

Edge Intelligence: Architectures, Challenges, and Applications

Dianlei Xu, Tong Li, Yong Li, *Senior Member, IEEE*, Xiang Su, *Member, IEEE*, Sasu Tarkoma, *Senior Member, IEEE*, Tao Jiang, *Fellow, IEEE*, Jon Crowcroft, *Fellow, IEEE*, and Pan Hui, *Fellow, IEEE*

Abstract—Edge intelligence refers to a set of connected systems and devices for data collection, caching, processing, and analysis proximity to where data is captured based on artificial intelligence. Edge intelligence aims at enhancing data processing and protect the privacy and security of the data and users. Although recently emerged, spanning the period from 2011 to now, this field of research has shown explosive growth over the past five years. In this paper, we present a thorough and comprehensive survey on the literature surrounding edge intelligence. We first identify four fundamental components of edge intelligence, i.e. edge caching, edge training, edge inference, and edge offloading based on theoretical and practical results pertaining to proposed and deployed systems. We then aim for a systematic classification of the state of the solutions by examining research results and observations for each of the four components and present a taxonomy that includes practical problems, adopted techniques, and application goals. For each category, we elaborate, compare and analyse the literature from the perspectives of adopted techniques, objectives, performance, advantages and drawbacks, etc. This article provides a comprehensive survey to edge intelligence and its application areas. In addition, we summarise the development of the emerging research fields and the current state-of-the-art and discuss the important open issues and possible theoretical and technical directions.

Index Terms—Artificial intelligence, edge computing, edge caching, model training, inference, offloading

I. INTRODUCTION

WITH the breakthrough of Artificial Intelligence (AI), we are witnessing a booming increase in AI-based applications and services. AI technology, e.g., machine learning (ML) and deep learning (DL), achieves state-of-the-art performance in various fields, ranging from facial recognition [1], [2], natural language processing [3], [4], computer vision [5], [6], traffic prediction [7], [8], and anomaly detection [9], [10]. Benefiting from the services provided by these intelligent

applications and services, our lifestyles have been dramatically changed.

However, existing intelligent applications are computation-intensive, which present strict requirements on resources, e.g., CPU, GPU, memory, and network, which makes it impossible to be available anytime and anywhere for end users. Although current end devices are increasingly powerful, it is still insufficient to support some deep learning models. For example, most voice assistants, e.g., Apple Siri and Google Microsoft's Cortana, are based on cloud computing and they would not function if the network is unavailable. Moreover, existing intelligent applications generally adopt centralised data management, which requires users to upload their data to central cloud based data-centre. However, there is giant volume of data which has been generated and collected by billions of mobile users and Internet of Thing (IoT) devices distributed at the network edge. According to Cisco's forecast, there will be 850 ZB of data generated by mobile users and IoT devices by 2021 [11]. Uploading such volume of data to the cloud consumes significant bandwidth resources, which would also result in unacceptable latency for users. On the other hand, users increasingly concern their privacy. The European Union has promulgated General Data Protection Regulation (GDPR) to protect private information of users [12]. If mobile users upload their personal data to the cloud for a specific intelligent application, they would take the risk of privacy leakage, i.e., the personal data might be extracted by malicious hackers or companies for illegal purposes.

Edge computing [13]–[17] emerges as an extension of cloud computing to push cloud services to the proximity of end users. Edge computing offers virtual computing platforms which provide computing, storage, and networking resources, which are usually located at the edge of networks. The devices that provide services for end devices are referred to as edge servers, which could be IoT gateways, routers, and micro data centres in mobile network base stations, on vehicles, and amongst other places. End devices, such as mobile phones, IoT devices, and embedded devices that requests services from edge servers are called edge devices. The main advantages of the edge computing paradigm could be summarised into three aspects. (i) Ultra-low latency: computation usually takes place in the proximity of the source data, which saves substantial amounts of time on data transmission. Edge servers provides nearly real-time responses to end devices. (ii) Saving energy for end devices: since end devices could offload computing tasks to edge servers, the energy consumption on end devices would significantly shrink. Consequently, the battery life

D. Xu, T. Li, X. Su, and P. Hui are with the Department of Computer Science, University of Helsinki, 00014 Helsinki, Finland.(e-mail: dianlei.xu@helsinki.fi, t.li@connect.ust.hk, xiang.su@helsinki.fi, sasu.tarkoma@helsinki.fi, panhui@cse.ust.hk.)

D. Xu and Y. Li are with Beijing National Research Center for Information Science and Technology (BNRist), Department of Electronic Engineering, Tsinghua University, Beijing 100084, China.(e-mail: liyong07@tsinghua.edu.cn.)

T. Li and P. Hui are also with the Department of Computer Science and Engineering, Hong Kong University of Science and Technology, Hong Kong.

T. Jiang is with the School of Electronics Information and Communications, Huazhong University of Science and Technology, Wuhan 430074, China. (e-mail: taojiang@ieee.org)

J. Crowcroft is with the Computer Laboratory, University of Cambridge, William Gates Building, 15 JJ Thomson Avenue, Cambridge CB3 0FD, UK. (e-mail: Jon.Crowcroft@cl.cam.ac.uk.)

of end devices would be extended. (iii) Scalability: cloud computing is still available if there are no enough resource on edge devices or edge servers. In such a case, the cloud server would help to perform tasks. In addition, end devices with idle resources could communicate amongst themselves to collaboratively finish a task. The capability of the edge computing paradigm is flexible to accommodate different application scenarios.

Edge computing addresses the critical challenges of AI based applications and the combination of edge computing and AI provides a promising solution. This new paradigm of intelligence is called edge intelligence [18], [19], also named mobile intelligence [20]. Edge intelligence refers to a set of connected systems and devices for data collection, caching, processing, and analysis proximity to where data is collected, with the purpose of enhancing the quality and speed of data processing and to protect the privacy and security of data. Compared with traditional cloud-based intelligence that requires end devices to upload generated or collected data to the remote cloud, edge intelligence processes and analyses data locally, which effectively protects users' privacy, reduces response time, and saves on bandwidth resources [21], [22]. Moreover, users could also customise intelligent applications by training ML/DL models with self-generated data [23], [24]. It is predicted that edge intelligence will be a vital component in 6G network [25]. It is also worth noting that AI could also be a powerful assistance for edge computing. This paradigm is called intelligent edge [26], [27], which is different from edge intelligence. The emphasis of edge intelligence is to realize intelligent applications in edge environment with the assistance of edge computing and protect users' privacy, while intelligent edge focuses on solving problems of edge computing with AI solutions, e.g., resource allocation optimization. Intelligent edge is out of our scope in this survey.

There exists lots of works which have proved the feasibility of edge intelligence by applying an edge intelligence paradigm to practical application areas. Yi *et al.* implement a face recognition application across a smartphone and edge server [28]. Results show that the latency is reduced from 900ms to 169ms, compared with cloud based paradigm. Ha *et al.* use a cloudlet to help a wearable cognitive assistance execute recognition tasks, which saves energy consumption by 30%-40% [29]. Some researchers pay attention to the performance of AI in the context of edge computing. Lane *et al.* successfully implement a constrained DL model on smartphones for activity recognition [30]. The demo achieves a better performance than shallow models, which demonstrates that ordinary smart devices are qualified for simple DL models. Similar verification is also done on wearable devices [31] and embedded devices [32]. The most famous edge intelligence application is Google G-board, which uses federated learning [33] to collaboratively train the typing prediction model on smartphones. Each user uses their own typing records to train G-board. Hence, the trained G-board could be used immediately, powering experiences personalised by the way users use this application.

This paper aims at providing a comprehensive survey to the development and the state-of-the-art of edge intelligence.

As far as we know, there exist few recent efforts [26], [34]–[38] in this direction, but they have very different focuses from our survey. Table I summarizes the comparison among these works. Specifically, Yang *et al.* provide a survey on federated learning, in which they mainly focus on the architecture and applications of federated learning [34]. The authors divide literature of federated learning into three classifications: horizontal federated learning, vertical federated learning, and federated transfer learning. Federated learning is also involved as a collaborative training structure in our survey. We present how federated learning is applied in edge environment with the consideration of communication and privacy/security issues. The focus of [35] is how to realize the training and inference of DL models on a single mobile device. They briefly introduce some challenges and existing solutions from the perspective of training and inference. By contrast, we provide a more comprehensive and deeper review on solutions from the perspective of model design, model compression, and model acceleration. We also survey how to realize model training and inference with collaboration of edge devices and edge servers, even the assistance from the cloud server, in addition to solo training and inference at edge. Mohammadi *et al.* review works on IoT big data analytic with DL approaches [36]. Edge intelligence is not necessary in this work. The emphasis of survey [37] is how to use DL techniques to deal with the problems in wireless networks, e.g., spectrum resource allocation, which has no overlap with our work.

To our best knowledge, ref. [26] and [38] are two most relevant articles to our survey. The focus of [26] is the inter-availability between edge computing and DL. Hence the scope of ref. [26] includes two parts: DL for edge computing, and edge computing for DL. The former part focuses on some optimisation problems at edge with DL approaches, whilst the latter part focus on applying DL in the context of edge computing (i.e., techniques to perform DL at edge). The authors analyse these two parts from a macro view. By contrast, we pay more attention to the implementation of AI based applications and services (including ML and DL) with the assistance of edge resources from the micro view. More specifically, we provide more comprehensive and detailed classification and comparison on existing works of this research area from multi-dimensions. Not only the implementation of AI based applications and services (including both training and inference), but also the management of edge data and the required computing power are involved in our work. Moreover, statistically, there are only 40 coincident surveyed papers between [26] and our work. Similarly, the survey [38] analyses the implementation of edge intelligence on different layers from a macro view. They propose a six-level rating to describe edge intelligence. This is also involved in our work. Different from [38], we analyse its implementation from a micro view, e.g., offloading strategies and caching strategies.

Our survey focuses on how to realise edge intelligence in a systematic way. There exist three key components in AI, i.e. data, model/algorithm¹, and computation. A complete process of implementing AI applications involves data collection and

¹model and algorithm are interchangeable in this article

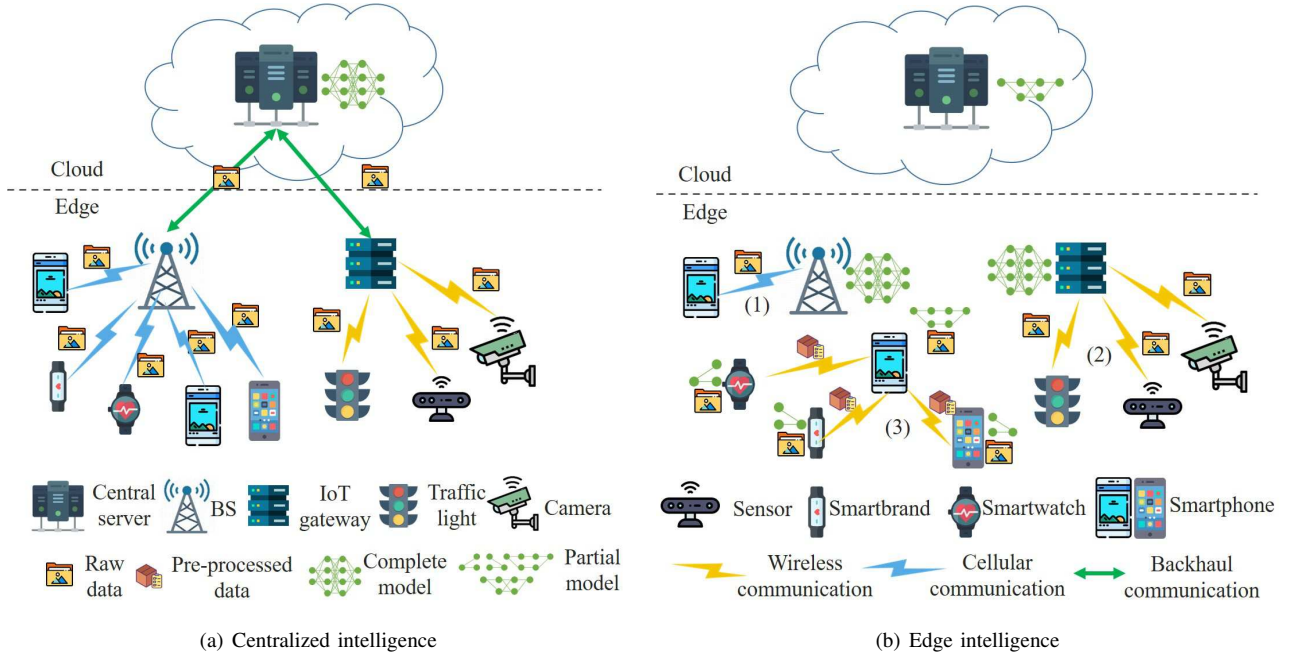


Fig. 1. The comparison of traditional intelligence and edge intelligence from the perspective of implementation. In traditional intelligence, all data must be uploaded to a central cloud server, whilst in edge intelligence, intelligent application tasks are done at the edge with locally-generated data in a distributed manner.

TABLE I
COMPARISON OF RELATIVE SURVEYS.

Ref.	Year	Domain	Scope	Analysing perspective
[34]	2019	Federated learning	Horizontal federated learning, vertical federated learning, and federated transfer learning	Macro-perspective
[35]	2018	DL-based mobile applications	Training and inference on single mobile device	Micro-perspective
[36]	2018	IoT big data	DL in IoT applications, and DL on IoT devices	Micro-perspective
[37]	2019	Intelligent wireless network	Algorithms that enables DL in wireless networks Applications ranging from traffic analytic to security	Micro-perspective
[26]	2019	Edge intelligence Intelligent edge	Training and inference systems DL for optimizing edge, and DL application on edge	Macro-perspective
[38]	2019	Edge intelligence	Cloud-edge-device coordination architecture Optimisation technologies in training and inference	Macro-perspective
Our work	2020	Edge intelligence	Edge caching, edge training, edge inference, and edge offloading	Micro-perspective

management, model training, and model inference. Computation plays an essential role throughout the whole process. Hence, we limit the scope of our survey on four aspects, including how to cache data to fuel intelligent applications (i.e., edge caching), how to train intelligent applications at the edge (i.e., edge training), how to infer intelligent applications at the edge (edge inference), and how to provide sufficient computing power for intelligent applications at the edge (edge offloading). Our contributions are summarized as following:

- We survey recent research achievements on edge intelligence and identify four key components: edge caching, edge training, edge inference, and edge offloading. For each component, we outline a systematical and comprehensive classification from a multi-dimensional view, e.g., practical challenges, solutions, optimisation goals, etc.
- We present thorough discussion and analysis on relevant papers in the field of edge intelligence from multiple views, e.g., applicable scenarios, methodology, performance, etc. and summarise their advantages and shortcomings.

- We discuss and summarise open issues and challenges in the implementation of edge intelligence, and outline five important future research directions and development trends, i.e., data scarcity, data consistency, adaptability of model/algorithms, privacy and security, and incentive mechanisms.

The remainder of this article is organized as follow. Section II overviews the research on edge intelligence, with considerations of the essential elements of edge intelligence, as well as the development situation of this research field. We present detailed introduction, discussion, and analysis on the development and recent advances of edge caching, edge training, edge inference, and edge offloading in Section III to Section VI, respectively. Finally, we discuss the open issue and possible solutions for future research in Section VII, and conclude the paper in Section VIII.

II. OVERVIEW

As an emerging research area, edge intelligence has received broad interests in the past few years. With benefits from edge

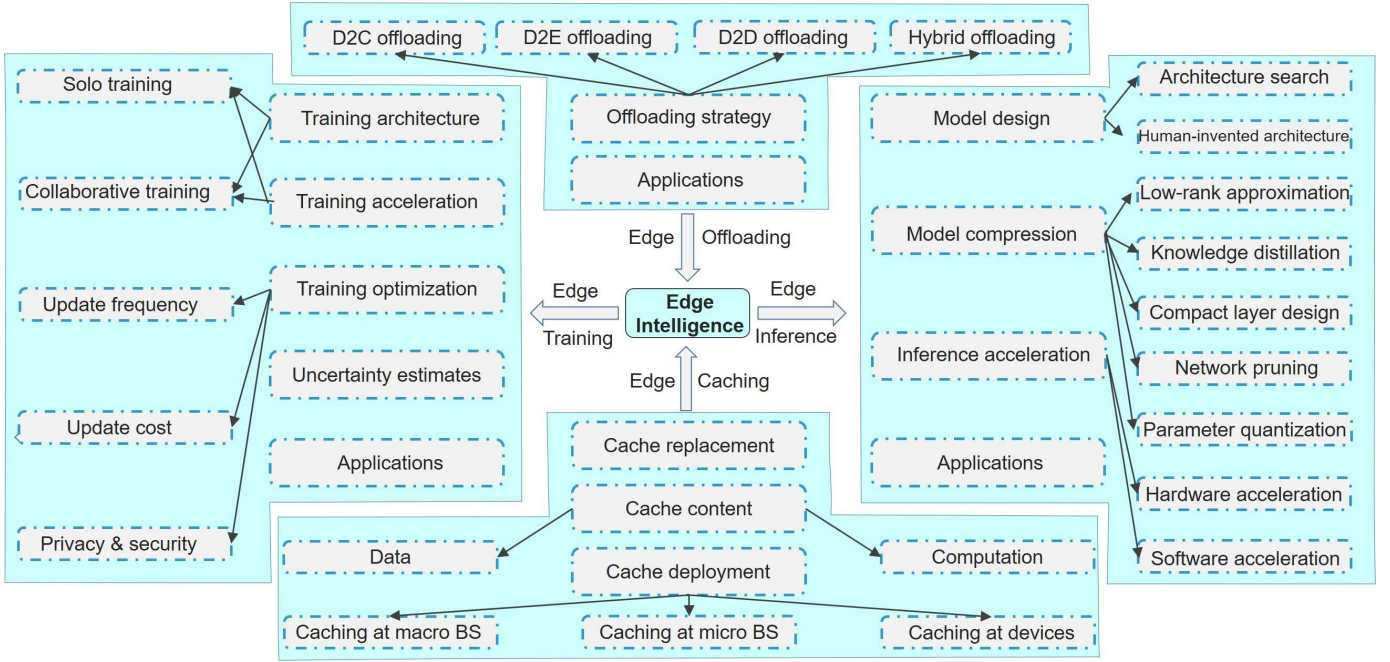


Fig. 2. The classification of edge intelligence literature.

computing and artificial intelligence techniques, combination of their contributions enables easy-to-use intelligent applications for users in daily lives and less dependent on the centralised cloud.

For convenience, we present the comparison between traditional centralised intelligence with edge intelligence from the perspective of implementation in Fig. 1. Traditional centralised intelligence is shown in Fig. 1(a), where all edge devices first upload data to the central server for intelligent tasks, e.g., model training or inference. The central server/data-centre is usually, but not necessarily, located in remote cloud. After the processing on the central server, results, e.g., recognition or prediction results, are transmitted back to edge devices. Fig. 1(b) demonstrates the implementation of edge intelligence, where a task, e.g., recognition and prediction is either done by edge servers and peer devices, or with the edge-cloud cooperation paradigm. A very small amount, or none of the data is uploaded to the cloud. For example, in area (1) and (2), cloudlet, i.e. BS and IoT gateway could run complete intelligent models/algorithms to provide services for edge devices. In area (3), a model is divided into several parts with different functions, which are performed by several edge devices. These edge devices work together to finish the task.

