, L. Qendro, and F. Kawsar, "Deepx: A software accelerator for low-power deep learning inference on mobile devices," in *2016 15th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN)*. IEEE, 2016, pp. 1–12.
- [20] A. Antoniou and P. Angelov, "A general purpose intelligent surveillance system for mobile devices using deep learning," in *2016 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2016, pp. 2879–2886.

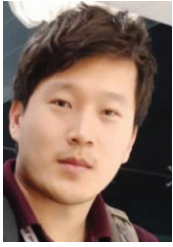
- [21] A. Xiao, R. Chen, D. Li, Y. Chen, and D. Wu, "An indoor positioning system based on static objects in large indoor scenes by using smartphone cameras," *Sensors*, vol. 18, no. 7, p. 2229, 2018.
- [22] N. Parikh, I. Shah, and S. Vahora, "Android smartphone based visual object recognition for visually impaired using deep learning," in *2018 International Conference on Communication and Signal Processing (ICCSP)*. IEEE, 2018, pp. 0420–0425.
- [23] V. Mandal, L. Uong, and Y. Adu-Gyamfi, "Automated road crack detection using deep convolutional neural networks," in *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018, pp. 5212–5215.
- [24] L. Tobías, A. Ducournau, F. Rousseau, G. Mercier, and R. Fablet, "Convolutional neural networks for object recognition on mobile devices: A case study," in *2016 23rd International Conference on Pattern Recognition (ICPR)*. IEEE, 2016, pp. 3530–3535.
- [25] S. Bazrafkan, T. Nedelcu, P. Filipczuk, and P. Corcoran, "Deep learning for facial expression recognition: A step closer to a smartphone that knows your moods," in *2017 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 2017, pp. 217–220.
- [26] I. Song, H.-J. Kim, and P. B. Jeon, "Deep learning for real-time robust facial expression recognition on a smartphone," in *2014 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 2014, pp. 564–567.
- [27] K. Z. Haider, K. R. Malik, S. Khalid, T. Nawaz, and S. Jabbar, "Deepgender: real-time gender classification using deep learning for smartphones," *Journal of Real-Time Image Processing*, vol. 16, no. 1, pp. 15–29, 2019.
- [28] B. A. Ashqar and S. S. Abu-Naser, "Image-based tomato leaves diseases detection using deep learning," 2018.
- [29] Y. Gu, Q. Wang, and S. Kamijo, "Intelligent driving data recorder in smartphone using deep neural network-based speedometer and scene understanding," *IEEE Sensors Journal*, vol. 19, no. 1, pp. 287–296, 2018.
- [30] H. Maeda, Y. Sekimoto, and T. Seto, "Lightweight road manager: smartphone-based automatic determination of road damage status by deep neural network," in *Proceedings of the 5th ACM SIGSPATIAL International Workshop on Mobile Geographic Information Systems*, 2016, pp. 37–45.
- [31] X. J. Tang, Y. H. Tay, N. A. Siam, and S. C. Lim, "Mywood-id: Automated macroscopic wood identification system using smartphone and macro-lens," in *Proceedings of the 2018 International Conference on Computational Intelligence and Intelligent Systems*, 2018, pp. 37–43.
- [32] A. Rattani and R. Derakhshani, "On fine-tuning convolutional neural networks for smartphone based ocular recognition," in *2017 IEEE international joint conference on biometrics (IJCB)*. IEEE, 2017, pp. 762–767.
- [33] H. Maeda, Y. Sekimoto, T. Seto, T. Kashiya, and H. Omata, "Road damage detection and classification using deep neural networks with smartphone images," *Computer-Aided Civil and Infrastructure Engineering*, vol. 33, no. 12, pp. 1127–1141, 2018.
- [34] N. Hnoohom and S. Yuenyong, "Thai fast food image classification using deep learning," in *2018 International ECTI northern section conference on electrical, electronics, computer and telecommunications engineering (ECTI-NCON)*. IEEE, 2018, pp. 116–119.
- [35] G. G. De Angelo, A. G. Pacheco, and R. A. Krohling, "Skin lesion segmentation using deep learning for images acquired from smartphones," in *2019 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2019, pp. 1–8.
- [36] W. Niu, X. Ma, Y. Wang, and B. Ren, "26ms inference time for resnet-50: Towards real-time execution of all dnns on smartphone," *arXiv preprint arXiv:1905.00571*, 2019.
- [37] A. Alfarrarjeh, D. Trivedi, S. H. Kim, and C. Shahabi, "A deep learning approach for road damage detection from smartphone images," in *2018 IEEE International Conference on Big Data (Big Data)*. IEEE, 2018, pp. 5201–5204.
- [38] G. Mittal, K. B. Yagnik, M. Garg, and N. C. Krishnan, "Spotgarbage: smartphone app to detect garbage using deep learning," in *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, 2016, pp. 940–945.
- [39] S. Yao, Y. Zhao, A. Zhang, L. Su, and T. Abdelzaher, "Deepiot: Compressing deep neural network structures for sensing systems with a compressor-critic framework," in *Proceedings of the 15th ACM Conference on Embedded Network Sensor Systems*, 2017, pp. 1–14.
- [40] S. Han, H. Mao, and W. J. Dally, "Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding," *arXiv preprint arXiv:1510.00149*, 2015.
- [41] A. Paszke, A. Chaurasia, S. Kim, and E. Culurciello, "Enet: A deep neural network architecture for real-time semantic segmentation," *arXiv preprint arXiv:1606.02147*, 2016.
- [42] C. N. Duong, K. G. Quach, I. Jalata, N. Le, and K. Luu, "Mobiface: A lightweight deep learning face recognition on mobile devices," *arXiv preprint arXiv:1811.11080*, 2018.
- [43] R. Nagpal, C. K. Paturu, V. Ragavan, R. Bhat, D. Ghosh *et al.*, "Real-time traffic sign recognition using deep network for embedded platforms," *Electronic Imaging*, vol. 2019, no. 15, pp. 33–1, 2019.