It is known that three most important elements for an intelligent application are: data, model, and computation. Suppose that an intelligent application is a ‘human’, model would be the ‘body’, and computation is the ‘heart’ which powers the ‘body’. Data is then the ‘book’. The ‘human’ improves their abilities by learning knowledge extracted from the ‘book’. After learning, the ‘human’ starts to work with the learned knowledge. Correspondingly, the complete deployment of most intelligent applications (unsupervised learning based application is not included) includes three components: data

collection and management (preparing the ‘book’), training (learning), and inference (working). Computation is a hidden component that is essential for the other three components. Combined with an edge environment, these three obvious components turn into edge cache (data collection and storage at edge), edge training (training at edge), and edge inference (inference at edge), respectively. Note that edge devices and edge servers are usually not powerful. Computation at edge usually is done via offloading. Hence, the hidden component turn into edge offloading (computation at edge). Our classification is organised around these four components, each of which features multidimensional analysis and discussion. The global outline of our proposed classification is shown in Fig. 2. For each component, we identify key problems in practical implementation and further break down these problems into multiple specific issues to outline a tiered classification. Next, we present an overview of these modules shown as Fig. 2.

A. Edge Caching

In edge intelligence, edge caching refers to a distributed data system proximity to end users, which collects and stores the data generated by edge devices and surrounding environments, and the data received from the Internet to support intelligent applications for users at the edge. Fig. 3 presents the essential idea of edge caching. Data is distributed at the edge. For example, mobile users’ information generated by themselves is stored in their smartphones. Edge devices such as monitoring devices and sensors record the environmental information. Such data is stored at reasonable places and used for processing and analysis by intelligent algorithms to provide services for end users. For example, the video captured by cameras could be cached on vehicles for aided driving [39]. BS caches the data that users recently accessed from the Internet

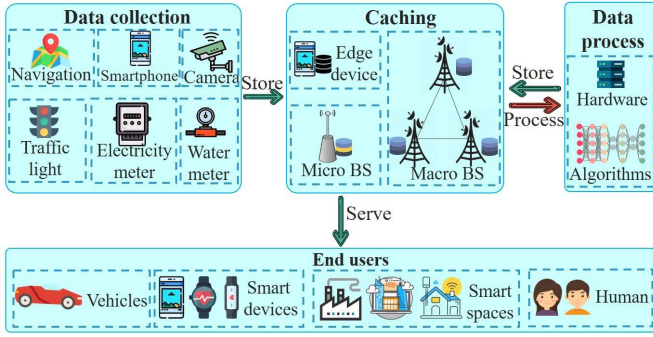


Fig. 3. The illustration of edge caching. Data generated by mobile users and collected from surrounding environments is collected and stored on edge devices, micro BSs, and macro BSs. Such data is processed and analysed by intelligent algorithms to provide services for end users.

to characterise users' interests for better a recommendation service [40]. To implement edge caching, we answer three questions: (i) what to cache, (ii) where to cache, and (iii) how to cache. The structure of this section is organised as the bottom module in Fig. 2.

For the first problem, what to cache, we know that caching is based on the redundancy of requests. In edge caching, the collected data is inputted into intelligent applications and results are sent back to where data is cached. Hence, there are two kinds of redundancy: data redundancy and computation redundancy. Data redundancy, also named communication redundancy, means that the inputs of an intelligent application may be the same or partially the same. For example, in continuous mobile vision analysis, there are large amounts of similar pixels between consecutive frames. Some resource-constrained edge devices need to upload collected videos to edge servers or the cloud for further processing. With cache, edge devices only need to upload different pixels or frames. For the repeated part, edge devices could reuse the results to avoid unnecessary computation. Ref. [41]–[44] have investigated the pattern of data redundancy. Caching based on such redundancy could effectively reduce computation and accelerate the inference. Computation redundancy means that the requested computing tasks of intelligent applications may be the same. For example, an edge server provides image recognition services for edge devices. The recognition tasks from the same context may be the same, e.g., the same tasks of flower recognition from different users of the same area. Edge servers could directly send the recognition results achieved previously back to users. Such kind of caching could significantly decrease computation and execution time. Some practical applications based on computation redundancy are developed in [45]–[47].

For the second problem, where to cache, existing works mainly focus on three places to deploy caches: macro BSs, micro BSs, and edge devices. The work in [48], [49] have discussed the advantages of caching at macro BSs from the perspective of covering range and hit probability. Some researchers also focus on the cached content at macro BSs. According to statistics, two kinds of content are considered: popular files [50]–[54] and intelligent models [55]–[59]. In

edge intelligence, macro BSs usually work as edge servers, which provide intelligent services with cached data. In addition, some works [14], [60]–[62] consider how to improve the performance of caching with the cooperation among macro BSs. Compared with macro BSs, micro BSs provide smaller coverage but higher quality of experience [63]–[69]. Existing efforts on this area mainly focus on two problems: how to deliver the cached content, and what to cache. For the aspect of delivery, research mainly focuses on two directions: delivery from single BS [70], and delivery from multiple BSs based on the cooperation amongst them [71]–[75]. Considering the small coverage of micro BSs and the mobility of mobile users, research on handover and users' mobility for better delivery service [76]–[79] is also carried out. In addition, the optimal content to cache, i.e., data redundancy based content [80]–[87] and computation redundancy based content [45], [88]–[94] is thoroughly investigated. Edge devices are usually of limited resources and high mobility, compared with macro BSs and micro BSs. Therefore, only few efforts pay attention to the problem of caching on a single edge device. For example, [39], [39], [95]–[97] studies the problem of what to cache based on communication and computation redundancy in some specific applications, e.g., computer vision. Most researchers adopt collaborative caching amongst edge devices, especially in the network with dense users. They usually formulate the caching problem into an optimisation problem on the content replacement [98]–[107], association policy [76], [108]–[110], [110]–[113], and incentive mechanisms [114], [115].

Since the storage capacity of macro BSs, micro BSs, and edge devices is limited, the content replacement must be considered. Works on this problem focus on designing replacement policies to maximise the service quality, such as popularity based schemes [116], [117], and ML based schemes [117], [118].

B. Edge Training

Edge training refers to a distributed learning procedure that learns the optimal values for all the weights and bias, or the hidden patterns based on the training set cached at the edge. For example, Google develops an intelligent input application, named G-board, which learns user's input habits with the user's input history and provides more precise prediction on the user's next input [33]. The architecture of edge training is shown as Fig. 4. Different from traditional centralised training procedures on powerful servers or computing clusters, edge training usually occurs on edge servers or edge devices, which are usually not as powerful as centralised servers or computing clusters. Hence, in addition to the problem of training set (caching), four key problems should be considered for edge training: (i) how to train (the training architecture), (ii) how to make the training faster (acceleration), (iii) how to optimise the training procedure (optimisation), and (iv) how to estimate the uncertainty of the model output (uncertainty estimates). The structure of this section is organised as the left module in Fig. 2.

For the first problem, researchers design two training architectures: solo training [30]–[32], [119] and collaborative

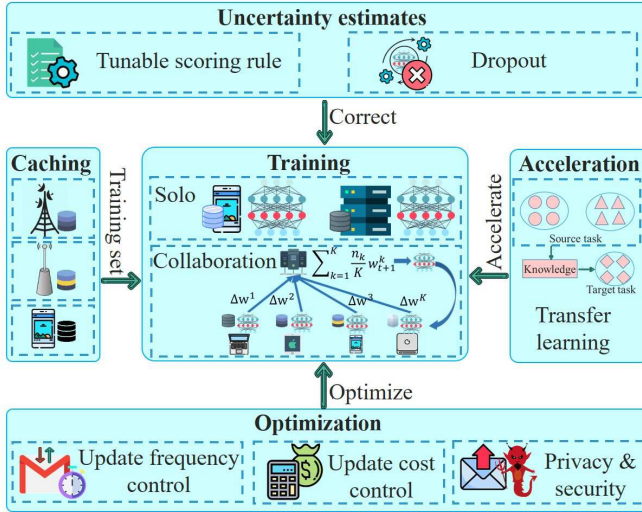


Fig. 4. The illustration of edge training. The model/algorithm is trained either on a single device (solo training), or by the collaboration of edge devices (collaborative training) with training sets cached at the edge. Acceleration module speeds up the training, whilst the optimisation module solves problems in training, e.g., update frequency, update cost, and privacy and security issues. Uncertainty estimates module controls the uncertainty in training.

training [33], [120]–[123]. Solo training means training tasks are performed on a single device, without assistance from others, whilst collaborative training means that multiple devices cooperate to train a common model/algorithm. Since solo training has a higher requirement on the hardware, which is usually unavailable, most existing literature focuses on collaborative training architectures.

Different from centralised training paradigms, in which powerful CPUs and GPUs could guarantee a good result with a limited training time, edge training is much slower. Some researchers pay attention to the acceleration of edge training. Corresponding to training architecture, works on training acceleration are divided into two categories: acceleration for solo training [119], [124]–[129], and collaborative training [123], [130], [131].

Solo training is a closed system, in which only iterative computation on single devices is needed to get the optimal parameters or patterns. In contrast, collaborative training is based on the cooperation of multiple devices, which requires periodic communication for updating. Update frequency and update cost are two factors which affect the performance of communication efficiency and training result. Researchers on this area mainly focus on how to maintain the performance of the model/algorithm with lower update frequency [55], [55], [132], [132]–[143], and update cost [55], [130], [144]–[147]. In addition, the public nature of collaborative training is vulnerable to malicious users. There is also some literature which focuses on the privacy [133], [148]–[158] and security [159]–[163], [163]–[165], [165]–[168] issues.

In DL training, the output results may be erroneously interpreted as model confidence. Estimating uncertainty is easy on traditional intelligence, whilst it is resource-consuming for edge training. Some literature [169], [170] pays attention to this problem and proposes various kinds of solutions to reduce

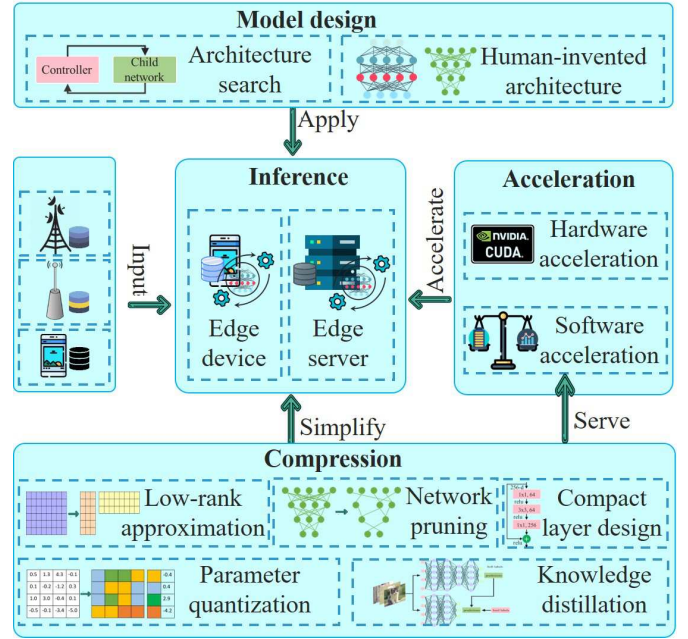


Fig. 5. The illustration of edge inference. AI models/algorithms are designed either by machines or humans. Models could be further compressed through compression technologies: low-rank approximation, network pruning, compact layer design, parameter quantisation, and knowledge distillation. Hardware and software solutions are used to accelerate the inference with input data.

computation and energy consumption.

We also summarised some typical applications of edge training [33], [55], [132], [133], [171]–[176], [176]–[180] that adopt the above-mentioned solutions and approaches.

C. Edge Inference

Edge inference is the stage where a trained model/algorithm is used to infer the testing instance by a forward pass to compute the output on edge devices and servers. For example, developers have designed a face verification application based DL, and employ on-device inference [181], [182], which achieves high accuracy and low computation cost. The architecture of edge inference is shown as Fig. 5. Most existing AI models are designed to be implemented on devices which have powerful CPUs and GPUs, this is not applicable in an edge environment. Hence, the critical problems of employing edge inference are: (i) how to make models applicable for their deployment on edge devices or servers (design new models, or compress existing models), and (ii) how to accelerate edge inference to provide real-time responses. The structure of this section is organised as the right module in Fig. 2.

For the problem of how to make models applicable for the edge environment, researchers mainly focus on two research directions: design new models/algorithms that have less requirements on the hardware, naturally suitable for edge environments, and compress existing models to reduce unnecessary operation during inference. For the first direction, there are two ways to design new models: let machines themselves design optimal models, i.e., architecture search [183], [184], [184], [184]–[187]; and human-invented architectures with

the application of depth-wise separable convolution [188]–[190] and group convolution [191], [192]. We also summarise some typical applications based on these architectures, including face recognition [181], [182], [193], human activity recognition (HAR) [194]–[202], vehicle driving [203]–[206], and audio sensing [207], [208]. For the second direction, i.e., model compression, researchers focus on compressing existing models to obtain thinner and smaller models, which are more computation- and energy-efficient with negligible or even no loss on accuracy. There are five commonly used approaches on model compression: low-rank approximation [209]–[214], knowledge distillation [215]–[223], compact layer design [224]–[232], network pruning [233]–[247], and parameter quantisation [210], [247]–[263]. In addition, we also summarise some typical applications [264]–[269] that are based on model compression.

Similar to edge training, edge devices and servers are not as powerful as centralised servers or computing clusters. Hence, edge inference is much slower. Some literature focuses on solving this problem by accelerating edge inference. There are two commonly used acceleration approaches: hardware acceleration and software acceleration. Literature on hardware acceleration [96], [270]–[279], [279]–[294] mainly focuses on the parallel computing which is available as hardware on devices, e.g., CPU, GPU, and DSP. Literature on software acceleration [39], [95], [96], [295]–[302] focus on optimising resource management, pipeline design, and compilers, based on compressed models.

D. Edge offloading

As a necessary component of edge intelligence, edge offloading refers to a distributed computing paradigm, which provides computing service for edge caching, edge training, and edge inference. If a single edge device does not have enough resource for a specific edge intelligence application, it could offload application tasks to edge servers or other edge devices. The architecture of edge offloading is shown as Fig. 6. Edge offloading layer transparently provides computing services for the other three components of edge intelligence. In edge offloading, Offloading strategy is of utmost importance, which should give full play to the available resources in edge environment. The structure of this section is organised as the top module in Fig. 2.

Available computing resources are distributed in cloud servers, edge servers, and edge devices. Correspondingly, existing literature mainly focuses on four strategies: device-to-cloud (D2C) offloading, device-to-edge server (D2E) offloading, device-to-device (D2D) offloading, and hybrid offloading. Works on the D2C offloading strategy [303]–[318] prefer to leave pre-processing tasks on edge devices and offload the rest of the tasks to a cloud server, which could significantly reduce the amount of uploaded data and latency. Works on D2E offloading strategy [19], [319]–[324] also adopt such operation, which could further reduce latency and the dependency on cellular network. Most works on D2D offloading strategy [325]–[333] focus on smart home scenarios, where IoT devices, smartwatches and smartphones collaboratively

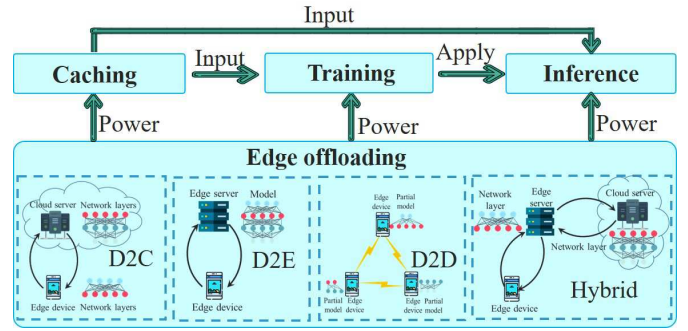


Fig. 6. The illustration of edge offloading. Edge offloading is located at the bottom layer in edge intelligence, which provides computing services for edge caching, edge training, and edge inference. The computing architecture includes D2C, D2E, D2D, and hybrid computing.

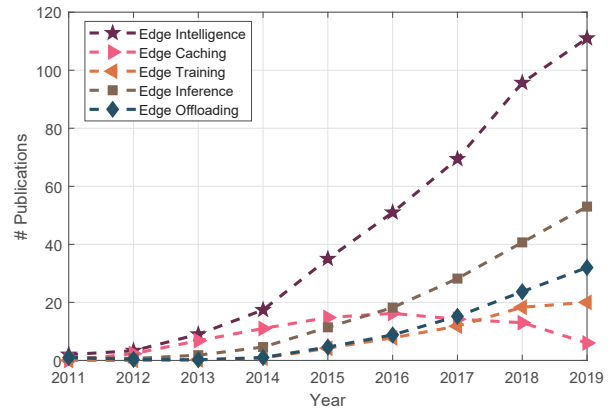


Fig. 7. Publication volume over time. These curves show the trend of publication volume in edge caching, edge training, edge computing, edge inference, and edge intelligence, respectively.

perform training/inference tasks. Hybrid offloading schemes [334]–[336] have the strongest ability of adaptiveness, which makes the most of all the available resources.

We also summarise some typical applications that are based on these offloading strategies, including intelligent transportation [337], smart industry [338], smart city [339], and healthcare [340] [341].

E. Summary

In our survey, we identify four key components of edge intelligence, i.e. edge caching, edge training, edge inference, and edge offloading. Edge intelligence shows an explosive developing trend with a huge amount of researcher have been carried out to investigate and realise edge intelligence over the past five years. We count the publication volume of edge intelligence, as shown in Fig. 7.

We see that this research area started from 2011, then grew at a slow pace before reaching 2014. Most strikingly, after 2014, there is a rapid rise in the publication volume of edge training, edge inference, and edge offloading. Meanwhile, the publication of edge caching is gradually winding down. Overall, the publication volume of edge intelligence is booming,

which demonstrates a research field replete with activity. Such prosperity of this research field owes to the following three reasons.

First, it is the booming development of intelligent techniques, e.g., deep learning and machine learning techniques that provides a theoretical foundation for the implementation of edge intelligence [342]–[344]. Intelligent techniques achieve state-of-the-art performance on various fields, ranging from voice recognition, behaviour prediction, to automatic piloting. Benefiting from these achievements, our life has been dramatically changed. People hope to enjoy smart services anywhere and at any time. Meanwhile, most existing intelligent services are based on cloud computing, which brings inconvenience for users. For example, more and more people are using voice assistant on smartphone, e.g., MI AI and Apple Siri. However, such applications can not work without networks.

Second, the increasing big data distributed at the edge, which fuels the performance of edge intelligence [345]–[347]. We have entered the era of IoT, where a giant amount of IoT devices collect sensory data from surrounding environment day and night and provide various kinds services for users. Uploading such giant amount of data to cloud data centres would consume significant bandwidth resources. Meanwhile, more and more people are concerned about privacy and security issues behind the uploaded data. Pushing intelligent frontiers is a promising solution to solve these problems and unleash the potential of big data at the edge.

Third, the maturing of edge computing systems [14], [16] and peoples' increasing demand on smart life [348], [349] facilitate the implementation of edge intelligence. Over the past few years, the theories of edge computing have moved towards application, and various kinds of applications have been developed to improve our life, e.g., augmented reality [350]–[352]. At the same time, with the wide spreading of 5G networks, more and more IoT devices are implemented to construct smart cities. People are increasingly reliant on the convenient service provided from a smart life. Large efforts from both academia and industry are enacted to realise these demands.