- [44] X. Zhang, X. Zhou, M. Lin, and J. Sun, "Shufflenet: An extremely efficient convolutional neural network for mobile devices," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 6848–6856.
- [45] Z. Qin, Z. Li, Z. Zhang, Y. Bao, G. Yu, Y. Peng, and J. Sun, "Thundernet: Towards real-time generic object detection on mobile devices," in *Proceedings of the IEEE International Conference on Computer Vision*, 2019, pp. 6718–6727.
- [46] A. Sehgal and N. Kehtarnavaz, "A convolutional neural network smartphone app for real-time voice activity detection," *IEEE Access*, vol. 6, pp. 9017–9026, 2018.
- [47] G. S. Bhat, N. Shankar, C. K. Reddy, and I. M. Panahi, "A real-time convolutional neural network based speech enhancement for hearing impaired listeners using smartphone," *IEEE Access*, vol. 7, pp. 78 421–78 433, 2019.
- [48] J. Y. Han, W. Z. Zheng, R. J. Huang, Y. Tsao, and Y. H. Lai, "Hearing aids app design based on deep learning technology," in *2018 11th International Symposium on Chinese Spoken Language Processing (ISCSLP)*. IEEE, 2018, pp. 495–496.
- [49] N. D. Lane, P. Georgiev, and L. Qendro, "Deeppear: robust smartphone audio sensing in unconstrained acoustic environments using deep learning," in *Proceedings of the 2015 ACM International Joint Conference on Pervasive and Ubiquitous Computing*, 2015, pp. 283–294.
- [50] J. M. Martín-Doñas, A. M. Gomez, I. López-Espejo, and A. M. Peinado, "Dual-channel dnn-based speech enhancement for smartphones," in *2017 IEEE 19th International Workshop on Multimedia Signal Processing (MMSP)*. IEEE, 2017, pp. 1–6.
- [51] E. Malykh, S. Novoselov, and O. Kudashev, "On residual cnn in text-dependent speaker verification task," in *International Conference on Speech and Computer*. Springer, 2017, pp. 593–601.
- [52] S. Novoselov, O. Kudashev, V. Shchemelinin, I. Kremnev, and G. Lavrentyeva, "Deep cnn based feature extractor for text-prompted speaker recognition," in *2018 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2018, pp. 5334–5338.
- [53] R. Prabhavalkar, O. Alsharif, A. Bruguier, and L. McGraw, "On the compression of recurrent neural networks with an application to lvc sr acoustic modeling for embedded speech recognition," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 5970–5974.
- [54] Z. Lu, V. Sindhwani, and T. N. Sainath, "Learning compact recurrent neural networks," in *2016 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2016, pp. 5960–5964.
- [55] G. Chen, C. Parada, and G. Heigold, "Small-footprint keyword spotting using deep neural networks," in *2014 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2014, pp. 4087–4091.
- [56] M. Wang, T. Sirlapu, A. Kwasniewska, M. Szankin, M. Bartscherer, and R. Nicolas, "Speaker recognition using convolutional neural network with minimal training data for smart home solutions," in *2018 11th International Conference on Human System Interaction (HSI)*. IEEE, 2018, pp. 139–145.
- [57] K. Nan, S. Liu, J. Du, and H. Liu, "Deep model compression for mobile platforms: A survey," *Tsinghua Science and Technology*, vol. 24, no. 6, pp. 677–693, 2019.
- [58] C. Chen, P. Zhang, H. Zhang, J. Dai, Y. Yi, H. Zhang, and Y. Zhang, "Deep learning on computational-resource-limited platforms: A survey," *Mobile Information Systems*, vol. 2020, 2020.
- [59] W. Dai and D. Berleant, "Benchmarking contemporary deep learning hardware and frameworks: A survey of qualitative metrics," in *2019 IEEE First International Conference on Cognitive Machine Intelligence (CogMI)*. IEEE, 2019, pp. 148–155.
- [60] Z. Wang, K. Liu, J. Li, Y. Zhu, and Y. Zhang, "Various frameworks and libraries of machine learning and deep learning: A survey," *Archives of computational methods in engineering*, pp. 1–24, 2019.
- [61] C. Zhang, P. Patras, and H. Haddadi, "Deep learning in mobile and wireless networking: A survey," *IEEE Communications surveys & tutorials*, vol. 21, no. 3, pp. 2224–2287, 2019.

- [62] J. Chen and X. Ran, "Deep learning with edge computing: A review," *Proceedings of the IEEE*, vol. 107, no. 8, pp. 1655–1674, 2019.
- [63] S. Deng, H. Zhao, W. Fang, J. Yin, S. Dustdar, and A. Y. Zomaya, "Edge intelligence: The confluence of edge computing and artificial intelligence," *IEEE Internet of Things Journal*, vol. 7, no. 8, pp. 7457–7469, 2020.
- [64] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam, "Mobilenets: Efficient convolutional neural networks for mobile vision applications," *arXiv preprint arXiv:1704.04861*, 2017.
- [65] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.
- [66] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," in *Advances in neural information processing systems*, 2012, pp. 1097–1105.
- [67] A. K. Gupta, K. Gupta, J. Jadhav, R. V. Deolekar, A. Nerurkar, and S. Deshpande, "Plant disease prediction using deep learning and iot," in *2019 6th International Conference on Computing for Sustainable Global Development (INDIACom)*. IEEE, 2019, pp. 902–907.
- [68] B. Blanco-Filgueira, D. García-Lesta, M. Fernández-Sanjurjo, V. M. Brea, and P. López, "Deep learning-based multiple object visual tracking on embedded system for iot and mobile edge computing applications," *IEEE Internet of Things Journal*, vol. 6, no. 3, pp. 5423–5431, 2019.
- [69] J. Su, V. D. Vasconcellos, S. Prasad, S. Daniele, Y. Feng, and K. Sakurai, "Lightweight classification of iot malware based on image recognition," in *2018 IEEE 42nd Annual Computer Software and Applications Conference (COMPSAC)*, vol. 2. IEEE, 2018, pp. 664–669.
- [70] T. Na, J. H. Ko, and S. Mukhopadhyay, "Noise-robust and resolution-invariant image classification with pixel-level regularization," in *International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, 2018.
- [71] M. Song, K. Zhong, J. Zhang, Y. Hu, D. Liu, W. Zhang, J. Wang, and T. Li, "In-situ ai: Towards autonomous and incremental deep learning for iot systems," in *2018 IEEE International Symposium on High Performance Computer Architecture (HPCA)*. IEEE, 2018, pp. 92–103.
- [72] C. M. Dourado Jr, S. P. P. da Silva, R. V. M. da Nóbrega, A. C. d. S. Barros, P. P. Reboucas Filho, and V. H. C. de Albuquerque, "Deep learning iot system for online stroke detection in skull computed tomography images," *Computer Networks*, vol. 152, pp. 25–39, 2019.
- [73] Y. Wu, X. Zhang, Y. Xiao, and J. Feng, "Attention neural network for water image classification under iot environment," *Applied Sciences*, vol. 10, no. 3, p. 909, 2020.
- [74] Y. Chen and Y. Xue, "A deep learning approach to human activity recognition based on single accelerometer," in *2015 IEEE International Conference on Systems, Man, and Cybernetics*. IEEE, 2015, pp. 1488–1492.
- [75] B. Almaslukh, A. M. Artoli, and J. Al-Muhtadi, "A robust deep learning approach for position-independent smartphone-based human activity recognition," *Sensors*, vol. 18, no. 11, p. 3726, 2018.
- [76] M. M. Hassan, M. Z. Uddin, A. Mohamed, and A. Almogren, "A robust human activity recognition system using smartphone sensors and deep learning," *Future Generation Computer Systems*, vol. 81, pp. 307–313, 2018.
- [77] M. Kim, C. Y. Jeong, and H. C. Shin, "Activity recognition using fully convolutional network from smartphone accelerometer," in *2018 International Conference on Information and Communication Technology Convergence (ICTC)*. IEEE, 2018, pp. 1482–1484.
- [78] B. Almaslukh, J. AlMuhtadi, and A. Artoli, "An effective deep autoencoder approach for online smartphone-based human activity recognition," *Int. J. Comput. Sci. Netw. Secur.*, vol. 17, no. 4, pp. 160–165, 2017.
- [79] M.-C. Kwon, M. Ju, and S. Choi, "Classification of various daily behaviors using deep learning and smart watch," in *2017 Ninth International Conference on Ubiquitous and Future Networks (ICUFN)*. IEEE, 2017, pp. 735–740.
- [80] Q. Zou, Y. Wang, Q. Wang, Y. Zhao, and Q. Li, "Deep learning based gait recognition using smartphones in the wild," *arXiv preprint arXiv:1811.00338*, 2018.
- [81] S. Cortés, A. Solin, and J. Kannala, "Deep learning based speed estimation for constraining strapdown inertial navigation on smartphones," in *2018 IEEE 28th International Workshop on Machine Learning for Signal Processing (MLSP)*. IEEE, 2018, pp. 1–6.
- [82] D. Ravi, C. Wong, B. Lo, and G.-Z. Yang, "Deep learning for human activity recognition: A resource efficient implementation on low-power devices," in *2016 IEEE 13th international conference on wearable and implantable body sensor networks (BSN)*. IEEE, 2016, pp. 71–76.
- [83] S. Wan, Y. Liang, Y. Zhang, and M. Guizani, "Deep multi-layer perceptron classifier for behavior analysis to estimate parkinson's disease severity using smartphones," *IEEE Access*, vol. 6, pp. 36 825–36 833, 2018.
- [84] C. Stamate, G. D. Magoulas, S. Küppers, E. Nomikou, I. Daskalopoulos, M. U. Luchini, T. Moussouri, and G. Roussos, "Deep learning parkinson's from smartphone data," in *2017 IEEE International Conference on Pervasive Computing and Communications (PerCom)*. IEEE, 2017, pp. 31–40.
- [85] A. Shrestha and M. Won, "Deepwalking: Enabling smartphone-based walking speed estimation using deep learning," in *2018 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2018, pp. 1–6.
- [86] W. Shao, H. Luo, F. Zhao, C. Wang, A. Crivello, and M. Z. Tunio, "Depedo: Anti periodic negative-step movement pedometer with deep convolutional neural networks," in *2018 IEEE international conference on communications (ICC)*. IEEE, 2018, pp. 1–6.
- [87] T. Yoshida, J. Nozaki, K. Urano, K. Hiroi, T. Yonezawa, and N. Kawaguchi, "Gait dependency of smartphone walking speed estimation using deep learning (poster)," in *Proceedings of the 17th Annual International Conference on Mobile Systems, Applications, and Services*, 2019, pp. 641–642.
- [88] W. Jiang and Z. Yin, "Human activity recognition using wearable sensors by deep convolutional neural networks," in *Proceedings of the 23rd ACM international conference on Multimedia*, 2015, pp. 1307–1310.
- [89] C. A. Ronao and S.-B. Cho, "Human activity recognition with smartphone sensors using deep learning neural networks," *Expert systems with applications*, vol. 59, pp. 235–244, 2016.
- [90] J. Yang, K. Cheng, J. Chen, B. Zhou, and Q. Li, "Smartphones based online activity recognition for indoor localization using deep convolutional neural network," in *2018 Ubiquitous Positioning, Indoor Navigation and Location-Based Services (UPINLBS)*. IEEE, 2018, pp. 1–7.