III. EDGE CACHING

Initially, the concept of caching comes from computer systems. The cache was designed to fill the throughput gap between the main memory and registers [353] by exploring correlations of memory access patterns. Later, the caching idea was introduced in networks to fill the throughput gap between core networks and access networks. Nowadays, the cache is deployed in edge devices, like various base stations and end devices. By leveraging the spatiotemporal redundancy of communication and computation tasks, caching at the edge can significantly reduce transmission and computation latency and improve users' QoE [354]–[356].

From existing studies, the critical issues of caching technologies in edge networks fall into three aspects: the cached content, caching places, and caching strategies. Next, we discuss and analyse relevant literature in edge caching in

terms of the above three perspectives. The related subjects include the preliminary of caching, cache deployment, and cache replacement.

A. Preliminary of Caching

The critical idea of edge caching technologies is to explore the spatiotemporal redundancy of users' requests. The redundancy factor largely determines the feasibility of caching techniques. Generally, there are two categories, i.e., communication redundancy, and computation redundancy.

1) *Communication Redundancy*: Communication redundancy is caused by the repetitive access of popular multimedia files, such as audio, video, and webpages. Content with high popularity tends to be requested by mobile users many times. Thus, the network needs to transmit the content to these users over and over again. In this case, caching popular content at edge devices can eliminate enormous duplicated transmissions from core networks.

To better understand communication redundancy, many existing studies investigate content request patterns of mobile users. For example, Crane *et al.* [41] analyse 5 million videos on YouTube and regards the number of daily views as the popularity of videos. They then discover four temporal popularity patterns of videos and demonstrate the temporal communication redundancy of popular videos from the angle of content providers'. Meanwhile, Adar *et al.* [42] analyse five weeks of webpage interaction logs collected from over 612,000 users. They show that temporal revisitation also exists at the individual level. In [43], Traverso *et al.* characterise the popularity profiles of YouTube videos by using the access logs of 60,000 end-users. They propose a realistic arrival process model, i.e., the Shot Noise Model (SNM), to model the temporal revisitation of online videos. Moreover, Dernbach *et al.* [44] exhibit the existence of regional movie preferences, i.e., the spatial communication redundancy of content by analysing the MovieLens dataset, which contains 6,000 user viewing logs. Consequently, the above studies apply large-scale and real-world datasets demonstrating communication redundancy from both the temporal and spatial dimensions. These studies support the feasibility of edge caching ideas and pave the way for employing edge caching technologies in real-world scenarios.

2) *Computation Redundancy*: Computation redundancy is caused by commonly used high computational complexity applications or AI models. In the wave of AI, we are now surrounded by various intelligent edge devices such as smartphones, smart watches, and smart brands. These intelligent edge devices provide diverse applications to augment users' understanding of their surroundings [357], [358]. For example, speech-recognition based AI assistants, e.g., Siri and Cortana, and song identification enabled music applications, have been widely used in peoples' daily lives. Such AI-based applications are of high computational complexity and cause high power consumption of the device [350], [351], [359].

Meanwhile, some researchers have discovered the computation redundancy in AI-based applications. For example, in a park, nearby visitors may use their smartphones to recognise

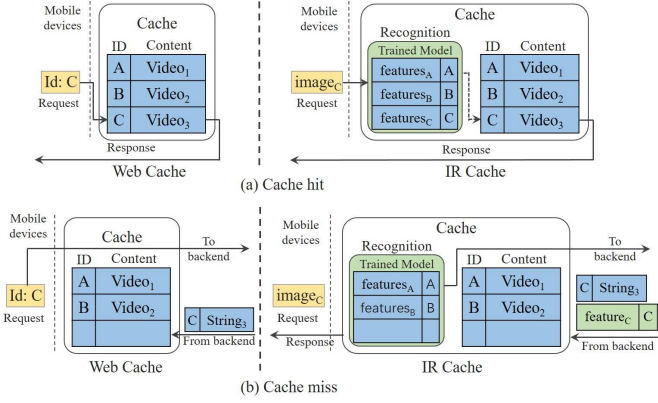


Fig. 8. Illustration of image recognition (IR) cache, as well as the comparison with traditional web cache. The request/input for IR is an image. The system run the trained model to recognize the image, which will be labelled with an identifier. Then, the recognized image’s identifier is used to find relevant content within cache. If there is a cache hit, content is returned. Otherwise, IR modelling is performed and the result is placed in the cache.

the same flowers and then search the information accordingly. In this case, there are a lot of unnecessary computations across devices. Therefore, if we offload such a painting recognition task to edges and cache the computation results, redundant computations can be further eliminated [45], [360]. In [45], Guo *et al.* crawl street views by using the Google Streetview API [46] and builds an ‘outdoor input image set’. They then find that around 83% of the images exhibit redundancy, which leads to a large number of potential unnecessary computations for image reorganisation application. Also, they analyse NAVVIS [47], an indoor view dataset, and observed that nearly 64% of indoor images exhibit redundancy. To our knowledge, this work is the earliest one using real-world datasets to demonstrate the existence of computation redundancy.

In the elimination of computation redundancy, an important step is to capture and quantify the similarity of users’ requests. As shown in Fig. 8, in the case of communication redundancy, a unique identifier can identify users’ request content, e.g., Universal Resource Identifier, URI. However, for computation redundancy, we first need to obtain the features of users’ requests and then find the best match of the computation results according to the extracted features. It is notable that in computation redundancy, the cached content is computation results instead of requested files.

B. Cache Deployment

As shown in Fig. 9, in edge networks, there are three main places to deploy cache units, i.e., macro base stations, small base stations, and end devices. Caching at different places have different characteristics, and we now provide a detailed discussion.

1) *Caching at Macro Base Stations*: The main purpose of deploying cache units in macro base stations (MBSs) is to relax the burden of backhaul [48] by exploring the communication redundancy and cache machine learning models by reducing the computation redundancy. Caching popular files and models in MBSs, content could be directly fetched from MBSs instead of core networks. In this way, redundant

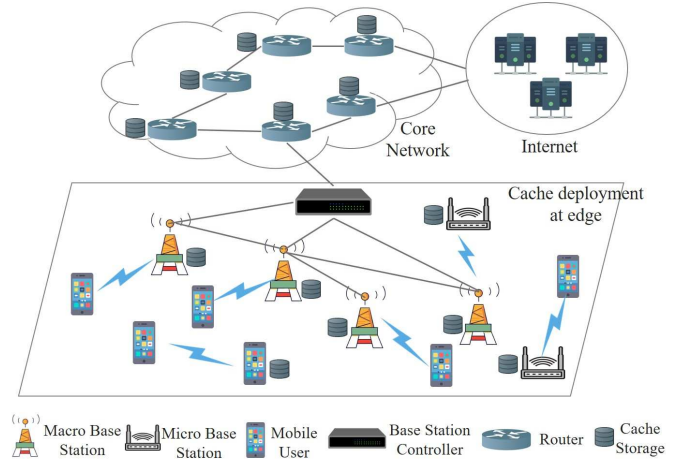


Fig. 9. Cache deployment at the edge. There are three places to deploy caches: macro base stations, micro base stations, and end devices.

transmission and computation are eliminated. Compared with other cache places in the edge networks, MBSs have the most extensive coverage range and the most massive cache spaces. The typical coverage radius of an MBS is nearly 500 meters [49]. Due to its broad coverage range and vast cache spaces, MBS can serve more users. Thus, caching at MBSs can explore more of both communication redundancy and computation redundancy, and obtain a better cache hit probability. Also, since MBSs are deployed by operators, the topology structure of MBSs is stable and does not change over time. The comparison of different cache places is summarised in Table II.

TABLE II
COMPARISON OF DIFFERENT CACHE PLACES.

Cache places	MBSs	SBSs	Devices
Coverage radius	500m	20 ~ 200m	10m
Cache spaces	Large	Medium	Small
Served users	Massive	Small	Few
Topology structure	Stable	Change slightly	Change dramatically
Redundancy potential	High	Medium	Low
Computational power	High	Medium	Low

Some researchers first studied the most fundamental problem, i.e., what files should be cached at MBSs to improve end-users’ QoE. One straightforward idea is to cache the most popular files. However, in [50], Blaszczyzyn *et al.* find that always caching the most popular contents maybe not the optimal strategy. Specifically, they assume the distribution of MBSs follows a Poisson Poisson point process [51], and both users’ requests and the popularity of content follow Zipf’s law [52]. By maximising users’ cache hit ratio, they derived the optimal scheme, in which less popular contents are also cached. Furthermore, Ahlehagh and Dey [53] take user preference profiles into account, and the videos of user-preferred video categories have a high priority for caching. Alternatively, Chatzieftheriou *et al.* [54] explore the effect of content recommendations on the MBSs caching system. They discover that caching the recommended content can improve the cache hit ratio. Still, the recommendations of

service providers may distort user preferences.

Apart from popular files, caching machine learning models at MBSs is another promising field. In 2016, Google proposed a novel machine learning framework, i.e., federated learning with the objective of training a global model, e.g., deep neural networks, across multiple devices [55]. In detail, the multiple devices train local models on local data samples. And they only exchange parameters, e.g., the weights of a deep neural network, between these local models and finally converge to a global model. Note that data is only kept in local devices and not exchanged across devices. Thus, federated learning can ensure data privacy and data secrecy and attract widespread attention from both industry [56] and academia [57]. Because MBSs can serve more mobile devices and have more powerful computation units, they are usually acted as the central server to orchestrate the different steps of the federated learning algorithm by caching the global model and collecting parameters from multiple devices [58], [59].

To fully explore the communication and computation redundancy, MBSs are allowed to cooperatively cache content, including files, computation results, and AI-models. In other words, a mobile user can be neighbouring base stations [60]. In [61], Peng *et al.* consider a collaborative caching network where all base stations are managed by a controller. They discover that contents with the highest popularity should be cached first in the case of long latency through backhaul network. Otherwise, caches should keep content diversity, i.e., to cache as much different content as possible. Also, Tuyen *et al.* [14] propose a collaborative computing framework across multiple MBSs. The computation workload at different base stations usually exhibits spatial diversity [62]. Therefore, offloading computation tasks to nearby idle MBSs and cache the computation results or trained-AI models cooperatively can facilitate the performance of mobile networks.

2) *Caching at Small Base Stations*: Small base stations (SBSs) are a set of low-powered mobile access points that have a range of 20 metres to 200 metres, e.g., microcell, picocell, and femtocell [63]. By deploying small base stations on hot spots, mobile users will have a better quality of experience, such as high end-rate, due to the benefit from spatial spectrum reuse [64], [361]. Therefore, densely deploying small base stations [65], [66] is a promising approach in future mobile networks and also brings huge potential to reduce both communication and computation redundancy by caching at SBSs [67]–[69].

In [362], Bacstug *et al.* use a stochastic geometry method to model the caching network of small cells, where users' most popular files or computation results are stored on SBSs. They theoretically demonstrate that employing storage units in SBSs indeed brings gains in terms of the average delivery rate. Unlike [362] which considers a centralised control method, Chen *et al.* [363] propose a distributed caching strategy where each SBS only considers local users' request pattern instead of global popularity. Each SBS maintains a content list by sampling requested files.

Compared with MBS, SBS has less redundancy potential due to its smaller coverage range, the number of served users, and even cache spaces. Thus, various technologies are adopted

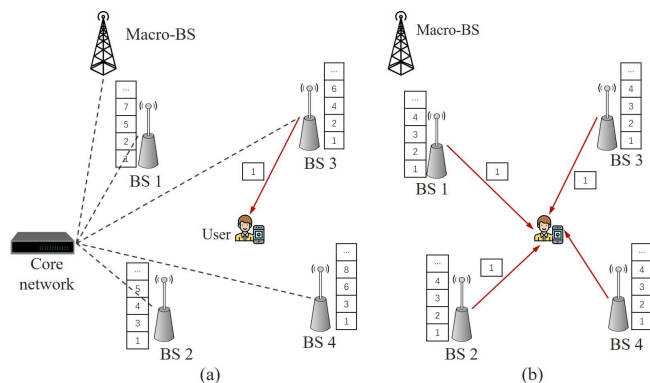


Fig. 10. Cache model at SBS. (a) mobile use is associated with one SBS. (b) mobile use is associated with multiple SBS, which send content to the user with beamforming and CoMP.

to explore caching benefits fully. A commonly used technique is the cooperative transmission. In SBS networks, the coverage area of SBSs is usually overlapped with each other, especially for the dense deployment scenario. In other words, a mobile user is able to receive content from multiple SBSs, making cooperative caching of multiple SBSs possible. There are a lot of transmission technologies applied in SBS networks, such as multicast, beamforming, cooperative multi-point (CoMP), and so on.

In [364], Liao *et al.* investigate caching enabled SBS networks. By exploring the potential of multicast transmission from SBSs to mobile users, their approach could reduce the backhaul cost in SBS networks. Similarly, Poularakis *et al.* [365] also delve into multicast opportunities in cache-enabled SBS networks. Different from [364], they assume that both MBS and SBSs are able to use multicast. Each SBS can create multicast transmissions to end-users, while each MBS can provide multicast of popular content to SBSs within the coverage area. In terms of simulation results, the serving cost reduction up to 88% compared to unicast transmissions.

In SBS networks, since a mobile user has the potential to receive signals from multiple SBSs, there are two association cases, as shown in Fig.10. In the first case, the mobile user is only associated with one SBS. Alternatively, the mobile user is associated with multiple SBSs. By using the beamforming [366] and CoMP technology [367], the associated SBSs can jointly send content for downlink transmission.

In [70], Pantisano *et al.* consider the first association case and presents a cache-aware user association approach. They adjust users' association to improve the local cache hit ratio based on whether the associated base station contains files and AI models required by the user. The problem is formulated as a one-to-many matching game and they propose a distributed algorithm based on the deferred acceptance scheme to solve it.

Alternatively, some scholars focus on the second association case and explore the power of cooperative transmission in cache-enabled SBS networks [71]–[75]. In [71], Shanmugam *et al.* study the problem of content deployment in SBS networks. They assume that mobile users could communicate with multiple cache units and formulate the optimisation prob-

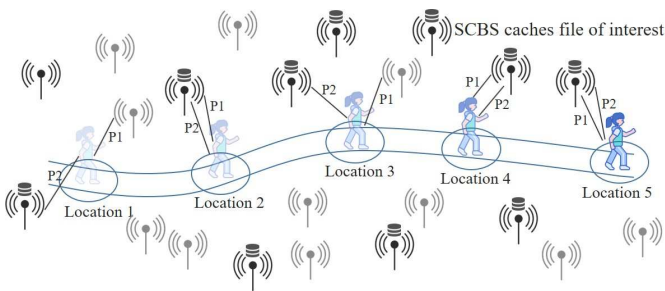


Fig. 11. User gets file of interest from SCBS using policy 1, denoted as P1, and policy 2, denoted as P2, while moving from location 1 to location 5. P1 means connecting to the SCBS that provides the highest average received power. P2 means connecting to the nearest SCBS that could provide files of interest.

lem to minimise the average downloading delay. Liu *et al.* [72] explore the potential of energy efficiency of cache-enabled SBS networks where all SBSs use CoMP to transmit cached content cooperatively. By maximising energy efficiency, the optimal transmit power of SBSs are worked out. In [73] and [74], Ao *et al.* study a distributed caching strategy in SBS networks where CoMP technology is applied. In the CoMP enabled networks, caching strategy can bring two different gains. On the one hand, diverse contents could be cached in nearby BSs to maximise the cache hit. On the other hand, caching the same content in nearby BSs can let corresponding BSs transmit concurrently and bring multiplexing gain. By trading off both gains, they devised a near-optimal strategy to maximise the system throughput. Moreover, they find when content is of skewed popularity distribution, caching multiple copies of popular files and AI models yields larger caching gains. Further, Chen *et al.* [75] consider a similar system. Unlike [73] and [74], where all nearby SBSs can employ CoMP, Chen *et al.* first group SBSs into multiple disjointed clusters and only the SBSs in the same cluster are able to transmit content cooperatively. To trade off the parallel transmission and joint transmission, they divide the cache space into two parts. One is in charge of caching content with less popularity to improve content diversity, while the other is used to cache contents with the highest popularity. Then, they optimize the problem of space assignment.

Since the coverage range of each SBS is too small, mobile users will go through multiple SBSs within a short time, as shown in Fig. 11. This frequent handover behaviour will cause the degradation of caching performance. In [76], Krishnan *et al.* investigate the retransmission in cache-enabled SBS networks. Because of the frequent handover between different SBSs, sometimes, when none of the SBSs in the user's vicinity has cached the requested file or AI models, the file transmission will be interrupted. Nevertheless, the retransmission will be triggered when the requested file or AI models are cached at vicinity SBSs. By using stochastic geometry to analyse the cache hit probability, Krishnan *et al.* find that SBSs should cache content diversely for mobile users. In [77], Guan *et al.* assume that users' preferences for content and mobility patterns are known prior, and users' preferences remain constant over a short period. They then formulate an

optimisation problem with the objective of maximising the utility of caching and devise a heuristic caching strategy. In [78] and [79], the same caching system model is investigated where mobile users migrate between multiple SBSs. Due to the limited transmission time, users may not be able to download the complete requested files or parameters of AI models from the associated SBS, and the requests can be redirected to MBS. In [78], Poularakis *et al.* use random walks to model user movements and formulated an optimisation problem based on the Markov chain aiming to maximize the probability of serving by SBSs. They further propose two caching strategies, i.e., a centralised solution for the small-scale system and a distributed solution for the large-scale system. Unlike [78], Ozfatura *et al.* propose a distributed greedy algorithm to minimise the amount of data downloaded from MBS [79]. Requests with deadlines below a given threshold are responded by SBSs while other requests are served by MBS.

In SBS caching systems, predicting users' requests to improve the cache hit ratio is also a commonly used method, which falls into the field of artificial intelligence applications. By applying AI technology to historical users' request logs, we can profile user preference or content popularity patterns. Next, we can predict users' requests or content popularity, respectively [80]. In [81], Kader *et al.* design a big data platform and collects mobile traffic data from a telecom operator in Turkey. They then used collaborative filtering, a common machine learning method, to estimate content popularity. The simulation results demonstrate that the caching benefits are further explored with the help of content popularity prediction. Similar to [81], in [82], Pantisano *et al.* also apply collaborative filtering to predict content popularity. They then devised a user-SBS association scheme based on estimated popularity and the current cache composition to minimise the backhaul bandwidth allocation. In [83], Bastug *et al.* focus on individual content request probability instead of global content request probability. They propose to use the Bayesian learning method to predict personal preferences and then incorporate this crucial information into the caching strategy. If we lack the historical data of user request logs, how can we predict content popularity? In [84], Bastug *et al.* investigate this open issue and proposes a transfer learning-based caching procedure. Specifically, they exploit contextual information, e.g., social networks, and referred to it as a source domain. Then the prior information in the source domain is incorporated in the target domain to estimate content popularity.

Also, SBSs are used to cache the data from end devices, like smartphones and IoT. In recent times, IoT devices are widely distributed in homes, streets, and even whole cities to allow users to monitor the ambient environment [85], [352]. By collecting and analysing the big data on IoT devices, a smarter physical world can be built. Considering the demand of real-time data analysis, caching and processing the data at the edge is a common and promising method. In [86], Quevedo *et al.* introduce a caching system for IoT data and proof that the caching system could reduce the energy consumption of IoT sensors. In [87], Sharma *et al.* propose a collaborative edge and cloud data processing framework for IoT networks where SBSs are in charge of caching IoT data, extracting

useful features and uploading features to cloud part.