- [91] V. Radu, N. D. Lane, S. Bhattacharya, C. Mascolo, M. K. Marina, and F. Kawsar, "Towards multimodal deep learning for activity recognition on mobile devices," in *Proceedings of the 2016 ACM International Joint Conference on Pervasive and Ubiquitous Computing: Adjunct*, 2016, pp. 185–188.
- [92] P. Vepakomma, D. De, S. K. Das, and S. Bhansali, "A-wristocracy: Deep learning on wrist-worn sensing for recognition of user complex activities," in *2015 IEEE 12th International conference on wearable and implantable body sensor networks (BSN)*. IEEE, 2015, pp. 1–6.
- [93] P. Sundaramoorthy, G. K. Gudur, M. R. Moorthy, R. N. Bhandari, and V. Vijayaraghavan, "Harnet: Towards on-device incremental learning using deep ensembles on constrained devices," in *Proceedings of the 2nd International Workshop on Embedded and Mobile Deep Learning*, 2018, pp. 31–36.
- [94] C. Xu, D. Chai, J. He, X. Zhang, and S. Duan, "Innohar: a deep neural network for complex human activity recognition," *Ieee Access*, vol. 7, pp. 9893–9902, 2019.
- [95] Z. Peng, S. Gao, Z. Li, B. Xiao, and Y. Qian, "Vehicle safety improvement through deep learning and mobile sensing," *IEEE network*, vol. 32, no. 4, pp. 28–33, 2018.
- [96] X. Liang and G. Wang, "A convolutional neural network for transportation mode detection based on smartphone platform," in *2017 IEEE 14th international conference on mobile Ad Hoc and sensor systems (MASS)*. IEEE, 2017, pp. 338–342.
- [97] X. Liang, Y. Zhang, G. Wang, and S. Xu, "A deep learning model for transportation mode detection based on smartphone sensing data," *IEEE Transactions on Intelligent Transportation Systems*, 2019.
- [98] S.-R. G. Christopoulos, S. Kanarachos, and A. Chronos, "Learning driver braking behavior using smartphones, neural networks and the sliding correlation coefficient: Road anomaly case study," *IEEE Transactions on Intelligent Transportation Systems*, vol. 20, no. 1, pp. 65–74, 2018.
- [99] S.-H. Fang, Y.-X. Fei, Z. Xu, and Y. Tsao, "Learning transportation modes from smartphone sensors based on deep neural network," *IEEE Sensors Journal*, vol. 17, no. 18, pp. 6111–6118, 2017.
- [100] H. Zhao, C. Hou, H. Alrobassy, and X. Zeng, "Recognition of transportation state by smartphone sensors using deep bi-lstm neural network," *Journal of Computer Networks and Communications*, vol. 2019, 2019.
- [101] W. Yuan, Y. Jiang, H. Li, and M. Cai, "A lightweight on-device detection method for android malware," *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2019.
- [102] E. B. Karbab, M. Debbabi, A. Derhab, and D. Mouheb, "Maldozer: Automatic framework for android malware detection using deep learning," *Digital Investigation*, vol. 24, pp. S48–S59, 2018.



- [103] N. McLaughlin, J. Martinez del Rincon, B. Kang, S. Yerima, P. Miller, S. Sezer, Y. Safaei, E. Trickle, Z. Zhao, A. Doupé *et al.*, “Deep android malware detection,” in *Proceedings of the Seventh ACM on Conference on Data and Application Security and Privacy*, 2017, pp. 301–308.
- [104] R. Feng, S. Chen, X. Xie, L. Ma, G. Meng, Y. Liu, and S.-W. Lin, “Mobidroid: A performance-sensitive malware detection system on mobile platform,” in *2019 24th International Conference on Engineering of Complex Computer Systems (ICECCS)*. IEEE, 2019, pp. 61–70.
- [105] A. Gharib and A. Ghorbani, “Dna-droid: A real-time android ransomware detection framework,” in *International Conference on Network and System Security*. Springer, 2017, pp. 184–198.
- [106] M. Abhijith, B. K. Priya, and N. Ramasubramanian, “Malware detection in android using machine learning on chip,” in *Computing in Engineering and Technology*. Springer, 2020, pp. 287–295.
- [107] A. Rago, G. Piro, G. Boggia, and P. Dini, “Multi-task learning at the mobile edge: An effective way to combine traffic classification and prediction,” *IEEE Transactions on Vehicular Technology*, vol. 69, no. 9, pp. 10 362–10 374, 2020.
- [108] G. Aceto, D. Ciunzo, A. Montieri, and A. Pescapé, “Mobile encrypted traffic classification using deep learning: Experimental evaluation, lessons learned, and challenges,” *IEEE Transactions on Network and Service Management*, vol. 16, no. 2, pp. 445–458, 2019.
- [109] S. Yao, T. Wang, J. Li, and T. Abdelzaher, “Stardust: A deep learning serving system in iot: demo abstract,” in *Proceedings of the 17th Conference on Embedded Networked Sensor Systems*, 2019, pp. 402–403.
- [110] Z. Zhao, K. M. Barijough, and A. Gerstlauer, “Deepthings: Distributed adaptive deep learning inference on resource-constrained iot edge clusters,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 37, no. 11, pp. 2348–2359, 2018.
- [111] C. Gao, S. Braun, I. Kiselev, J. Anumula, T. Delbruck, and S.-C. Liu, “Real-time speech recognition for iot purpose using a delta recurrent neural network accelerator,” in *2019 IEEE International Symposium on Circuits and Systems (ISCAS)*. IEEE, 2019, pp. 1–5.
- [112] P. Wang, F. Ye, X. Chen, and Y. Qian, “Datanet: Deep learning based encrypted network traffic classification in sdn home gateway,” *IEEE Access*, vol. 6, pp. 55 380–55 391, 2018.