Meanwhile, since SBSs are often deployed at hot points, the requested computation tasks from served users will exhibit spatiotemporal locality. Therefore, by caching computation results at SBSs, the redundant computation tasks can be eliminated. Drolia *et al.* [88] propose a caching strategy, Cachier, to cache the recognition results on edge servers by release repetitive recognition computation. Specifically, Cachier first extracts features of the requested recognition task, and then tries to match a similar object from the cache. If there is a cache hit, the corresponding computation results would be sent back to the mobile device. Otherwise, the request would be sent to the cloud. To identify similar recognition tasks, they used a Locality Sensitive Hashing (LSH) algorithm [89] to determine the best match. Furthermore, to overcome the unbalanced and time-varying distribution of users' requested tasks, Guo *et al.* [45] design an Adaptive LSH-Homogenized kNN joint algorithm which outperforms LSH in terms of evaluation results. Drolia *et al.* further introduce a proactive caching strategy into their system by predicting the requirement of users and proactively caching parts of models on SBSs server for pre-processing to further reduce the latency [90]. Such a strategy is also used in [91] to deal with unstructured data at SBSs.

Moreover, in some task-fickle scenarios, multiple different kinds of tasks, e.g., voice recognition and object recognition, are offloaded from devices to SBSs. By pre-caching multiple kinds of deep learning models at SBSs for different kinds of tasks, we can reduce the computation time and further improve users' QoE. Taylor *et al.* propose an adaptive model selection scheme to select the best model for users [92]. They use a supervised learning method to train a predictor model offline and then deploy it on an edge server. When a request arrives, the predictor will select an optimal model for the task. In [93], Zhao *et al.* propose a system, Zoo, to compose different models to provide a satisfactory service for users. Ogden *et al.* propose a deep inference platform, MODI, to determine what model to cache and what model to use for specific tasks [94]. There is a decision engine inside MODI, which aggregates previous results to decide what new models are required to cache.

3) *Caching at Devices*: Caching at devices exploits the available storage space of end equipment, like mobile phones and IoT devices. These devices can leverage the communication and computation redundancy locally. Furthermore, they can fetch the requested content or computation results from other devices in proximity through device-to-device (D2D) communication [368], [369].

First of all, end devices can explore the communication and computation redundancy locally. For instance, in some static continuous computer vision applications, such as monitoring, the captured consecutive images are similar to some extent. Therefore, the results of previous images could be reused for the latter inference. In [39], [95], [96], they cache the results of previous frames to reduce redundant computation and latency. In some mobile continuous computer vision applications, such as driving assistance, the system is required to provide high trackability. The system needs to recognise, locate, and label

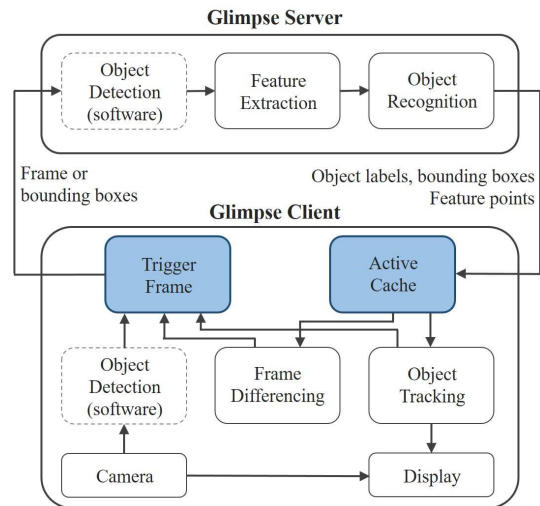


Fig. 12. The architecture of Glimpse. Edge device, i.e., glimpse client, only uploads trigger frames to the cloud to save bandwidth resources. Glimpse server transmits the recognition results and features back to edge devices. Edge devices deal with local frames with these features.

the tracked object, e.g., road signs, on the screen in real-time. The recognised object would repeatedly appear in multiple images for a period. Chen *et al.* develop an active cache based continuous object recognition system, called Glimpse, to achieve real-time trackability [97]. The structure of Glimpse is shown as Fig. 12. Glimpse caches frames locally and only uploads trigger frames to the cloud server. Trigger frames refers to the frames, for which the recognition from the server is different from current local tracking. The cloud server sends back the recognised object, its labels, bounding boxes, and features, which would be cached locally on devices. Then the devices would track the object with the labels, bounding boxes, and features locally on captured frames. A similar approach is also adopted in CNNCache [39].

On the other hand, compared with MBSs and SBSs, devices have very limited cache spaces and coverage range, due to the cost constraint of end devices and their low transmission power. Although these limitations seem to only take the small caching benefits for local devices, the situation will be changed when the networks are with dense users. The benefits of caching will be amplified with the number of users increases. In [370], Chen *et al.* study the difference between caching at SBSs and devices where content is cached according to a joint probability distribution. By applying stochastic geometry, they derive the closed-form expression of hit probability and request density. Although the cache hit probability of device caching is always lower than SBS caching due to the small cache spaces, the request density of device caching is much higher than SBS caching. This is because, in device caching, more concurrent links are allowed, compared with the case of SBS caching, especially in the dense user scenario. Similarly, in [371], Gregori *et al.* investigate caching at both devices and SBSs as well. However, they do not compare these two different scenarios and only design joint transmission and caching policies for them to minimise energy consumption

separately.

In the device caching system, the number of coexisting D2D links affect the device caching performance dramatically. The D2D links are the fundamental requirement for end-devices sharing files, models, and computation results and further reducing communication and computation redundancy amongst devices. First of all, the establishment of a D2D connection depends on content placement. In other words, when a device discovers that the requested content is placed in a nearby device within the D2D transmission range, the D2D link can be built, and the content will be transmitted directly. Therefore, there are a lot of scholars trying to maximise caching performance by optimising the content placement [98]–[102]. In [101], Malak *et al.* model senders and receivers as members of Poisson Point Process and compute the probability of delivery in D2D networks. Considering the low transmission noise case, they find that the optimal content allocation could be approximately achieved by Benford’s law when the path loss exponent equals 4. A similar system model is applied in [102] but with a different performance metric. In [102], Peng *et al.* analyse the outage probability of D2D transmission in cache-enabled D2D networks. They then obtain the optimal caching strategy by using a gradient descent method. In [98], Chen *et al.* aim to maximise successful offloading probability. Different from [101], [102] which did not consider the time divided transmission, Chen *et al.* divide time into multiple slots, and each transmitter independently chooses a time slot to transmit files. Employing the gradient descent method, they design a distributed caching policy. Unlike the above studies where each end-user applies the same caching strategy, Giatsoglou *et al.* [99] divide the $2K$ most popular contents into two groups of the same size. K is the cache capacity of each user. Then randomly allocate these two groups to users, i.e., some users cache group A, whilst others cache group B.

Apart from content placement, association policy is also important to the establishment of D2D links. In [108], Golrezaei *et al.* optimise the collaboration distance for D2D communications with distributed caching, where the collaboration distance is the maximum allowable distance for D2D communication. They assume each user employs a simple randomised caching policy. In [109], Naderializadeh *et al.* propose a greedy association method, i.e., greedy closest-source policy. In this association policy, starting from the first user, each user chooses the closest user with the desired file forming a D2D pair. They assume that each file is randomly cached in devices and derive a lower bound on the cached content reuse.

Generally, end devices are controlled by end consumers who are able to decide whether to cache and share content or to not do so. Therefore, the incentive mechanism is introduced in D2D networks to encourage users to exploit the storage space of their equipment and share cached content, e.g., files, AI models and computation results, with other users. In [114], Chen *et al.* propose an incentive mechanism where the base station rewards the users who shares content with others via D2D communications. Since the base station will determine the reward to minimise its total cost while users would like to maximise their reward by choosing the caching

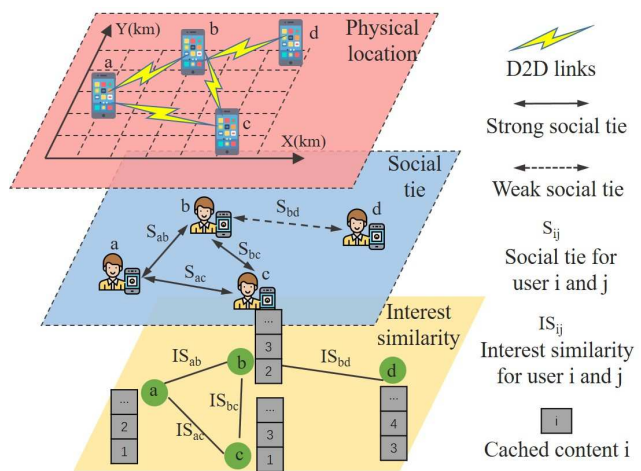


Fig. 13. Social-aware caching at edge devices. In location based framework, users close to each other could exchange cached content, e.g., user a, b, and c. In social tie based framework, user could also exchange cached content with others if they have strong social tie. In interest based framework, interest similarity is used to estimate the social tie among users.

policy, Chen *et al.* model this conflict as a Stackelberg game and proposed an iterative gradient algorithm to obtain the Stackelberg Equilibrium. In [115], Taghizadeh *et al.* consider a similar case where content providers pay the download cost to encourage users to download and share content. However, they do not model the conflict and merely design the caching strategy to minimise content provisioning costs.

Since end devices are bound up with users and affected by user attributes to some extent, some researchers focus on exploring the knowledge of user attributes, like social ties and interests, to assist device caching as shown in Fig. 13. In [111], Bastug *et al.* propose to let influential users cache content, such that these users could disseminate the cached contents to others through their social ties. The influential users are determined by their social networks. First, a social graph is built based on past action history of users’ encounters and file requests. Then, the influence of users is measured in terms of the centrality metric [112]. Apart from social ties, Bai *et al.* [113] consider users’ interests as well. They use the hypergraph to model the relationships among social ties, common interests, and spectrum resources and design an optimal caching strategy to maximise the cache hit ratio.

C. Cache Replacement

In practice, the request distribution of content varies with time, and new content is constantly being created. Hence, it is critical to update caches at intervals. The cache update process generally takes place when new content is delivered and needs to be cached, but all cache units are occupied. Hence, some cached old content needs to be replaced. Therefore, the cache update process is called cache replacement as well.

Several conventional cache replacement strategies have been proposed, such as first-in first-out (FIFO), least frequently used (LFU), least recently used (LRU), and their variants [381]. FIFO evicts the content in terms of cached time without any regard to how often or how many times it was accessed before. LFU keeps the most frequently requested content, while LRU

TABLE III
COMPARISON OF DIFFERENT CACHE DEPLOYMENT STRATEGIES.

Ref.	Cache places	Performance metrics	Mathematical tools	Control methods	Transmission Cooperativity
[50]	MBSs	Hit probability	Stochastic geometry	Centralised	Non-cooperative
[54]	MBSs	Cache hit ratio	Optimisation	Centralised	Non-cooperative
[372] [373]	MBSs	Energy efficiency	Stochastic geometry	Centralised	Non-cooperative
[53]	MBSs	The number of concurrent videos	Optimisation	Distributed	Non-cooperative
[61]	MBSs	The average download delay	Optimisation	Centralised	Cooperative
[374]	MBSs	Storage space	Optimisation	Centralised	Cooperative
[375]	MBSs	Aggregate operational cost	Optimisation	Centralised	Cooperative
[376]	MBSs	The aggregated caching and download cost	Optimisation	Centralised	Cooperative
[377]	MBSs	Cache failure probability	Optimisation	Centralised	Cooperative
[362]	SBSs	Outage probability Content delivery rate	Contract theory	Centralised	Non-cooperative
[363]	SBSs	The cache service probability	Stochastic geometry	Distributed	Non-cooperative
[378]	SBSs	The profit of NSP and VPs	Stochastic geometry	Centralised	Non-cooperative
[379]	SBSs	The number of requests served by SBSs	Optimisation	Centralised	Non-cooperative
[364]	SBSs	Backhaul cost	Optimisation	Centralised	Non-cooperative
[365]	SBSs	Servicing cost	Optimisation	Centralised	Non-cooperative
[70]	SBSs	Caching utility	Matching theory	Distributed	Non-cooperative
[71]	SBSs	Downloading time of files	Optimisation	Centralised	Cooperative
[72]	SBSs	Energy efficiency	Optimisation	Centralised	Cooperative
[73] [74]	SBSs	System throughput	Optimisation	Distributed	Cooperative
[75]	SBSs	Cache hit probability and energy efficiency	Stochastic geometry	Centralised	Cooperative
[76]	SBSs	Cache hit probability	Stochastic geometry	Centralised	Non-cooperative
[77]	SBSs	Maximise caching utility	Optimisation	Centralised	Non-cooperative
[78]	SBSs	The probability of response from MBS	Optimisation	Centralised & Distributed	Non-cooperative
[79]	SBSs	The amount of data downloaded from MBS	Optimisation	Distributed	Non-cooperative
[81]	SBSs	Backhaul load	Machine learning	Centralised	Non-cooperative
[82]	SBSs	Backhaul bandwidth allocation	Machine learning	Distributed	Non-cooperative
[83]	SBSs	System throughput	Machine learning	Centralised	Non-cooperative
[84]	SBSs	Backhaul offloading gains	Machine learning	Centralised	Non-cooperative
[370]	Devices	Cache hit ratio Density of cache-served requests	Stochastic Geometry	Distributed	Non-cooperative
[371]	Devices	Energy consumption	Optimisation	Distributed	Non-cooperative
[101]	Devices	The probability of successful content delivery	Stochastic Geometry	Distributed	Non-cooperative
[102]	Devices	Outage probability	Optimisation	Distributed	Non-cooperative
[98]	Devices	offloading probability	Optimisation	Distributed	Non-cooperative
[99]	Devices	offloading gain	Stochastic Geometry	Centralised	Non-cooperative
[108]	Devices	number of D2D links	Optimisation	Distributed	Non-cooperative
[109]	Devices	Spectral reuse	Optimisation	Distributed	Non-cooperative
[103] [104]	Devices	System throughput	Optimisation	Centralised	Non-cooperative
[105] [106]	Devices	Network throughput	Optimisation	Centralised	Cooperative
[107]	Devices	Coverage probability	Stochastic Geometry	Centralised	Non-cooperative
[110]	Devices	Service success probability	Stochastic Geometry	Distributed	Non-cooperative
[380]	Devices	Coverage probability	Stochastic Geometry	Distributed	Non-cooperative
[114]	Devices	Caching reward	Game theory	Distributed	Non-cooperative
[115]	Devices	Content provisioning costs	Optimisation	Distributed	Non-cooperative
[111]	Devices	Backhaul costs	Graph Theory	Centralised	Non-cooperative
[113]	Devices	Cache hit ratio	Graph Theory	Centralised	Non-cooperative

keeps the most recently accessed content. However, these replacement strategies merely consider content request features in a short time window and may not obtain the global optimal solutions.

Another popular method is to replace the content based on its popularity. In [116], Blasco *et al.* divide time into periods and within each period there is a cache replacement phase. During each cache replacement phase, the content of the lowest popularity is discarded. Apart from historical popularity, Bacstug *et al.* [117] take the future content popularity into consideration as well. They propose a proactive popularity caching (PropCaching) method to estimate content popularity and then determine which content should be evicted.

Mathematically, the cache replacement problem could be formulated as a Markov Decision Process (MDP) [117], [118]. The MDP model can be represented into a tuple $(\mathcal{S}, \mathcal{A}, R(s, a))$. \mathcal{S} refers to the set of possible states for caches. \mathcal{A} is the set of eviction actions. $R(s, a)$ is the reward function that determines the reward when cache performers action a in the state s . The reward is usually modelled as the cache hit or the changes in transmission cost. In [117], Bacstug *et al.* obtain the cache replacement actions based on Q-learning. In [118], the method is upgraded. Wang *et al.* apply deep reinforcement learning to solve it.

IV. EDGE TRAINING

The standard learning approach requires centralising training data on one machine, whilst edge training relies on distributed training data on edge devices and edge servers, which is more secure and robust for data processing. The main idea of edge training is to perform learning tasks where the data is generated or collected with edge computing resources. It is not necessary to send users' personal data to a central server, which effectively solves the privacy problem and saves network bandwidth.

Training data could be solved through edge caching. We discuss how to train an AI model in edge environments in this section. Since the computing capacity on edge devices and edge servers is not as powerful as central servers, the training style changes correspondingly in the edge environment. The major change is the distributed training architecture, which must take the data allocation, computing capability, and network into full consideration. New challenges and problems, e.g., training efficiency, communication efficiency, privacy and security issues, and uncertainty estimates, come along with the new architecture. Next, we discuss these problems in more detail.

A. Training architecture

Training architecture depends on the computing capacity of edge devices and edge servers. If one edge device/server is powerful enough, it could adopt the same training architecture as a centralised server, i.e., training on a single device. Otherwise, cooperation with other devices is necessary. Hence, there are two kinds of training architectures: solo training, i.e., perform training tasks on a single edge device/server, and collaborative training, i.e., few devices and servers work collaboratively to perform training tasks.

1) *Solo training*: Early researchers mainly focus on verifying the feasibility of directly training deep learning models on mobile platforms. Chen *et al.* find that the size of neural network and the memory resource are two key factors that affect training efficiency [119]. For a specific device, training efficiency could be improved significantly by optimising the model. Subsequently, Lane *et al.* successfully implement a constrained deep learning model on smartphones for activity recognition and audio sensing [30]. The demonstration achieves a better performance than shallow models, which demonstrates that ordinary smart devices are qualified for simple deep learning models. Similar verification is also done on wearable devices [31] and embedded devices [32].

2) *Collaborative training*: The most common collaborative training architecture is the master-slave architecture. Federated learning [33] is a typical example, in which a server employs multiple devices and allocates training tasks for them. Li *et al.* develop a mobile object recognition framework, named DeepCham, which collaboratively trains adaptation models [120]. The DeepCham framework consists of one master, i.e., edge server and multiple workers, i.e., mobile devices. There is a training instance generation pipeline on workers that recognises objects in a particular mobile visual domain. The master trains the model using the training instance generated by workers. Huang *et al.* consider a more complex framework with additional scheduling from the cloud [121]. Workers with training instances first uploads a profile of the training instance and requests to the cloud server. Then, the cloud server appoints an available edge server to perform the model training.

Peer-to-peer is another collaborative training architecture, in which participants are equal. Valerio *et al.* adopt such training architecture for data analysis [122]. Specifically, participants first perform partial analytic tasks separately with their own data. Then, participants exchange partial models and refine them accordingly. The authors use an activity recognition model and a pattern recognition model to verify the proposed architecture and find that the trained model could achieve similar performance with the model trained by a centralised server. Similar training architecture is also used in [123] to enable knowledge transferring amongst edge devices.

B. Training Acceleration

Training a model, especially deep neural networks, is often too computationally intensive, which may result in low training efficiency on edge devices, due to their limited computing capability. Hence, some researchers focus on how to accelerate the training at edge. Table IV summaries existing literature on training acceleration.

Chen *et al.* find that the size of a neural network is an important factor that affects the training time [119]. Some efforts [124], [125] investigate transfer learning to speed up the training. In transfer learning, learned features on previous models could be used by other models, which could significantly reduce the learning time. Valery *et al.* propose to transfer features learned by the trained model to local models, which would be re-trained with the local training instances

TABLE IV
LITERATURE SUMMARY OF MODEL ACCELERATION IN TRAINING.