- [113] B. Mao, F. Tang, Z. M. Fadlullah, N. Kato, O. Akashi, T. Inoue, and K. Mizutani, “A novel non-supervised deep-learning-based network traffic control method for software defined wireless networks,” *IEEE Wireless Communications*, vol. 25, no. 4, pp. 74–81, 2018.
- [114] F. Tang, B. Mao, Z. M. Fadlullah, N. Kato, O. Akashi, T. Inoue, and K. Mizutani, “On removing routing protocol from future wireless networks: A real-time deep learning approach for intelligent traffic control,” *IEEE Wireless Communications*, vol. 25, no. 1, pp. 154–160, 2017.
- [115] G. Marín, P. Casas, and G. Capdehourat, “Rawpower: Deep learning based anomaly detection from raw network traffic measurements,” in *Proceedings of the ACM SIGCOMM 2018 Conference on Posters and Demos*, 2018, pp. 75–77.
- [116] M. Lotfollahi, M. J. Siavoshani, R. S. H. Zade, and M. Saberian, “Deep packet: A novel approach for encrypted traffic classification using deep learning,” *Soft Computing*, vol. 24, no. 3, pp. 1999–2012, 2020.
- [117] A. Ignatov, R. Timofte, W. Chou, K. Wang, M. Wu, T. Hartley, and L. Van Gool, “Ai benchmark: Running deep neural networks on android smartphones,” in *Proceedings of the European Conference on Computer Vision (ECCV) Workshops*, 2018, pp. 0–0.
- [118] A. Ignatov, R. Timofte, A. Kulik, S. Yang, K. Wang, F. Baum, M. Wu, L. Xu, and L. Van Gool, “Ai benchmark: All about deep learning on smartphones in 2019,” in *2019 IEEE/CVF International Conference on Computer Vision Workshop (ICCVW)*. IEEE, 2019, pp. 3617–3635.
- [119] S. Venticinque, “Benchmarking physical and virtual iot platforms,” in *2019 IEEE International Conference on Cloud Engineering (IC2E)*. IEEE, 2019, pp. 247–252.
- [120] H. Kong, S. Huai, D. Liu, L. Zhang, H. Chen, S. Zhu, S. Li, W. Liu, M. Rastogi, R. Subramaniam *et al.*, “Edlab: A benchmark for edge deep learning accelerators,” *IEEE Design & Test*, 2021.
- [121] W. Niu, X. Ma, S. Lin, S. Wang, X. Qian, X. Lin, Y. Wang, and B. Ren, “Patdnn: Achieving real-time dnn execution on mobile devices with pattern-based weight pruning,” in *Proceedings of the Twenty-Fifth International Conference on Architectural Support for Programming Languages and Operating Systems*, 2020, pp. 907–922.
- [122] X. Cheng, W. Lu, X. Huang, S. Hu, and K. Chen, “Haflo: Gpu-based acceleration for federated logistic regression,” *arXiv preprint arXiv:2107.13797*, 2021.
- [123] J. Qiu, J. Wang, S. Yao, K. Guo, B. Li, E. Zhou, J. Yu, T. Tang, N. Xu, S. Song *et al.*, “Going deeper with embedded fpga platform for convolutional neural network,” in *Proceedings of the 2016 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2016, pp. 26–35.
- [124] C. Zhang, P. Li, G. Sun, Y. Guan, B. Xiao, and J. 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*, 2015, pp. 161–170.
- [125] F. Akopyan, J. Sawada, A. Cassidy, R. Alvarez-Icaza, J. Arthur, P. Merolla, N. Imam, Y. Nakamura, P. Datta, G.-J. Nam *et al.*, “Truenorth: Design and tool flow of a 65 mw 1 million neuron programmable neurosynaptic chip,” *IEEE transactions on computer-aided design of integrated circuits and systems*, vol. 34, no. 10, pp. 1537–1557, 2015.
- [126] S. Rivas-Gomez, A. J. Pena, D. Moloney, E. Laure, and S. Markidis, “Exploring the vision processing unit as co-processor for inference,” in *2018 IEEE International Parallel and Distributed Processing Symposium Workshops (IPDPSW)*. IEEE, 2018, pp. 589–598.
- [127] S. Han, X. Liu, H. Mao, J. Pu, A. Pedram, M. A. Horowitz, and W. J. Dally, “Eie: efficient inference engine on compressed deep neural network,” *ACM SIGARCH Computer Architecture News*, vol. 44, no. 3, pp. 243–254, 2016.
- [128] T. Chen, Z. Du, N. Sun, J. Wang, C. Wu, Y. Chen, and O. Temam, “Diannao: A small-footprint high-throughput accelerator for ubiquitous machine-learning,” *ACM SIGARCH Computer Architecture News*, vol. 42, no. 1, pp. 269–284, 2014.
- [129] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer, “Squeezenet: Alexnet-level accuracy with 50x fewer parameters and < 0.5 mb model size,” *arXiv preprint arXiv:1602.07360*, 2016.
- [130] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, “Ssd: Single shot multibox detector,” in *European conference on computer vision*. Springer, 2016, pp. 21–37.
- [131] TensorFlow. (2020) Object detection. [Online]. Available: [https://www.tensorflow.org/lite/models/object\\_detection/overview](https://www.tensorflow.org/lite/models/object_detection/overview)
- [132] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, 2015, pp. 91–99.
- [133] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, “You only look once: Unified, real-time object detection,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 779–788.
- [134] J. Redmon and A. Farhadi, “Yolo9000: better, faster, stronger,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 7263–7271.
- [135] D. Anguita, A. Ghio, L. Oneto, X. Parra, and J. L. Reyes-Ortiz, “A public domain dataset for human activity recognition using smartphones,” in *Esann*, vol. 3, 2013, p. 3.
- [136] J. R. Kwapisz, G. M. Weiss, and S. A. Moore, “Activity recognition using cell phone accelerometers,” *ACM SigKDD Explorations Newsletter*, vol. 12, no. 2, pp. 74–82, 2011.
- [137] N. Mostafa and J. Slay, “Unsw-nb15: a comprehensive data set for network intrusion detection systems (unsw-nb15 network data set),” in *2015 military communications and information systems conference (MilCIS)*. IEEE, 2015, pp. 1–6.
- [138] M. Ring, S. Wunderlich, D. Grüdl, D. Landes, and A. Hotho, “Flow-based benchmark data sets for intrusion detection,” in *Proceedings of the 16th European conference on cyber warfare and security*, 2017, pp. 361–369.
- [139] N. D. Lane, S. Bhattacharya, A. Mathur, C. Forlivesi, and F. Kawsar, “Dxtk: Enabling resource-efficient deep learning on mobile and embedded devices with the deepx toolkit,” in *MobiCASE*, 2016, pp. 98–107.
- [140] J. Wang, B. Cao, P. Yu, L. Sun, W. Bao, and X. Zhu, “Deep learning towards mobile applications,” in *2018 IEEE 38th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2018, pp. 1385–1393.
- [141] S. Duan, D. Zhang, Y. Wang, L. Li, and Y. Zhang, “Jointrec: A deep-learning-based joint cloud video recommendation framework for mobile iot,” *IEEE Internet of Things Journal*, vol. 7, no. 3, pp. 1655–1666, 2019.
- [142] B. Shi, M. Sun, C.-C. Kao, V. Rozgic, S. Matsoukas, and C. Wang, “Compression of acoustic event detection models with low-rank matrix factorization and quantization training,” *arXiv preprint arXiv:1905.00855*, 2019.
- [143] F. Tung and G. Mori, “Clip-q: Deep network compression learning by in-parallel pruning-quantization,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7873–7882.
- [144] G. Hinton, O. Vinyals, and J. Dean, “Distilling the knowledge in a neural network,” *arXiv preprint arXiv:1503.02531*, 2015.

- [145] G. Chen, W. Choi, X. Yu, T. Han, and M. Chandraker, "Learning efficient object detection models with knowledge distillation," in *Advances in Neural Information Processing Systems*, 2017, pp. 742–751.
- [146] A. Ashok, N. Rhinehart, F. Beainy, and K. M. Kitani, "N2n learning: Network to network compression via policy gradient reinforcement learning," *arXiv preprint arXiv:1709.06030*, 2017.
- [147] B. Moons, B. De Brabandere, L. Van Gool, and M. Verhelst, "Energy-efficient convnets through approximate computing," in *2016 IEEE Winter Conference on Applications of Computer Vision (WACV)*. IEEE, 2016, pp. 1–8.
- [148] K. Wang, Z. Liu, Y. Lin, J. Lin, and S. Han, "Hq: Hardware-aware automated quantization with mixed precision," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2019, pp. 8612–8620.
- [149] I. Markovsky and K. Usevich, *Low rank approximation*. Springer, 2012, vol. 139.
- [150] A. Yaguchi, T. Suzuki, S. Nitta, Y. Sakata, and A. Tanizawa, "Scalable deep neural networks via low-rank matrix factorization," *arXiv preprint arXiv:1910.13141*, 2019.
- [151] W. Shi, Y. Hou, S. Zhou, Z. Niu, Y. Zhang, and L. Geng, "Improving device-edge cooperative inference of deep learning via 2-step pruning," *arXiv preprint arXiv:1903.03472*, 2019.
- [152] Y. Weiss, B. Schölkopf, and J. Platt, "Advances in neural information processing systems 18: Proceedings of the 2005 conference," in *Nineteenth Annual Conference on Neural Information Processing Systems (NIPS 2005)*. Mit Press, 2006.
- [153] X. Lin, C. Zhao, and W. Pan, "Towards accurate binary convolutional neural network," *arXiv preprint arXiv:1711.11294*, 2017.
- [154] R. Ding, T.-W. Chin, Z. Liu, and D. Marculescu, "Regularizing activation distribution for training binarized deep networks," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 11408–11417.
- [155] Q. V. Le, J. Ngiam, A. Coates, A. Lahiri, B. Prochnow, and A. Y. Ng, "On optimization methods for deep learning," in *ICML*, 2011.
- [156] S. R. Young, D. C. Rose, T. P. Karnowski, S.-H. Lim, and R. M. Patton, "Optimizing deep learning hyper-parameters through an evolutionary algorithm," in *Proceedings of the workshop on machine learning in high-performance computing environments*, 2015, pp. 1–5.
- [157] S. S. Latifi Oskouei, H. Golestani, M. Hashemi, and S. Ghiasi, "Cn-droid: Gpu-accelerated execution of trained deep convolutional neural networks on android," in *Proceedings of the 24th ACM international conference on Multimedia*, 2016, pp. 1201–1205.
- [158] C. Holmes, D. Mawhirter, Y. He, F. Yan, and B. Wu, "Grnn: Low-latency and scalable rnn inference on gpus," in *Proceedings of the Fourteenth EuroSys Conference 2019*, 2019, pp. 1–16.
- [159] P. Hill, A. Jain, M. Hill, B. Zamirai, C.-H. Hsu, M. A. Laurenzano, S. Mahlke, L. Tang, and J. Mars, "Deftnn: Addressing bottlenecks for dnn execution on gpus via synapse vector elimination and near-compute data fission," in *Proceedings of the 50th Annual IEEE/ACM International Symposium on Microarchitecture*, 2017, pp. 786–799.
- [160] T. Luo, S. Liu, L. Li, Y. Wang, S. Zhang, T. Chen, Z. Xu, O. Temam, and Y. Chen, "Dadiannao: A neural network supercomputer," *IEEE Transactions on Computers*, vol. 66, no. 1, pp. 73–88, 2016.
- [161] Y.-H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolutional neural networks," *IEEE journal of solid-state circuits*, vol. 52, no. 1, pp. 127–138, 2016.