Ref.	Model	Approach	Learning method	Object	Performance
[119]	DNN	Hardware acceleration	Transfer learning	Review training factors	N/A
[124]	CNN	Hardware acceleration	Transfer learning	Alleviate memory constraint	Faster than Caffe-OpenCL trained
[125]	CNN	Hardware acceleration parameter quantisation	Transfer learning	Alleviate memory constraint	Faster than Caffe-OpenCL trained
[127]	DNN	Analog memory	Transfer learning	Better energy-efficiency	Close to software baseline of 97.9
[128]	RF, ET, NB LR, SVM	Human annotation	Incremental learning	Investigate iML for HAR	93.3% accuracy
[129]	Naive Bayes	Human annotation	Incremental learning	Reduce limitations in learning	6-8 hours to train a model
[123]	CNN	Software acceleration	Transfer learning	Reduce required labelled data	50× faster
[130]	Statistical model	Software acceleration	Federated learning	Address statistical challenges	Outperform global, local manners
[131]	GCN	Software-hardware Co-optimization	Supervised learning	Accelerate GCN training on heterogeneous platform	An order of magnitude faster

[124]. Meanwhile, they exploit the shared memory of the edge devices to enable the collaboration between CPU and GPU. This approach could reduce the required memory and increase computing capacity. Subsequently, the authors further accelerate the training procedure by compressing the model by replacing float-point with 8-bit fixed point [125].

In some specific scenarios, interactive machine learning (iML) [382], [383] could accelerate the training. iML engages users in generating classifiers. Users iteratively supply information to the learning system and observe the output to improve the subsequent iterations. Hence, model updates are more fast and focused. For example, Amazon often asks users targeted questions about their preferences for products. Their preferences are promptly incorporated into a learning system for recommendation services. Some efforts [126], [128] adopt such approach in model training on edge devices. Shahmohammadi *et al.* apply iML on human activity recognition, and find that only few training instances are enough to achieve a satisfactory recognition accuracy [128]. Based on such theory, Flutura *et al.* develop DrinkWatch to recognise drink activities based on sensors on smartwatch [129].

In a collaborative training paradigm, edge devices are enabled to learn from each other to increase learning efficiency. Xing *et al.* propose a framework, called RecycleML, which uses cross modal transfer to speed up the training of neural networks on mobile platforms across different sensing modalities in the scenario that the labelled data is insufficient [123]. They design an hourglass model for knowledge transfer for multiple edge devices, as shown in Fig. 14. The bottom part denotes lower layers of multiple specific models, e.g., AudioNet, IMUNet, and VideoNet. The middle part represents the common layers of these specific models. These models project their data into the common layer for knowledge transfer. The upper part represents the task-specific layers of different models, which are trained in a targeted fashion. Experiments show that the framework achieves 50x speedup for the training. Federated learning could be also applied to accelerate the training of models on distributed edge devices. Smith *et al.* propose a systems-aware framework to optimise the setting of federated learning (e.g., update cost and stragglers) and to speed up the training [130].

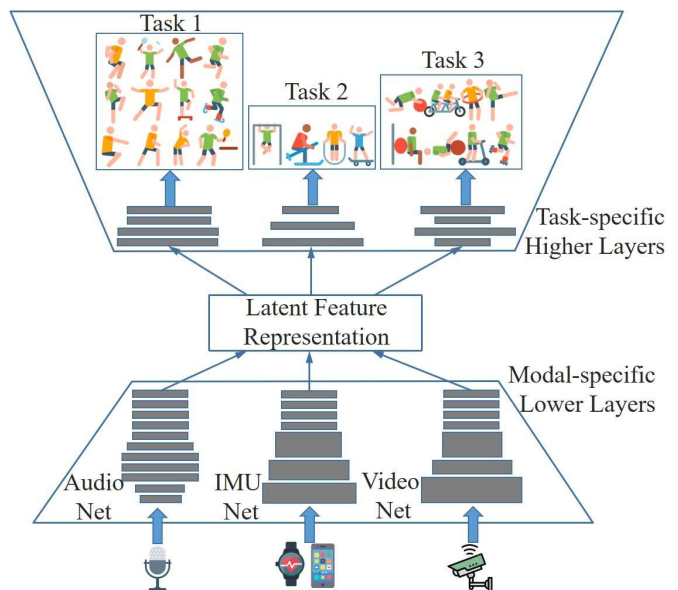


Fig. 14. Illustration of the hourglass model. The lower part represents lower layers of specific sensing models. The latent feature representation part is the common layer. Lower layers project their data into this layer for knowledge transfer. The upper part represents task-specific higher layers, which are trained for specific recognition tasks.

C. Training optimisation

Training optimisation refers to optimising the training process to achieve some given objectives, e.g., energy consumption, accuracy, privacy-preservation, security-preservation, etc. Since solo training is similar to training on a centralised server to a large extent, existing work mainly focuses on collaborative training. Federated learning is the most typical collaborative training architecture, and almost all literature on collaborative training is relevant to this topic.

Federated learning is a kind of distributed learning [122], [384], [385], which allows training sets and models to be located in different, non-centralised positions, and learning can occur independent of time and places. This training architecture is first proposed by Google, which allows smartphones to collaboratively learn a shared model with their local training data, instead of uploading all data to a central cloud server [33]. The learning process of federated learning

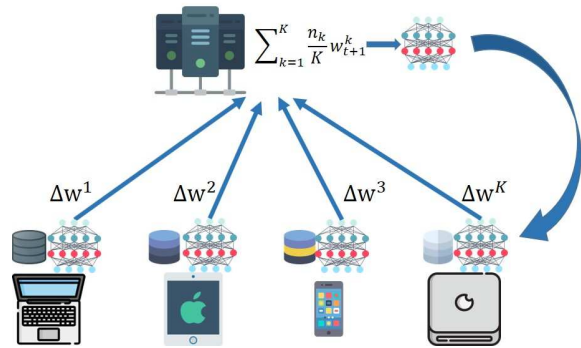


Fig. 15. The illustration of federated learning. Each training participant trains the shared model with cached data. After training, an update, i.e., Δw will be uploaded to the central server. All received updates from training participants would be aggregated to update the shared model. Then, the new shared model would be sent to all edge devices for the next round of learning.

is shown as Fig. 15. There is a untrained shared model on the central server, which will be allocated training participants for training. Training participants, i.e., edge devices train the model with the local data. After local learning, changes of the model are summarised as a small focused update, which will be sent to the central server through encrypted communication. The central server averages received changes from all mobile devices and updates the shared model with the averaged result. Then, mobile devices download the update for their local model and repeats the procedure to continuously improve the shared model. In this learning procedure, only the encrypted changes are uploaded to the cloud and the training data of each mobile user remains on mobile devices. Transfer learning and edge computing are combined to learn a smarter model for mobile users. In addition, since learning occurs locally, federated learning could effectively protect user privacy, when compared with a traditional centralised learning approach.

Typical edge devices in federated learning are smartphones with unreliable and slow network connections. Moreover, due to the unknown mobility, these devices may be intermittently available for working. Hence, the communication efficiency between smartphones and the central server is of the utmost importance to the training. Specifically, there are factors affecting the communication efficiency: communication frequency and communication cost. In addition, the update from edge devices is vulnerable to malicious users. Hence privacy and security issues should also be considered. We discuss these problems in detail next. Table V summarises literature on training optimisation.

1) *Communication Frequency*: In federated learning, communication between edge devices and the cloud server is the most important operation, which uploads the updates from edge devices to the cloud server and downloads the aggregated update from the shared model to local models. Due to the possible unreliable network condition of edge devices, minimising the number of update rounds, i.e., communication frequency between edge devices and cloud server is necessary. Jakub *et al.* are the first to deploy federated learning framework

and propose the setting for federated optimisation [132]. In [55], the authors characterise the training data as massively distributed (data points are stored across massive edge devices), non-IID (training set on devices may be drawn from different distributions), and unbalanced (different devices have different number of training samples). In each round, each device sends an encrypted update to the central server. Then they propose a federated stochastic variance reduced gradient (FSVRG) algorithm to optimise the federated learning. They find that the central shared model could be trained with a small number of communication rounds.

McMahan *et al.* propose a federated averaging algorithm (FedAvg) to optimise federated learning in the same scenario with [55], [132] and further evaluate the framework with five models and four datasets to prove the robustness of the framework [133]. Although FedAvg could reduce the number of communication rounds for certain datasets, Zhao *et al.* find that using this algorithm to train CNN models with highly skewed non-IID dataset would result in the significant reduction of the accuracy [134]. They find the accuracy reduction results from the weight divergence, which refers to the difference of learned weights between two training processes with the same weight initialisation. Earth mover's distance (EMD) between the distribution over classes on each mobile device and the distribution of population are used to quantify the weight divergence. They then propose to extract a subset of data, which is shared by all edge devices to increase the accuracy.

Strategies that reduce the number of updates should be on the premise of not compromising the accuracy of the shared model. Wang *et al.* propose a control algorithm to determine the optimal number of global aggregations to maximise the efficiency of local resources [137]. They first analyse the convergence bound of SGD based federated learning. Then, they propose an algorithm to adjust the aggregation frequency in real-time to minimise the resource consumption on edge devices, with the joint consideration of data distribution, model characteristics, and system dynamics.

Above-mentioned works adopt a synchronous updating method, where in each updating round, updates from edge devices are first uploaded to the central server, and then aggregated to update the shared model. Then the central server allocates aggregated updates to each edge device. Some researchers think that it is difficult to synchronise the process. On one hand, edge devices have significantly heterogeneous computing resources, and the local model are trained asynchronously on each edge device. On the other hand, the connection between edge devices and the central server is not stable. Edge devices may be intermittently available, or response with a long latency due to the poor connection. Wang *et al.* propose an asynchronous updating algorithm, called CO-OP through introducing an age filter [138]. The shared model and downloaded model by each edge device would be labelled with ages. For each edge device, if the training is finished, it would upload its update to the central server. Only when the update is neither obsolete nor too frequent, it will be aggregated to the shared model. However, most works adopt synchronous approaches in federated learning, due to its effectiveness [133], [143].

2) *Communication cost*: In addition to communication frequency, communication cost is another factor that affects the communication efficiency between edge devices and the central server. Reducing the communication cost could significantly save bandwidth and improve communication efficiency. Konevny *et al.* propose and prove that the communication cost could be lessened through structured and sketched updates [55], [144]. The structured update means learning an update from a restricted space that could be parametrised with few variables through using low rank and random mask structure, while sketched update refers to compressing the update of the full model through quantisation, random rotations and sub-sampling.

Lin *et al.* find most of the gradient update between edge devices and the central server is redundant in SGD based federated learning [145]. Compressing the gradient could solve the redundancy problem and reduce the update size. However, compression methods, such as gradient quantisation and gradient sparsification would lead to the decreased accuracy. They propose a deep gradient compression (DGC) method to avoid the loss of accuracy, which use momentum correction and local gradient clipping on top of the gradient sparsification. Hardy *et al.* also try to compress the gradient and propose a compression algorithm, called AdaComp [146]. The basic idea of AdaComp is compute staleness on each parameter and remove a large part of update conflicts.

Smith *et al.* propose to combine multi-task learning and federated learning together, which train multiple relative models simultaneously [130]. It is quite cost-effective for a single model, during the training. They develop an optimisation algorithm, named MOCHA, for federated setting, which allows personalisation through learning separate but related models for each participant via multi-task learning. They also prove the theoretical convergence of this algorithm. However, this algorithm is inapplicable for non-convex problems.

Different from the client-to-server federated learning communication in [55], [145], [146], Caldas *et al.* propose to compress the update from the perspective of server-to-client exchange and propose Federated Dropout to reduce the update size [147]. In client-to-server paradigm, edge devices download the full model from the server, while in a server-to-client paradigm, each edge device only downloads a sub-model, which is a subset of the global shared model. This approach both reduces the update size and the computation on edge devices.

3) *Privacy and security issues*: After receiving updates from edge devices, the central server needs to aggregate these updates and construct an update for the shared global model. Currently, most deep learning models rely on variants of stochastic gradient descent (SGD) for optimisation. FedAvg, proposed in [133], is a simple but effective algorithm to aggregate SGD from each edge device through weighted averaging. Generally, the update from each edge device contains significantly less information of the users' local data. However, it is still possible to learn the individual information of a user from the update [148], [386]. If the updates from users are inspected by malicious hackers, participant edge users' privacy would be threatened. Bonawitz *et al.* propose Secure

Aggregation to aggregate the updates from all edge devices, which makes the participant updates un-inspectable by the central server [149]. Specifically, each edge device uploads a masked update, i.e., parameter vector to the server, and then the server accumulates a sum of the masked update vectors. As long as there is enough edge devices, the masks would be counteracted. Then, the server would be able to unmask the aggregated update. During the aggregation, all individual updates are non-inspectable. The server can only access the aggregated unmasked update, which effectively protect participants' privacy. Liu *et al.* introduce homomorphic encryption to federated learning for privacy protection [150]. Homomorphic encryption [151] is an encryption approach that allows computation on ciphertexts and generates an encrypted result, which, after decryption, is the same with the result achieved through direct computation on the plain text. The central server could directly aggregate the encrypted updates from participants.

Geyer *et al.* propose an algorithm to hide the contribution of participants at the clients' based on differential privacy [152]. Similar to differential privacy-preserving traditional approaches [154], [155], the authors add a carefully calibrated amount of noise to the updates from edge devices in federated learning. The approach ensures that attackers could not find whether an edge device participated during the training. A similar differential privacy mechanisms are also adopted in federated learning based recurrent language model and federated reinforcement learning in [156] and [157].

In federated learning, the participants could observe intermediate model states and contribute arbitrary updates to the global shared model. All aforementioned research assumes that the participants in federated learning are un-malicious, which provides a real training set and uploads the update based on the training set. However, if some of the participants are malicious, who uploads erroneous updates to the central server, the training process fails. In some cases, the attack would result in large economic losses. For example, in a backdoor attacked face recognition based authentication system, attackers could mislead systems to identify them as a person who can access a building through impersonation. According to their attack patterns, attacks could be classified into two categories: data-poisoning and model-poisoning attacks. Data-poisoning means compromising the behaviour and performance of the model through changing the training set, e.g., accuracy, whilst model-poisoning only change the model's behaviour on specific inputs, without impacting the performance on other inputs. The impact of data-poisoning attack is shown as Fig. 16.

The work in [159] tests the impact of a data-poisoning attack on SVM through injecting specially crafted training data, and find that the SVM's test error increases with the attack. Steinhardt *et al.* construct the approximate upper bound of the attack loss on SVM and provides a solution to eliminate the impact of the attack [160]. In particular, they first remove outliers residing outside a feasible bound, and then minimise the margin-based loss on the rest data.

Fung *et al.* evaluate the impact of sybil-based data-poisoning attack on federated learning and propose a defense scheme, FoolsGold, to solve the problem [161]. A sybil-

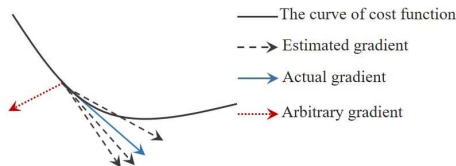


Fig. 16. The impact of data-poisoning attack. The black dashed arrows refers to the gradient estimates computed by honest participants, which are distributed around the actual gradient. The red dotted arrow indicates the arbitrary gradient computed by malicious participants, which hampers the convergence of the training.

based attack [162] means that a participant edge device has a wrong training dataset, in which the data is the same with other participants whilst its label is wrong. For example, in digit recognition, the digit ‘1’ is labelled with ‘7’. They find that attackers may overpower other honest participants by poisoning the model with sufficient sybils. The proposed defense system, FoolGold, is based on contribution similarity. Since sybils share a common objective, their updates appear more similar than honest participants. FoolGold eliminates the impact of sybil-based attacks through reducing the learning rate of participants that repeatedly upload the same updates.

Blanchard *et al.* evaluate the Byzantine resilience of SGD in federated learning [163]. Byzantine refers to arbitrary failures in federated learning, such as erroneous data and software bugs. They find that linear gradient aggregation has no tolerance for even one Byzantine failure. Then they propose a Krum algorithm for aggregation with the tolerance of f Byzantines out of n participants. Specifically, the central server computes pairwise distances amongst all updates from edge devices, and takes the sum of $n - f - 2$ closest distance for all updates. The update with the minimum sum would be used to update the global shared model. However, all updates from edge devices are inspectable during computation, which may result in the risk of privacy disclosure. Chen *et al.* propose to use the geometric median of gradients as the update in federated learning [164]. This approach could tolerate q Byzantine failures up to $2q(1 + \epsilon) \leq m$, in which q is the number of Byzantine failures, m refers to the headcount of participants, and ϵ is a small constant. This approach groups all participants into mini-batches. However, Yin *et al.* find that the approach fails if there is one Byzantine in each mini-batch [165]. They then propose a coordinate-wise median based approach to deal with the problem.

In fact, data-poisoning based attacks on federated learning is low in efficiency in the condition of small numbers of malicious participants. Because there are usually thousands of edge devices participating in the training in federated learning. The arbitrary update would be offset by averaging aggregation. In contrast, model-poisoning based attacks are more effective. Attackers directly poison the global shared model, instead of the updates from thousands of participants. Attackers introduce hidden backdoor functionality in the global shared model. Then, attackers use key, i.e., input with attacker-chosen features to trigger the backdoor. The model-poisoning based attack is shown as Fig. 17. Works on model-poisoning mainly

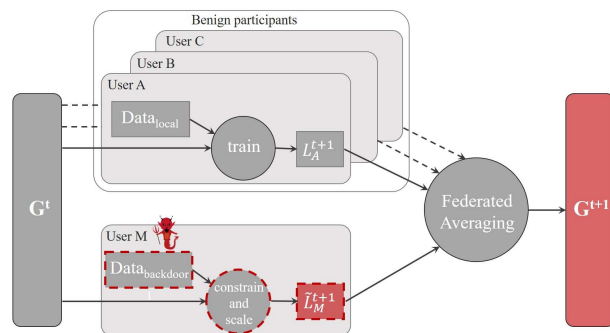


Fig. 17. Overview of model-poisoning based attack. Attackers train the backdoor model with local data. Then, attackers scale up the weight of the update to guarantee that the backdoor model would not be cancelled out by other updates.

focus on the problem of how backdoor functionality is injected in federated learning. Hence, we will focus on this direction as well.

Chen *et al.* evaluate the feasibility of conducting a backdoor in deep learning through adding few poisoning samples into the training set [166]. They find that only 5 poisoning samples out of 600,000 training samples are enough to create a backdoor. Bagdasaryan *et al.* propose a model replacement technique to open a backdoor to the global shared model [167]. As we aforementioned, the central server computes an update through averaging aggregation on updates from thousands of participants. The model replacement method scales up the weights of the ‘backdoored’ update to ensure that the backdoor survives the averaging aggregation. This is a single-round attack. Hence, such attack usually occurs during the last round update of federated learning. Different from [167], Bhagoji *et al.* propose to poison the shared model even when it is far from convergence, i.e., the last round update [168]. To prevent that, the malicious update is offset by updates from other participants, they propose a explicit boosting mechanism to negate the aggregation effect. They evaluate the attack technique against some famous attack-tolerant algorithms, i.e., Krum algorithm [163] and coordinate-wise median algorithm [165], and find that the attack is still effective.