- [162] Y. Ma, Y. Cao, S. Vrudhula, and J.-s. Seo, "Optimizing loop operation and dataflow in fpga acceleration of deep convolutional neural networks," in *Proceedings of the 2017 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, 2017, pp. 45–54.
- [163] TensorFlow. (2020) Tensorflow hub. [Online]. Available: <https://tfhub.dev/>
- [164] Keras. (2020) Keras applications. [Online]. Available: <https://keras.io/api/applications/>
- [165] Caffe2. (2020) Caffe2 model zoo. [Online]. Available: <https://caffe2.ai/docs/zoo.html>
- [166] (2020) Model optimization. [Online]. Available: [https://www.tensorflow.org/lite/performance/model\\_optimization](https://www.tensorflow.org/lite/performance/model_optimization)
- [167] (2020) Post-training quantization. [Online]. Available: [https://www.tensorflow.org/lite/performance/post\\_training\\_quantization](https://www.tensorflow.org/lite/performance/post_training_quantization)
- [168] (2020) Trim insignificant weights. [Online]. Available: [https://www.tensorflow.org/model\\_optimization/guide/pruning](https://www.tensorflow.org/model_optimization/guide/pruning)
- [169] PyTorch, "quantization," <https://pytorch.org/docs/stable/quantization>.
- [170] I. A. Lab, "Using distiller to prune a pytorch language model," <https://github.com/NervanaSystems/distiller/>
- [171] (2020) Caffe2 quantization. [Online]. Available: <https://github.com/pytorch/pytorch/tree/master/caffe2/quantization/server>
- [172] Microsoft, "Mmdnn," <https://github.com/microsoft/MMdnn>.
- [173] F. Microsoft, "Onnx," <https://onnx.ai/>.
- [174] Z. Xu, F. Yu, J. Xiong, and X. Chen, "Helios: Heterogeneity-aware federated learning with dynamically balanced collaboration," *arXiv preprint arXiv:1912.01684*, 2019.
- [175] M. Sankupellay and D. Konovalov, "Bird call recognition using deep convolutional neural network, resnet-50," in *Proc. ACOUSTICS*, vol. 7, 2018, pp. 1–8.
- [176] A. Imteaj and M. H. Amini, "Fedparl: Client activity and resource-oriented lightweight federated learning model for resource-constrained heterogeneous iot environment," *Frontiers in Communications and Networks*, vol. 2, p. 10, 2021.
- [177] A. Hard, K. Rao, R. Mathews, S. Ramaswamy, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [178] Q. Wu, K. He, and X. Chen, "Personalized federated learning for intelligent iot applications: A cloud-edge based framework," *IEEE Open Journal of the Computer Society*, vol. 1, pp. 35–44, 2020.
- [179] Y. Jiang, S. Wang, V. Valls, B. J. Ko, W.-H. Lee, K. K. Leung, and L. Tassiulas, "Model pruning enables efficient federated learning on edge devices," *arXiv preprint arXiv:1909.12326*, 2019.
- [180] S. Wu, G. Li, F. Chen, and L. Shi, "Training and inference with integers in deep neural networks," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=HJGXzmspb>
- [181] R. Raju, D. Gope, U. Thakker, and J. Beu, "Understanding the impact of dynamic channel pruning on conditionally parameterized convolutions," in *Proceedings of the 2nd International Workshop on Challenges in Artificial Intelligence and Machine Learning for Internet of Things*, 2020, pp. 27–33.
- [182] X. Gao, Y. Zhao, Ł. Dudziak, R. Mullins, and C.-z. Xu, "Dynamic channel pruning: Feature boosting and suppression," *arXiv preprint arXiv:1810.05331*, 2018.
- [183] X. Wang, F. Yu, Z.-Y. Dou, T. Darrell, and J. E. Gonzalez, "Skipnet: Learning dynamic routing in convolutional networks," in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 409–424.
- [184] Z. Yang, S. Hu, and K. Chen, "Fpga-based hardware accelerator of homomorphic encryption for efficient federated learning," *arXiv preprint arXiv:2007.10560*, 2020.
- [185] J. Mao, X. Chen, K. W. Nixon, C. Krieger, and Y. Chen, "Modnn: Local distributed mobile computing system for deep neural network," in *Design, Automation & Test in Europe Conference & Exhibition (DATE)*, 2017. IEEE, 2017, pp. 1396–1401.
- [186] R. Hadidi, J. Cao, M. Woodward, M. S. Ryoo, and H. Kim, "Distributed perception by collaborative robots," *IEEE Robotics and Automation Letters*, vol. 3, no. 4, pp. 3709–3716, 2018.
- [187] L. Zeng, X. Chen, Z. Zhou, L. Yang, and J. Zhang, "Coedge: Cooperative dnn inference with adaptive workload partitioning over heterogeneous edge devices," *IEEE/ACM Transactions on Networking*, vol. 29, no. 2, pp. 595–608, 2020.
- [188] Y. Meidan, M. Bohadana, Y. Mathov, Y. Mirsky, A. Shabtai, D. Breitenbacher, and Y. Elovici, "N-baiot—network-based detection of iot botnet attacks using deep autoencoders," *IEEE Pervasive Computing*, vol. 17, no. 3, pp. 12–22, 2018.
- [189] B. Roy and H. Cheung, "A deep learning approach for intrusion detection in internet of things using bi-directional long short-term memory recurrent neural network," in *2018 28th International Telecommunication Networks and Applications Conference (ITNAC)*. IEEE, 2018, pp. 1–6.
- [190] Y. Guo, Y. Liu, A. Oerlemans, S. Lao, S. Wu, and M. S. Lew, "Deep learning for visual understanding: A review," *Neurocomputing*, vol. 187, pp. 27–48, 2016.
- [191] V. Sze, Y.-H. Chen, T.-J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295–2329, 2017.
- [192] N. D. Lane, S. Bhattacharya, A. Mathur, P. Georgiev, C. Forlivesi, and F. Kawsar, "Squeezing deep learning into mobile and embedded devices," *IEEE Pervasive Computing*, vol. 16, no. 3, pp. 82–88, 2017.
- [193] F. Tu, W. Wu, S. Yin, L. Liu, and S. Wei, "Rana: Towards efficient neural acceleration with refresh-optimized embedded dram," in *2018 ACM/IEEE 45th Annual International Symposium on Computer Architecture (ISCA)*. IEEE, 2018, pp. 340–352.
- [194] D. Marculescu, D. Stamouli, and E. Cai, "Hardware-aware machine learning: modeling and optimization," in *Proceedings of the International Conference on Computer-Aided Design*, 2018, pp. 1–8.



**Tianming Zhao** received his B.E. from the Department of Software Engineering, Chongqing University, China and his M.S. degrees from the Department of computer science, Binghamton University, USA. Currently, he is a Ph.D. student of the Computer & Information Sciences Department at Temple University. His research interests include mobile computing and sensing, cyber security and privacy, and smart health. He is currently working with Prof. Yan Wang.



**Xiaonan Guo** received the Ph.D. degree in computer science and engineering from the Hong Kong University of Science and Technology, Hong Kong, in 2013. He is currently an assistant professor of the Department of Computer, Information and Technology at Indiana University-Purdue University Indianapolis (IUPUI). Prior to joining IUPUI, he was a research associate in the Electrical and Computer Engineering Department, Stevens Institute of Technology. His research interests include pervasive computing, mobile computing, and cyber security and privacy. He received the Best Paper Award from the ACM Conference on Information, Computer, and Communications Security (ASIACCS) 2016 and EAI International Conference on IoT Technologies for HealthCare (EAI Healthy IoT), 2019. He is a member of the IEEE.



**Yucheng Xie** is currently a Ph.D. student of the Electrical & Computer Engineering Department in the Purdue School of Engineering & Technology at Indiana University Purdue University Indianapolis. He received his M.S. degrees from the Department of Computer Science at Stevens Institute of Technology, USA. His research interests include smart healthcare, mobile sensing and computing, cyber security and privacy. He is currently working with Prof. Xiaonan Guo.



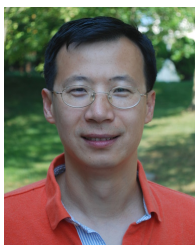
**Bin Hu** received the B. Sc. degree in mechanical engineering from Xi'an Jiaotong University, China in 2000. He received the two M. Sc. degrees in software engineering from Xi'an Jiaotong University, China in 2005 and Monmouth University, USA in 2018, respectively. From 2006 to 2016, he was an assistant professor at Xi'an University of Posts and Telecommunications, China. Currently, he is a Ph.D. student of Electrical and Computer Engineering at Rutgers University. His research interests include mobile computing and sensing, robotics, and machine learning. He is currently working with Prof. Yingying Chen.



**Yan Wang** is an Assistant Professor with Computer & Information Sciences Department at Temple University. Before that, he was with Department of Computer Science at SUNY, Binghamton. He received his Ph.D. degree in Electrical Engineering from Stevens Institute of Technology. His research interests include Cyber Security and Privacy, Internet of Things, Mobile and Pervasive Computing, and Smart Healthcare. His research is supported by the National Science Foundation (NSF). He is the recipient of the Best Paper Award from IEEE CNS 2018, IEEE SECON 2017 and ACM AsiaCCS 2016. He is the Winner of ACM MobiCom Student Research Competition, 2013. His research has been reported in numerous media outlets including MIT Technology Review, CNN, Fox News Channel, Wall Street Journal, National Public Radio, and IEEE Spectrum.



**Yingying (Jennifer) Chen** is a Professor of Electrical and Computer Engineering and Peter Cherasia Endowed Faculty Scholar at Rutgers University. She is the Associate Director of Wireless Information Network Laboratory (WINLAB). She also leads the Data Analysis and Information Security (DAISY) Lab. She is an IEEE Fellow. Her research interests include mobile sensing and computing, cyber security and privacy, Internet of Things, and smart healthcare. Her background is a combination of Computer Science, Computer Engineering and Physics. She had extensive industry experiences at Nokia previously. She has published over 200 journal articles and conference papers. She is the recipient of multiple Best Paper Awards from EAI HealthyIoT 2019, IEEE CNS 2018, IEEE SECON 2017, ACM AsiaCCS 2016, IEEE CNS 2014 and ACM MobiCom 2011. She is also the recipient of NSF CAREER Award and Google Faculty Research Award. She received NJ Inventors Hall of Fame Innovator Award and is also the recipient of IEEE Region 1 Technological Innovation in Academic Award. Her research has been reported in numerous media outlets including MIT Technology Review, CNN, Fox News Channel, Wall Street Journal, National Public Radio and IEEE Spectrum. She has been serving/served on the editorial boards of IEEE Transactions on Mobile Computing (IEEE TMC), IEEE Transactions on Wireless Communications (IEEE TWireless), IEEE/ACM Transactions on Networking (IEEE/ACM ToN) and ACM Transactions on Privacy and Security.



**Jerry Cheng** is an Assistant Professor of Computer Science at New York Institute of Technology. His research interests include big data analytics, statistical learning, Bayesian statistics, and their applications in computer systems and smart healthcare. He was an Assistant Professor in Robert Wood Johnson Medical School at Rutgers University. He was formerly a Postdoc Researcher in the Department of Statistics at Columbia University and had extensive industrial experiences as a Member of Technical Staff at AT&T Labs. His background is a combination of Computer Science, Statistics and Physics. His work has been reported by many new media including MIT Technology Review, Yahoo News, Digital World, FierceHealthcare, and WTOP Radio.