D. Uncertainty Estimates

Standard deep learning method for classification and regression could not capture model uncertainty. For example, in model for classification, obtained results may be erroneously interpreted as model confidence. Such problems exist as well in edge intelligence. Efficient and accurate assessment of the deep learning output is of crucial importance, since the erroneous output may lead to undesirable economy loss or safety consequence in practical applications.

In principle, the uncertainty could be estimated through extensive tests. [169] propose a theoretical framework that casts dropout training in DNNs as approximate Bayesian inference in deep Gaussian processes. The framework could be used to model uncertainty with dropout neural networks through extracting information from models. However, this process is computation intensive, which is not applicable on

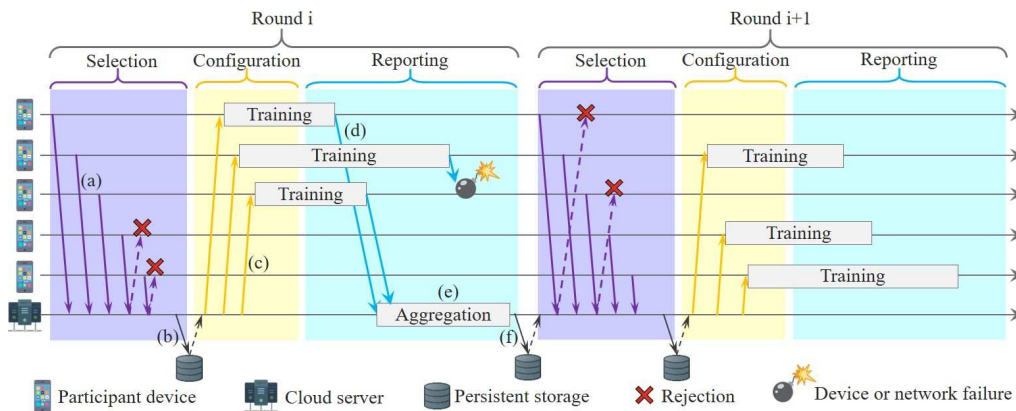


Fig. 18. The illustration of a TensorFlow-based federated learning system. (a) edge devices register to participate in federated training. Un-selected devices would be suggested to participate in the next round. (b) server reads the checkpoint of the model from storage. (c) server sends a shared model to each selected edge device. (d) edge devices train the model with local data and uploads their updates. (e) All received updates are aggregated. (f) the server save the checkpoint of the model.

mobile devices. This approach is based on sampling, which requires sufficient output samples for estimation. Hence, the main challenge to estimate uncertainty on mobile devices is the computational overhead. Based on the theory proposed in [169], Yao *et al.* propose RDeepSense, which integrates scoring rules as training criterion that measures the quality of the uncertainty estimation to reduce energy and time consumption [170]. RDeepSense requires to re-train the model to estimate uncertainty. The authors further propose ApDeepSense, which replaces the sampling operations with layer-wise distribution approximations following closed-form representations [387].

E. Applications

Bonawitz *et al.* develop a scalable product system for federated learning on mobile devices, based on TensorFlow [171]. In this system, each updating round consists of three phases: client selection, configuration, and reporting, as shown in Fig. 18. In the client selection phase, eligible edge devices, e.g., devices with sufficient energy and computing resources, periodically send messages to the server to report the liveness. The server selects a subset among them according to a given objective. In the configuration phase, the server sends a shared model to each selected edge device. In the reporting phase, each edge device reports the update to the server, which would be aggregated to update the shared model. This protocol presents a framework of federated learning, which could adopt multiple strategies and algorithms in each phase. For example, the client selection algorithm proposed in [172] could be used in the client selection phase. The communication strategy in [55], [132] could be used for updating, and the FedAvg algorithm in [133] is adopted as an aggregation approach.

Researchers from Google have been continuously working on improving the service of Gboard with federated learning. Gboard consists of two parts: text typing and a search engine. The text typing module is used to recognise users' input, whilst the search engine provides user relevant suggestions according to their input. For example, when you type 'let's eat', Gboard may display the information about nearby restaurants. Hard *et*

al. train a RNN language model using a federated learning approach to improve the prediction accuracy of the next-word for Gboard [173]. They compare the training result with traditional training methods on a central server. Federated learning achieves comparable accuracy with the central training approach. Chen *et al.* use federated learning to train a character-level RNN to predict high-frequency words on Gboard [174]. The approach achieves 90.56% precision on a publicly-available corpus. McMahan *et al.* undertake the first step to apply federated learning to enhance the search engine of Gboard [33]. When users search with Gboard, information about the current context and whether the clicked suggestion would be stored locally. Federated learning processes this information to improve the recommendation model. Yang *et al.* further improve the recommendation accuracy by introducing an additional *triggering model* [175]. Similarly, there are some works [176], [176] focusing on emoji prediction on mobile keyboards.

Federated learning has great potential in the medical imaging domain, where patient information is highly sensitive. Sheller *et al.* train a brain tumour segmentation model with data from multi-institution by applying federated learning [177]. The encrypted model is first sent to data owners, i.e., institutions, then the data owners decode, train, encrypt and upload the model back to the central aggregator. Roy *et al.* further develop an architecture of federated learning that uses peer-to-peer communications to replace the central aggregator towards medical applications [178].

Samarakoon *et al.* apply federated learning in vehicular networks to jointly allocate power and resources for ultra reliable low latency communication [179]. Vehicles train and upload their local models to the roadside unit (RSU), and RSU feeds back the global model to vehicles. Vehicles could use the model to estimate the queue length in city. Based on the queue information, the traffic system could reduce the queue length and optimise the resource allocation.

Nguyen *et al.* develop DIoT, a self-learning system to detect infected IoT devices by malware botnet in smart home envi-

TABLE V
LITERATURE SUMMARY OF TRAINING OPTIMISATION.

Ref.	Problem	Solution	Dataset	Performance
[132]	Communication efficiency	FSVRG	Google+ posts	Less rounds
[55]	Communication efficiency	FSVRG	Google+ posts	Less rounds
[133]	Communication efficiency	FedAvg	MNIST, CIFAR-10, KWS	10 – 100× less rounds
[134]	Communication efficiency	FedAvg, data sharing	MNIST, CIFAR-10, KWS	30% higher accuracy
[135]	Uncoordinated communication	Incentive mechanism Admission control	N/A	22% gain in reward
[136]	Incentive mechanism	Deep reinforcement learning	MNIST	Lower communication cost
[137]	Communication frequency	Aggregation control	MNIST	Near to the optimum
[138]	Communication frequency	CO-OP	MNIST	80% accuracy
[139]	Communication bandwidth	Beamforming design	CIFAR-10	Lower training loss, higher accuracy
[140]	Noisy communication	Successive convex approximation	MNIST	Approach to centralized method
[141]	Wireless fading channel	D-DSGD, CA-DSGD	MNIST	Converges faster, higher accuracy
[142]	Single point of failure	Server-less aggregation	Real-world sensing data	One order of magnitude less rounds
[55]	Communication cost	Structured update sketched update	CIFAR-10, Reddit	85% accuracy
[145]	Communication cost	DGC	ImageNet, Penn Treebank Cifar10, Librispeech Corpus	270 – 600× smaller update size
[146]	Communication cost	Compression staleness mitigation	MNIST	191× smaller update size
[130]	Multi-task learning Communication cost	MOCHA	Human Activity Recognition GLEAM, Vehicle Sensor	Lowest prediction error
[147]	Communication cost	Federated Dropout	MNIST, EMNIST, CIFAR-10	28× smaller update size
[149]	Information revealing	Secure Aggregation	N/A	1.98× expansion for 2^{14} users
[150]	Privacy protection	Homomorphic encryption	NUS-WIDE, Default-Credit	Little accuracy drop
[152]	Privacy protection	Differentially privacy	non-IID MNIST	Privacy maintained
[153]	Privacy protection	Differentially privacy K-client random scheduling	MNIST	Privacy maintained
[157]	Privacy protection	Gaussian differential	WHS, CT, WHG	F1 score 10% - 20% higher
[158]	Privacy protection	SecureBoost	Credit 1, Credit 2	Higher accuracy, F1-score
[156]	Privacy protection	Differentially privacy	Reddit posts	Similar to un-noised models
[161]	Sybil-based attack	FoolGold	MNIST, VGGFace2	Attacking rate <1%
[163]	Byzantine failure	Krum aggregation	MNIST, Spambase	Toleratable for 45% Byzantines
[164]	Byzantine failure	Batch gradients median	N/A	$2q(1 + \epsilon) \leq m$ Byzantines
[165]	Byzantine failure	Coordinate-wise median	N/A	Optimal statistical error rate
[168]	Backdoor attack	Explicit boosting	Fashion-MNIST, Adult Census	100% backdoor accuracy

ronments [180]. IoT devices connect to the Internet through a gateway. They design two models for IoT device identification and anomaly detection. These two models are trained through the federated learning approach.

V. EDGE INFERENCE

The exponential growth of network size and the associated increase in computing resources requirement have been become a clear trend. Edge inference, as an essential component of edge intelligence, is usually performed locally on edge devices, in which the performance, i.e., execution time, accuracy, energy efficiency, etc. would be bounded by technology scaling. Moreover, we see an increasingly widening gap between the computation requirement and the available computation capacity provided by the hardware architecture [180]. In this section, we discuss various frameworks and approaches that contribute to bridging the gap.

A. Model Design

Modern neural network models are becoming increasingly larger, deeper and slower, they also require more computation resources [183], [388], [389], which makes it quite difficult to directly run high performance models on edge devices with limited computing resources, e.g., mobile devices, IoT

terminals and embedded devices. Guo *emph* evaluate the performance of DNN on edge device and find inference on edge devices costs up to two orders of magnitude greater energy and response time than central server [390]. Many recent works have focused on designing lightweight neural network models, which could be performed on edge devices with less requirements on the hardware. According to the approaches of model design, existing literature could be divided into two categories: architecture search, and human-invented architecture. The former is to let machine automatically design the optimal architecture, while the latter is to design architectures by human.

1) *Architecture Search*: Designing neural network architectures is quite time-consuming, which requires substantial effort of human experts. One possible research direction is to use AI to enable machine search for the optimal architecture automatically. In fact, some automatically searched architectures, e.g., NASNet [183], AmoebaNet [184], and Adanet [185], could achieve competitive even much better performance in classification and recognition. However, these architectures are extremely hardware-consuming. For example, it requires 3150 GPU days of evolution to search for the optimal architecture for CIFAR-10 [184]. Mingxing *et al.* adopt reinforcement learning to design mobile CNNs, called MnasNet, which could balance accuracy and inference latency [186]. Different from

[183]–[185], in which only few kinds of cells are stacked, MnasNet cuts down per-cell search space and allow cells to be different. There are more 5×5 depthwise convolutions in MnasNet, which makes MnasNet more resource-efficient compared with models that only adopt 3×3 kernels.

Recently, a new research breakthrough of differentiable architecture search (DARTS) [187] could significantly reduce dependence on hardware. Only four GPU days are required to achieve the same performance as [184]. DARTS is based on continuous relaxation of the architecture representation and uses gradient descent for architecture searching. DARTS could be used for both convolutional and recurrent architectures.

Architecture search is hot research area and has a wide application future. Most literature on this area is not specially for edge intelligence. Hence, we will not further discuss on this field. Readers interested in this field could refer to [391], [392].

2) *Human-invented Architecture*: Although architecture search shows good ability in model design, its requirement on hardware holds most researchers back. Existing literature mainly focuses on human-invented architecture. Howard *et al.* use depth-wise separable convolutions to construct a lightweight deep neural network, MobileNets, for mobile and embedded devices [188]. In MobileNets, a convolution filter is factorised into a depth-wise and a point-wise convolution filter. The drawback of depth-wise convolution is that it only filters input channels. Depth-wise separable convolution, which combines depth-wise convolution and 1×1 point-wise convolution could overcome this drawback. MobileNet uses 3×3 depth-wise separable convolutions, which only requires 8-9 times less computation than standard ones. Moreover, depth-wise and point-wise convolutions could also be applied to implement keyword spotting (KWS) models [189] and depth estimation [190] on edge devices.

Group convolution is another way to reduce computation cost for model designing. Due to the costly dense 1×1 convolutions, some basic architectures, e.g., Xception [393] and ResNeXt [394] cannot be used on resource-constrained devices. Zhang *et al.* propose to reduce the computation complexity of 1×1 convolutions with pointwise group convolution [191]. However, there is a side effect brought on by group convolution, i.e., outputs of one channel are only derived from a small part of the input channels. The authors then propose to use a *channel shuffle* operation to enable information exchanging among channels, as shown in Fig 19.

Depth-wise convolution and group convolution are usually based on ‘sparsely-connected’ convolutions, which may hamper inter-group information exchange and degrades model performance. Qin *et al.* propose to solve the problem with merging and evolution operations [192]. In merging operation, features of the same location among different channels are merged to generate a new feature map. Evolution operation extracts the information of location from the new feature map and combines extracted information with the original network. Therefore, information is shared by all channels, so that the information loss problem of inter-groups is effectively solved.

3) *Applications*: A large number of models have been designed for various applications, including face recognition

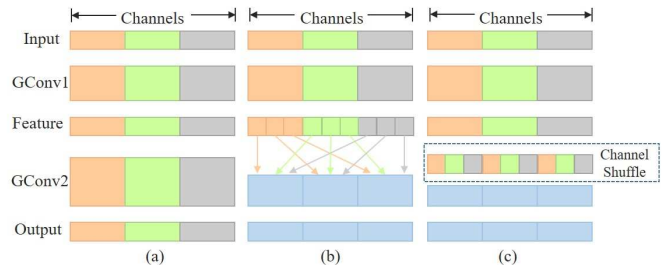


Fig. 19. Illustration of channel shuffle. GConv refers to group convolution. (a) Two stacked convolution layers. Each output channel is related with an input channel of the same group. (b) GConv2 takes data from different groups to make full relations with other channels. (c) The implementation of channel shuffle, which achieves the same effect with (b).

[181], [182], [193], human activity recognition (HAR) [194]–[202], vehicle driving [203]–[206], and audio sensing [207], [208]. We introduce such applications next.

Face verification is increasingly attracting interests in both academic and industrial areas, and it is widely used in device unlocking [395] and mobile payments [396]. Particularly, some applications, such as smartphone unlocking need to run locally with high accuracy and speed, which is challenging for traditional big CNN models due to constrained resources on mobile devices. Sheng *et al.* present a compact but efficient CNN model, MobileFaceNets, which uses less than 1 million parameters and achieves similar performance to the latest big models of hundreds MB size [181]. MobileFaceNets uses a global depth-wise convolution filter to replace the global average pooling filter and carefully design a class of face feature. Chi *et al.* further lighten the weight of MobileFaceNets and presents MobiFace [182]. They adopt the Residual Bottleneck block [193] with expansion layers. Fast downsampling is also used to quickly reduce the dimensions of layers over 14×14 . These two adopted strategies could maximise the information embedded in feature vectors and keep low computation cost.

Edge intelligence could be used to extract contextual information from sensor data and facilitate the research on Human Activity Recognition (HAR). HAR refers to the problem of recognising when, where, and what a person is doing [397], which could be potentially used in many applications, e.g., healthcare, fitness tracking, and activity monitoring [398], [399]. Table VI compares existing HAR technologies, regarding to their frameworks, models, ML methods, and objects. The challenges of HAR on edge platforms could summarised as follows.

- Commonly used classifiers for HAR, e.g., naive Bayes, SVM, DNN, are usually computation-intensive, especially when multiple sensors are involved.
- HAR requires to support near-real-time user experience in many applications.
- Very limited amount of labelled data is available for training HAR models.
- The data collected by on-device sensor includes noise and ambiguity.

Sourav *et al.* investigate how to deploy Restricted Boltzmann Machines (RBM)-based HAR models on smartwatch

TABLE VI
COMPARISON OF DIFFERENT HAR APPLICATIONS.

Ref.	Model	ML method	Objective	Dataset
[194]	RBM	Unsupervised Learning	Energy-efficiency, higher accurate	Opportunity dataset
[195]	CNN	Deep Learning	Improve accuracy	UCI & WISDM
[196]	CNN	Deep Learning	Improve accuracy	RealWorld HAR
[197]	LSTM	Incremental learning	Minimise resource consumption	Heterogeneity Dataset
[198]	CNN	Multimodal Deep Learning	Integrate sensor data	Opportunity dataset
[199]	Heuristic function	Supervised learning	Automatic labelling	38 day-long dataset
[200]	Random forest Naive bayes decision tree	Ensemble learning	Detect label errors	CIMON
[201]	CNN & RNN	Supervised learning	Reduce data noise	Opportunity dataset
[202]	CNN & RNN	Supervised Learning	Heterogeneous sensing quality	Opportunity dataset

platforms, i.e., the Qualcomm Snapdragon 400 [194]. They first test the complexity of a model that a smartwatch can afford. Experiments show that although a simple RBM-based activity recognition algorithm could achieve satisfactory accuracy, the resource consumption on a smartwatch platform is unacceptably high. They further develop pipelines of feature representation and RBM layer activation functions. The RBM model could effectively reduce energy consumption on smartwatches. Bandar *et al.* introduce time domain statistical features in CNN to improve the recognition accuracy [195]. In addition, to reduce the over-fitting problem of their model, they propose a data augmentation method, which applies a label-preserving transformation on raw data to create new data. The work is extended with extracting position features in [196].

Although deep learning could automatically extract features by exploring hidden correlations within and between data, pre-trained models sometimes cannot achieve the expected performance due to the diversities of devices and users, e.g., the heterogeneity of sensor types and user behaviour [400]. Prahalathan *et al.* propose to use on-device incremental learning to provide a better service for users [197]. Incremental learning [401] refers to a secondary training for a pre-trained model, which constrains newly learned filters to be linear combinations of existing ones. The re-trained model on mobile devices could provide personalised recognition for users.

Collecting fine-grained datasets for HAR training is challenging, due to a variety of available sensors, e.g., different sampling rates and data generation models. Valentin *et al.* propose to use RBM architecture to integrate sensor data from multiple sensors [198]. Each sensor input is processed by a single stacked restricted Boltzmann machine in RBM model. Afterwards, all outputted results are merged for activity recognition by another stacked restricted Boltzmann machine. Supervised machine learning is a most commonly utilised approach for activity recognition, which requires a large amount of labelled data. Manually labelling requires extremely large amounts of effort. Federico *et al.* propose a knowledge-driven automatic labelling method to deal with the data annotation problem [199]. GPS data and step count information are used to generate weak labels for the collected raw data. However, such an automatic annotation approach may create labelling errors, which impacts the quality of the collected data. There are three types of labelling errors, including inaccurate timestamps, mislabelling, and multi-action labels.

Multi-action labels means that individuals perform multiple different actions during the same label. Xiao *et al.* solve the last two labelling errors through an ensemble of four stratified trained classifiers of different strategies, i.e., random forest, naive bayes, and decision tree [200].

The data collected by on-device sensors maybe noisy and it is hard to eliminate [400], [402]. For example, in movement tracking application on mobile devices, the travelled distance is computed with the sensory data, e.g., acceleration, speed, and time. However, the sensory data maybe noisy, which will result in estimation errors. Yao *et al.* develop DeepSense, which could directly extract robust noise features of sensor data in a unified manner [201]. DeepSense combines CNN and RNN together to learn the noise model. In particular, the CNN in DeepSense learns the interaction among sensor modalities, while the RNN learn the temporal relationship among them based on the output of the CNN. The authors further propose QualityDeepSense with the consideration of the heterogeneous sensing quality [202]. QualityDeepSense hierarchically adds sensor-temporal attention modules into DeepSense to measure the quality of input sensory data. Based on the measurement, QualityDeepSense selects the input with more valuable information to provide better predictions.

Distracted driving is a key problem, as it potentially leads to traffic accidents [403]. Some researchers address this problem by implementing DL models on smartphones to detect distracted driving behaviour in real-time. Christopher *et al.* design DarNet, a deep learning based system to analyse driving behaviours and to detect distracted driving [203]. There are two modules in the system: data collection and analytic engine. There is a centralised controller in the data collection component, which collects two kinds of data, i.e., IMU data from drivers' phones and images from IoT sensors. The analytic engine uses CNN to process image data, and RNN for sensor data, respectively. The outputs of these two models are combined through an ensemble-based learning approach to enable near real-time distracted driving activity detection. Fig. 20 presents the architecture of DarNet. In addition to CNN and RNN models, there are also other models could be used to detect unsafe driving behaviours, such as SVM [204], HMM [205], and decision tree [206].

Audio sensing has become an essential component for many applications, such as speech recognition [404], emotion detection [405], and smart homes [406]. However, directly

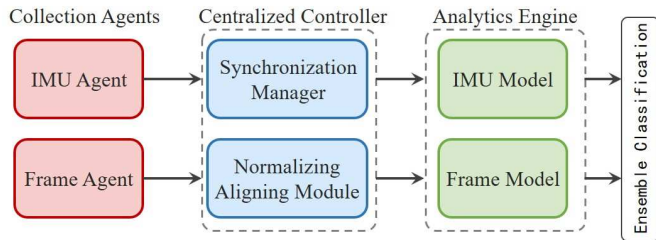


Fig. 20. Architecture of DarNet. IMU agent runs on IoT devices and frame agent runs on mobile devices. A centralised controller collects and pre-processes data for the analytic engine.

running audio sensing models, even just the inference, would introduce a heavy burden on the hardware, such as digital signal processing (DSP) and battery. Nicholas *et al.* develop DeepEar, a DNN based audio sensing prototype for the smartphone platform [207], including four coupled DNNs of stacked RBMs that collectively perform sensing tasks. These four DNNs share the same bottom layers, and each of them is responsible for a specific task, for example, emotion detection, and tone recognition. Experiments show that only 6% of the battery is enough to work through a day with the compromise of 3% accuracy drop. Petko *et al.* further improve the accuracy and reduces the energy consumption through applying multi-task learning and training shared deep layers [208]. The architecture of multi-task learning is shown as Fig. 21, in which the input and hidden layers are shared for audio analysis tasks. Each task has a distinct classifier. Moreover, the shared representation is more scalable than DeepEar, since there is no limitation in the integration of tasks.

B. Model Compression

Although neural networks are quite powerful in various promising applications, the increasing size of neural networks, both in depth and width, results in the considerable consumption of storage, memory and computing powers, which makes it challenging to run neural networks on edge devices. Moreover, statistic shows that the gaps between computational complexity and energy efficiency of deep neural networks and the hardware capacity are growing [407]. It has been proved that neural networks are typically over-parameterised, which makes deep learning models redundant [408]. To implement neural networks on powerless edge devices, large amounts of effort try to compress the models. Model compression aims to lighten the model, improve energy efficiency, and speed up the inference on resource-constraint edge devices, without lowering the accuracy. According to their approaches, we classify these works into five categories: low-rank approximation/matrix factorisation, knowledge distillation, compact layer, parameter quantising, and network pruning. Table VII summarises literature on model compression.

1) *Low-rank Approximation:* The main idea of low-rank approximation is to use the multiplication of low-rank convolutional kernels to replace kernels of high dimension. This is based on the fact that a matrix could be decomposed into the multiplication of multiple matrices of smaller size. For

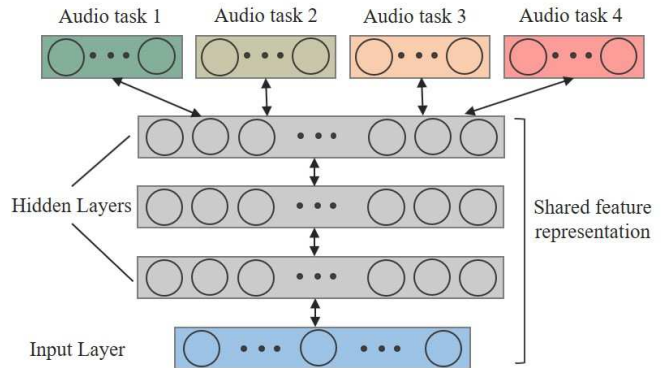


Fig. 21. Illustration of the multi-task audio sensing network.

example, there is a weight matrix W of $m \times k$ dimension. The matrix W could be decomposed into two matrices, i.e., X ($m \times d$) and Y ($d \times k$), and $W = UV$. The computational complexity of matrix W is $O(m \times k)$, while the complexity for the decomposed two matrices is $O(m \times d + d \times k)$. Obviously, the approach could effectively reduce the model size and computation, as long as d is small enough.

Jaderberg *et al.* decompose the matrix of convolution layer $d \times d$ into the multiplication of two matrices $d \times 1$ and $1 \times d$ to compress the CNNs [209]. The authors also propose two schemes to approximate the original filter. Fig. 22 presents the compression process. Fig. 22(a) shows a convolutional layer acting on a single-channel input. The convolutional layer consists of N filters. For the first scheme, they use the linear combination of M ($M < N$) filters to approximate the operation of N filters. For the second scheme, they factorise each convolutional layer into a sequence of two regular convolutional layers but with rectangular filters. The approach achieves a 4.5x acceleration with 1% drop in accuracy. This work is a rank-1 approximation. Maji *et al.* apply this rank-1 approximation on compressing CNN models on IoT devices, which achieves 9x acceleration of the inference [211]. Denton *et al.* explore the approximation of rank- k [210]. They use monochromatic and biclustering to approximate the original convolutional layer.

Kim *et al.* propose a whole network compression scheme with the consideration of entire convolutional and fully connected layers [212]. The scheme consists of three steps: rank selection, low-rank tensor decomposition, and fine-tuning. In particular, they first determine the rank of each layer through a global analytic solution of variational Bayesian matrix factorisation (VBMF). Then they apply Tucker decomposition to decompose the convolutional layer matrix into three components of dimension 1×1 , $D \times D$ (D is usually 3 or 5), and 1×1 , which differs from SVD in [210]. The approach achieves a 4.26x reduction in energy consumption. We note that the component of spatial size $w \times h$ still requires a large amount of computation. Wang *et al.* propose a Block Term Decomposition (BTD) to further reduce the computation in operating the network, which is based on low-rank approximation and group sparsity [213]. They decompose the original weight matrix into the sum of few low multilinear rank weight matrices,

TABLE VII
LITERATURE SUMMARY OF MODEL COMPRESSION.

Ref.	Model	Approach	Object	Performance	Type
[215]	NN	Knowledge distillation	Less resource requirement	Faster	Lossless
[216]	NN	Knowledge distillation	Compress model	80% improvement	Lossless
[217]	NN	Knowledge distillation	Generate thinner model	More accurate and smaller	Improved
[218]	CNN	Knowledge distillation Attention	Improve performance with shallow model	1.1% top-1 better	Improved
[219]	CNN	Knowledge distillation Regularisation	Reduce storage	33.28× smaller	Improved
[220]	CNN	Knowledge distillation	Less memory	40% smaller	Lossless
[221]	GooLeNet	Knowledge distillation	Less memory, acceleration	3× faster, 2.5× less memory	0.4% drop
[222]	CNN	Knowledge distillation	Improve training efficiency	6.4× smaller, 3× faster	Lossy
[223]	CNN	Knowledge distillation	Reconstruct training set	50% smaller	Lossy
[209]	CNN	Low-rank approximation	Reduce runtime	4.5× faster	Lossy
[210]	CNN	Low-rank approximation	Reduce computation	2× faster	Lossy
[211]	CNN	Low-rank approximation	Reduce computation	9× speedup	Lossless
[212]	CNN	Low-rank approximation	Reduce energy consumption	4.26× energy reduction	Lossy
[213]	CNN	Low-rank approximation Group sparsity	Reduce computation	5.91× faster	Improved
[214]	DNN CNN	Low-rank approximation Kernel separation	Use less resource	11.3× less memory 13.3× faster	Lossless
[224]	CNN	Compact layer design	Use less resources	3 – 10× faster	
[225]	ResNet	Compact layer design	Training acceleration	28% relative improvement	Improved
[226]	YOLO	Compact layer design	Reduce model complexity	15.1× smaller, 34% faster	Improved
[227]	CNN	Compact layer design	Reduce parameters	50× fewer parameter	
[228]	CNN	Compact layer design	Accelerates training	3.08% top-5 error	
[229]	CNN	Compact layer design Task decomposition	Utilise storage to trade for computing resources	5.17× smaller	Improved
[230]	CNN	Compact layer design	Simplify SqueezeNet	0.89MB total parameter	Lossy
[231]	RNN	Compact layer design	Improve compression rate	7.9× smaller	Lossy
[232]	CNN	Compressive sensing	Training efficiency	6x faster	Improved
[233]	NIN	Network pruning	On-device customisation	1.24× faster	3% Lossy
[234]	VGG-16	Network pruning	Reduce storage	13× fewer	Lossless
[235]	DNN	Network pruning	Higher energy efficiency	20× faster	Improved
[236]	CNN	Network pruning	Reduce iterations	33% fewer	Lossy
[237]	CNN	Network pruning	Speed up inference	10× faster	Lossy
[238]	CNN	Global filter pruning	Accelerate CNN	70% FLOPs reduction	Lossless
[239]	CNN	Network pruning	Energy-efficiency	3.7× reduction on energy	Lossy
[240]	RNN	Network pruning	Reduce model size	98.9% smaller, 94.5% faster 95.7% energy saved	Lossless
[241]	CNN	Network pruning	Reduce memory footprint	5× less computation	Lossless
[242]	CNN	Network pruning Data reuse	Maximise data reusability	1.43× faster, 34% smaller	Lossless
[243]	CNN	Channel pruning	Speed up CNN inference	2% higher top-1 accuracy	Improved
[244]	CNN	Progressive Channel Pruning	Effective pruning framework	Up to 44.5% FLOPs	Lossy
[245]	DNN	Debiased elastic group LASSO	Structured Compression of DNN	Several folder smaller	Lossless
[246]	CNN	Filter correlations	Minimal information loss	96.4% FLOPs pruning, 0.95% error	Lossless
[248]	CNN	Vector quantisation	Compress required storage	16 – 24× smaller	Lossy
[249]	NN	Hash function	Reduce model size	8× fewer	Lossy
[250]	VGG-16	Parameter quantisation Network pruning Huffman coding	Compress model	49× smaller	Lossless
[251]	CNN	Parameter quantisation	Compress model	20× smaller, 6× faster	Lossy
[252]	DNN	BinaryConnect	Compress model	State-of-the-art	Improved
[253]	DNN	Network Binarisation	Speed up training	State-of-the-art	Improved
[254]	DNN	Network Binarisation	Reduce model size	32× smaller, 58× faster	Lossy
[255]	DNN	Parameter quantisation Binary Connect	Compress model speed up training	Better than standard SGD	Improved
[210]	DNN	Parameter quantisation Binary Connect	Compress model speed up training	2 – 3× faster, 5 – 10× smaller	Lossy
[256]	HMM	Parameter quantisation	Speed up training	10× speedup at most	Lossless
[257]	LSTM	Quantisation aware training	Recover accuracy loss	4% loss recovered	8.1 % Lossy
[258]	Faster R-CNN	Parameter quantisation	Reduce model size	4.16× smaller	Improved
[259]	CNN	Parameter quantisation	Save energy	4.45fps, 6.48 watts	Lossy
[260]	CNN	Parameter quantisation	Reduce computation	1/10 memory shrinks	Improved
[261]	CNN	Posit number system	Reduce model size	36.4% memory shrinks	Lossy
[262]	MNN	Network Binarisation	Improve energy efficiency	State-of-the-art	Improved
[263]	NN	Network Binarisation	Improve energy efficiency	State-of-the-art	Improved
[247]	CNN	Non-parametric Bayesian	Improve quantisation efficiency	Better than RL methods	Lossy

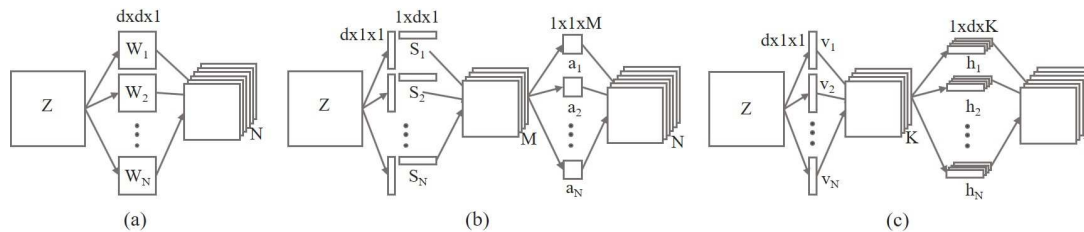


Fig. 22. The decomposition and approximation of a CNN. (a) The original operation of a convolutional layer acting on a single-channel input. (b) The approximation of the first scheme. (c) The approximation of the second scheme.

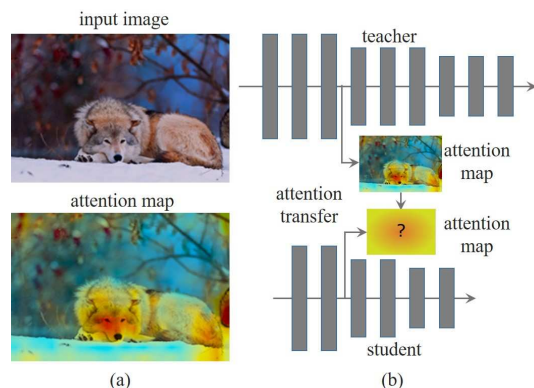


Fig. 23. The application of attention mechanism in teacher-student paradigm transfer learning. (a) The left image is an input and the right image is the corresponding spatial attention map of a CNN model which shows which feature affects the classification decision. (b) Schematic representation of attention transfer. The attention map of the teacher model is used to supervise the training of the student model.

which could approximately replace the original weight matrix. After fine-tuning, the compressed network achieves $5.91 \times$ acceleration on mobile devices with the network, and a less than 1% increase on the top-5 error.

Through optimising the parameter space of fully connected layers, weight factorisation could significantly reduce the memory requirement of DNN models and speed up the inference. However, the effect of the approach for CNN maybe not good, because there is a large amount of convolutional operations in CNN [32]. To solve the problem, Bhattacharya *et al.* propose a convolution kernel separation method, which optimises the convolution filters to significantly reduce convolution operations [214]. The authors verify the effectiveness of the proposed approach on various mobile platforms with popular models, e.g., audio classification and image recognition.

2) *Knowledge Distillation*: Knowledge distillation is based on transfer learning, which trains a neural network of smaller size with the distilled knowledge from a larger model. The large and complex model is called teacher model, whilst the compact model is referred as student model, which takes the benefit of transferring knowledge from the teacher network.

Bucilua *et al.* take the first step towards compressing models with knowledge distillation [409]. They first use a function learned by a high performing model to label pseudo data. Afterwards, the labelled pseudo data is utilised to train a

compact but expressive model. The output of the compact model is compatible with the original high performing model. This work is limited to shallow models. The concept of knowledge distillation is first proposed in [216]. Hinton *et al.* first train a large and complex neural model, which is an ensemble of multiple models. This complex model is the teacher model. Then they design a small and simple student model to learn its knowledge. Specifically, they collect a transfer dataset as the input of the teacher model. The data could be unlabelled data or the original training set of the teacher model. The temperature in softmax is raised to a high value in the teacher model, e.g., 20. Since the soft target of the teacher model is the mean result of multiple components of the teacher model, the training instances are more informative. Therefore, the student model could be trained on much less data than the teacher model. The authors prove the effectiveness on MNIST and speech recognition tasks. Sau *et al.* propose to supervise the training of the student model with multiple teacher models, with the consideration that the distilled knowledge from a single teacher may be limited [219]. They also introduce a noise-based regulariser to improve the health in the performance of the student model.

Romero *et al.* propose FitNet, which extends [216] to create a deeper and lighter student model [217]. Deeper models could better characterise the essence of the data. Both the output of the teacher model and the intermediate representations are used as hints to speed up training of the student model, as well as improve its performance. Opposite to [217], Zagoruyko *et al.* prove that shallow neural networks could also significantly improve the performance of a student model by properly defining attention [218]. Attention is considered as a set of spatial maps that the network focuses the most on in the input to decide the output decision. These maps could be represented as convolutional layers in the network. In the teacher-student paradigm, the spatial attention maps are used to supervise the student model, as shown in Fig. 23.

There are also some efforts focusing on how to design the student model. Crowley *et al.* propose to obtain the student model through replacing the convolutional layers of the teacher model with cheaper alternatives [220]. The new generated student model is then trained under the supervision of the teacher model. Li *et al.* design a framework, named DeepRebirth to merge the consecutive layers without weights, such as pooling and normalisation and convolutional layers vertically or horizontally to compress the model [221]. The

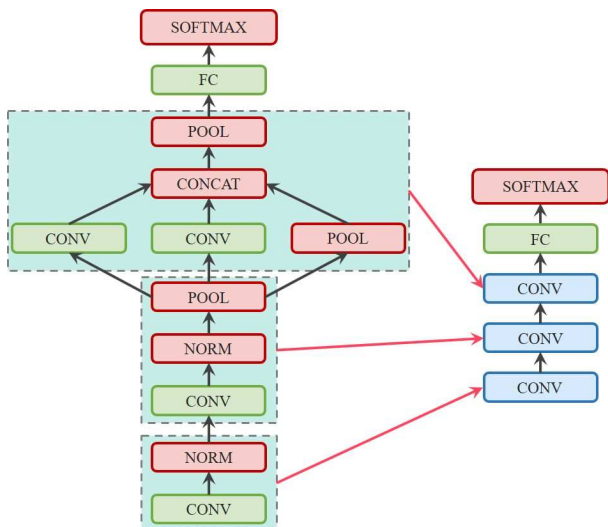


Fig. 24. The illustration of DeepRebirth. The upper model is the teacher model, while the lower is the student model. The highly correlated convolutional layer and non-convolutional layer are merged and become the new convolutional layer of the student model.

newly generated student model learns parameters through layer-wise fine-tuning to minimise the accuracy loss. Fig. 24 presents the framework of DeepRebirth. After compression, GoogLeNet achieves 3x acceleration and 2.5x reduction in runtime memory.

The teacher model is pre-trained in most relevant works. Nevertheless, the teacher model and the student model could be trained in parallel to save time. Zhou *et al.* propose a compression scheme, named Rocket Launching to exploit the simultaneous training of the teacher and student model [222]. During the training, the student model keeps acquiring knowledge learnt by the teacher model through the optimisation of the hint loss. The student model learns both the difference between its output and its target, and the possible path towards the final target learnt by the teacher model. Fig. 25 presents the structure of this framework.

When the teacher model is trained on a dataset concerning with privacy or safety, it is then difficult to train the student model. Lopes *et al.* propose an approach to distill the learned knowledge of the teacher model without accessing the original dataset, which only needs some extra metadata [223]. They first reconstruct the original dataset with the metadata of the teacher model. This step could find the images that best match these given by the network. Then they remove the noise of the image to approximate the activation records through gradients, which could partially reconstruct the original training set of the teacher model.

3) *Compact layer design*: In deep neural networks, if weights end up to be close to 0, the computation is wasted. A fundamental way to solve this problem is to design compact layers in neural networks, which could effectively reduce the consumption of resources, i.e., memories and computation power. Christian *et al.* propose to introduce sparsity and replace the fully connected layers in GoogLeNet [224]. Residual-Net replaces the fully connected layers with global

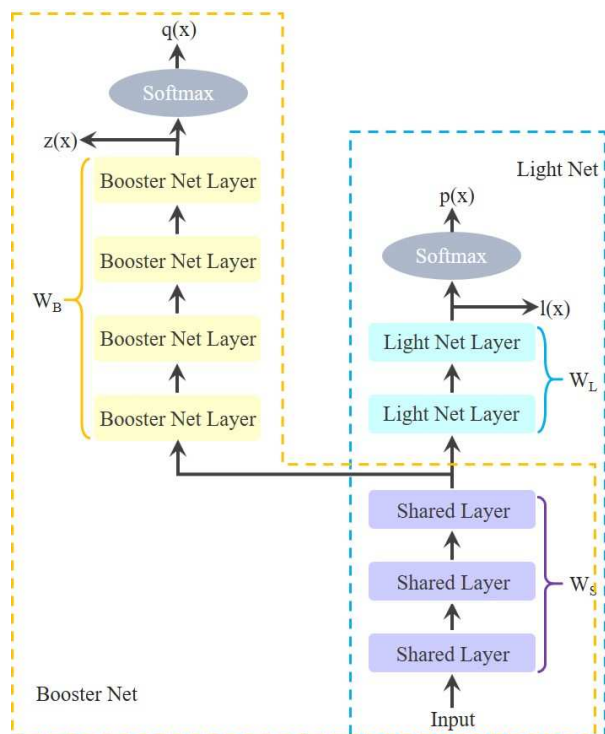


Fig. 25. The structure of Rocket Launching. W_S , W_L , and W_B denotes parameters. $z(x)$ and $l(x)$ represent the weighted sum before the softmax activation. $p(x)$ and $q(x)$ are outputs. Yellow layers are shared by the teacher and student.

average pooling to reduce the resource requirements [225]. Both GoogLeNet and Residual-Net achieve the best performance on multiple benchmarks.

Alex *et al.* propose a compact and lightweight CNN model, named YOLO Nano for image recognition [226]. YOLO Nano is a highly customised model with module-level macro- and micro-architecture. Fig. 26 shows the network architecture of YOLO Nano. There are three modules in YOLO Nano: expansion-projection (EP) macro-architecture, residual projection-expansion-projection (PEP) macro-architecture, and a fully-connected attention (FCA) module. PEP could reduce the architectural and computational complexity whilst preserving model expressiveness. FCA enables better utilisation of available network capacity.

Replacing a big convolution with multiple compact layers could effectively reduce the number of parameters and further reduce computations. Iandola *et al.* propose to compress CNN models with three strategies [227]. First, decomposing 3×3 convolution into 1×1 convolutions, since it has much fewer parameters. Second, cut down input channels in 3×3 convolutions. Third, downsample late to produce big feature maps. The larger feature maps could lead to higher classification accuracy. The first two strategies are used to decrease the quantity of parameters in CNN models and the last one is used to maximise the accuracy of the model. Based on three above mentioned strategies, the authors design SqueezeNet, which can achieve $50\times$ reduction in the number of parameters, whilst remaining the same accuracy as the complete AlexNet.

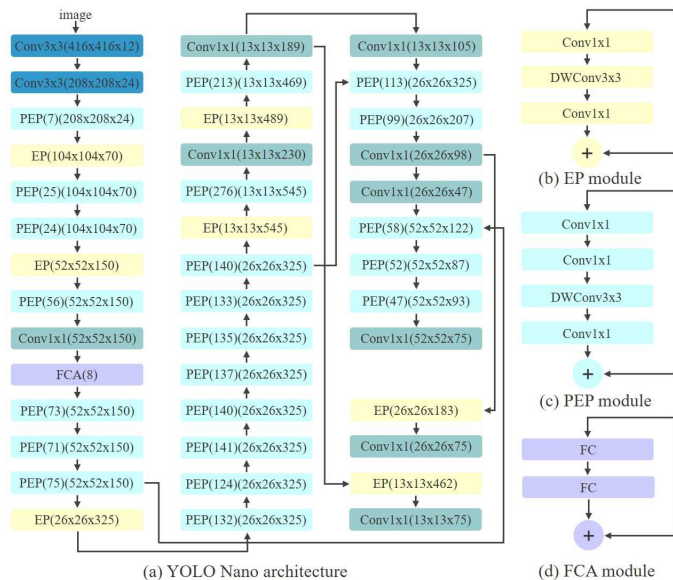


Fig. 26. The architecture of the YOLO Nano network. PEP(x) refers to x channels in PEP, while FCA(x) represents the reduction ratio of x.

Similar approaches are also used in [228]. Shafiee *et al.* modify SqueezeNet for applications with fewer target classes and they propose SqueezeNet v1.1, which could be deployed on edge devices [229]. Yang *et al.* propose to decompose a recognition task into two simple sub-tasks: context recognition and target recognition, and further design a compact model, namely cDeepArch [230]. This approach uses storage resource to trade for computing resources.

Shen *et al.* introduce Compressive Sensing (CS) to jointly modify the input layer and reduce nodes of each layer for CNN models [232]. CS [410] could be used to reduce the dimensionality of the original signal while preserving most of its information. The authors use CS to jointly reduce the dimensions of the input layer whilst extracting most features. The compressed input layer also enables the reduction of the number of parameters.

Besides the above-mentioned works about CNNs, Zhang *et al.* propose a dynamically hierarchy revolution (DirNet) to compress RNNs [231]. In particular, they mine dictionary atoms from original networks to adjust the compression rate with the consideration of different redundancy degrees amongst layers. They then adaptively change the sparsity across the hierarchical layers.

4) *Network pruning*: The main idea of network pruning is to delete unimportant parameters, since not all parameters are important in highly precise deep neural networks. Consequently, connections with less weights are removed, which converts a dense network into a sparse one, as shown in Fig. 27. There are some works which attempt to compress neural networks by network pruning.

The work [411] and [412] have taken the earliest steps towards network pruning. They prune neural networks to eliminate unimportant connections by using Hessian loss function. Experiment results prove the efficiency of pruning

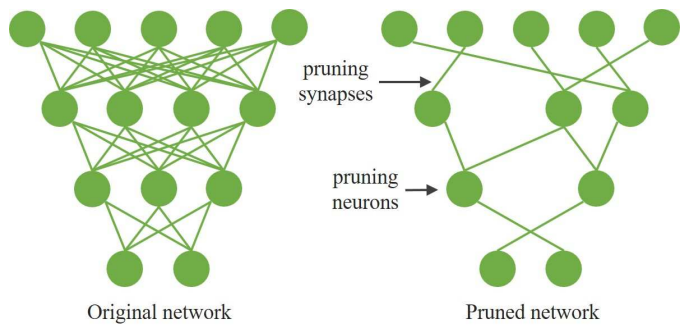


Fig. 27. Illustration of network pruning. Unimportant synapses and neurons would be deleted to generate a sparse network.

methods. Subsequent research focuses on how to prune the networks. Han *et al.* propose to prune networks based on a weight threshold [234]. Practically, they first train a model to learn the weights of each connection. The connections with lower weights than the threshold would then be removed. Afterwards, the network is retrained. The pruning approach is straightforward and simple. A similar approach is also used in [235]. In [235], the authors select and delete neurons of low performance, and then use a width multiplier to expand all layer sizes, which could allocate more resources to neurons of high performance. However, the assumption that connections with lower weights contribute less to the results may destroy the structure of the networks.

Identifying an appropriate threshold to prune neural networks usually takes iteratively trained networks, which consumes a lot of resources and time. Moreover, the threshold is shared by all the layers. Consequently, the pruned configuration maybe not the optimal, comparing with the case of identify thresholds for each layer. To break through these limitations, Manessi *et al.* propose a differentiability-based pruning method to jointly optimise the weights and thresholds for each layer [236]. Specifically, the authors propose a set of differentiable pruning functions and a new regulariser. Pruning could be performed during the back propagation phase, which could effectively reduce the training time.

Molchanov *et al.* propose a new criterion based on the Taylor expansion to identify unimportant neurons in convolutional layers [237]. Specifically, they use the change of cost function to evaluate the result of pruning. They formulate pruning as an optimisation problem, trying to find a weight matrix that minimises the change in cost function. The formulation is approximately converted to its first-degree Taylor polynomial. The gradient and feature map's activation could be easily computed during back-propagation. Therefore, the approach could train the network and prune parameters simultaneously. You *et al.* propose a global filter pruning algorithm, named Gate Decorator, which transforms a CNN module through multiplying its output by the channel-wise scaling factors [238]. If the scaling factor is set to be 0, the corresponding filter would be removed. They also adopt the Taylor expansion to estimate the change of the loss function caused by the changing of the scaling factor. They rank all global filters based on the estimation and prune according to

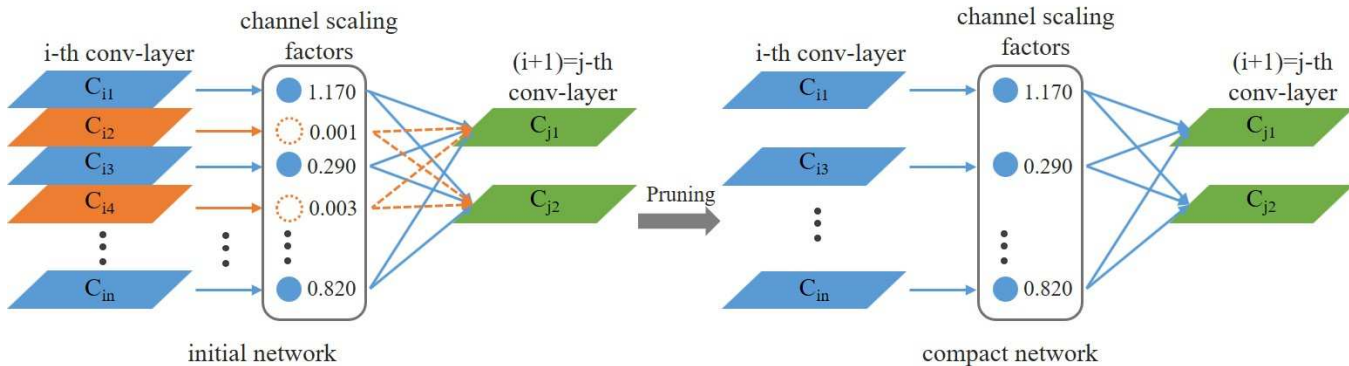


Fig. 28. Each channel is associated with a scaling factor γ in convolutional layers. Then the network is trained to jointly learn weights and scaling factors. After that, the channels with small scaling factors (in orange colour) are pruned, which results in a compact model.

the rank. Compared with [237], [238] does not require special operations or structures.

In addition to minimum weight and cost functions, there are efforts trying to prune with the metric of energy consumption. Yang *et al.* propose an energy-aware pruning algorithm to prune CNNs with the goal of minimising the energy consumption [239]. The authors model the relationship between data sparsity and bit-width reduction through extrapolating the detailed value of consumed energy from hardware measurements. The pruning algorithm identifies the parts of a CNN that consumes the most energy and prunes the weights to maximise energy reduction.

Yao *et al.* propose to minimise the number of non-redundant hidden elements in each layer whilst retaining the accuracy in sensing applications and propose DeepIoT [240]. In DeepIoT, the authors compress neural networks through removing hidden elements. This regularisation approach is called dropout. Each hidden element is dropped with a probability. The dropout probability is initialised with 0.5 for all hidden elements. DeepIoT develops a compressor neural network to learn the optimal dropout probabilities of all elements.

Liu *et al.* propose to identify important channels in CNN and remove unimportant channels to compress networks [241]. Specifically, they introduce a scaling factor γ for each channel. The output \hat{z} (also the input of the next layer) could be formulated as $\hat{z} = \gamma z + \beta$, where z is the input of the current layer and β is min-batch. Afterwards, they jointly train the network weight and scaling factors, with L1 regulation imposed on the latter. Following that, they prune the channels with the small scaling factor γ . Finally, the model is fine-tuned, which achieves a comparable performance with the full network. Fig. 28 presents this slimming process. However, the threshold of the scaling factor is not computed, which requires iterative evaluations to obtain a proper one.

Based on network pruning, the work in [242] investigates the data flow inside computing blocks and develops a data reuse scheme to alleviate the bandwidth burden in convolution layers. The data flow of a convolution layer is regular. If the common data could be reused, it is not necessary to load all data to a new computing block. The data reuse is used to parallelise computing threads and accelerate the inference of

a CNN model.

5) *Parameter quantisation*: A very deep neural network usually involves many layers with millions of parameters, which consumes a large amount of storage and slows down the training procedure. However, highly precise parameters in neural networks are not always necessary in achieving high performance, especially when these highly precise parameters are redundant. It has been proved that only a small number of parameters are enough to reconstruct a complete network [408]. In [408], the authors find that the parameters within one layer could be predicted by 5% of parameters, which means we could compress the model by eliminating redundant parameters. There are some works exploiting parameter quantisation for model compression.

Gong *et al.* propose to use vector quantisation methods to reduce parameters in CNN [248]. Vector quantisation is often used in lossy data compression, which is based on block coding [185]. The main idea of vector quantisation is to divide a set of points into groups, which are represented by their central points. Hence, these points could be denoted with fewer coding bits, which is the basis of compression. In [248], the authors use k-means to cluster parameters and quantise these clusters. They find that this method could achieve 16 – 24 \times compression rate of the parameters with the scarification of no more than 1% of the top-5 accuracy. In addition to k-means, hash method has been utilised in parameter quantisation. In [249], Chen *et al.* propose to use hash functions to cluster connections into different hash buckets uniformly. Connections in the same hash bucket share the same weight. Han *et al.* combine parameter quantisation and pruning to further compress the neural network without compromising the accuracy [250]. Specifically, they first prune the neural network through recognising the important connections through all connections. Unimportant connections are ignored to minimise computation. Then, they quantise the parameters, to save the storage of parameters. After these two steps, the model will be retrained. These remaining connections and parameters could be properly adjusted. Finally, they use Huffman coding to further compress the model. Huffman coding is a prefix coding, which effectively reduces the required storage of data [413]. Fig. 29 presents the three-step compression.

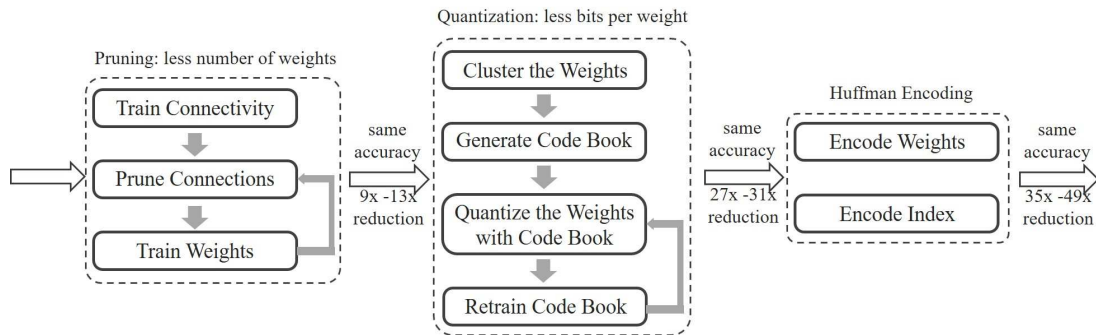


Fig. 29. Illustration of three-stage compression pipeline. First use pruning to reduce the number of weights by 10 \times , then use quantisation to further compress by 27 \times and 31 \times . Finally use Huffman coding to get more compression.

For most CNNs, the fully connected layers consume most storage in neural network. Compressing parameters of fully connected layers could effectively reduce the model size. The convolutional layers consume most of the times during training and inference. Wu *et al.* design Q-CNN to quantise both fully connected layers and convolutional layers to jointly compress and accelerate the neural network [251]. Similar to [248], the authors utilise k-means to optimally cluster parameters in fully connected and convolutional layers. Then, they quantise parameters by minimising the estimated error of response for each layer. They also propose a training scheme to suppress the accumulative error for the quantisation of multiple layers.

Enormous amount of floating point multiplications consumes significant times and computing resources in inference. There are two potential solutions to address this challenge. The first one is to replace floating point with fixed point, and the second one is to reduce the amount of floating point multiplications.

According to the evaluation of Xilinx, fixed point could achieve the same accuracy results as float [414]. Vanhoucke *et al.* evaluate the implementation of fixed point of an 8-bit integer on x86 platform [256]. Specifically, activation and the weights of intermediate layer are quantised into an 8-bit fixed point with the exception of biases that are encoded as 32-bit. The input layer remains floating point to accommodate possible large inputs. Through the quantisation, the total required memory shrinks 3–4 \times . Results show that the quantised model could achieve a 10 \times speedup over an optimised baseline and a 4 \times speedup over an aggressively optimised float point baseline without affecting the accuracy. Similarly, Nasution *et al.* convert floating point to 8 and 16 bits to represent weights and outputs of layers, which lowers the storage to 4.16 \times [258]. Peng *et al.* quantise an image classification CNN model into an 8-bit fixed-point at the cost of 1% accuracy drop [259]. Anwar *et al.* propose to use L2 error minimisation to quantise parameters [260]. They quantise each layer one by one to induce sparsity and retrain the network with the quantised parameters. This approach is evaluated with MNIST and CIRAR-10 dataset. The results shows that the approach could reduce the required memory by 1/10.

In addition to fixed point, posit number could also be utilised to replace floating point numbers to compress neural

networks. Posit number is a unique non-linear numerical system, which could represent all numbers in a dynamic range [415]. The posit number system represents numbers with fewer bits. Float point numbers could be converted into the posit number format to save storage. To learn more about the conversion, readers may refer to [416]. Langroudi *et al.* propose to use the posit number system to compress CNNs with non-uniform data [261]. The weights are converted into posit number format during the reading and writing operations in memory. During the training or inference, when computing operations are required, the number would be converted back to float point. Because this approach only converts the weight between two number systems, no quantisation occurs. The network does not require to be re-trained.

Network Binarisation is an extreme case of weight quantisation. Weight quantisation indicates that all weights are represented by two possible values (e.g., -1 or 1), which could overwhelmingly compress neural networks [252]. For example, the original network requires 32 bits to store one parameter, while in binary connect based network, only 1 bit is enough, which significantly reduces the model size. Another advantage of binary connect is that replacing multiply-accumulate operations by simple accumulations, which could drastically reduce computation in training. Courbariaux *et al.* extend the work [252] further and proposes Binary Neural Network (BNN), which completely changes the computing style of traditional neural networks [253]. Not only the weights, but also the input of each layer is binarised. Hence, during the training, all multiplication operations are replaced by accumulation operations, which drastically improves the power-efficiency. However, substantial experiments indicate that BNN could only achieve good performance on small scale datasets.

Rastegari *et al.* propose a XNOR-net to reduce storage and improve training efficiency, which is different with [253] in the binarisation method and network structure [254]. In Binary-Weight network, all weight values are approximately binarized, e.g., -1 or 1, which reduces the size of network by 32 \times . Convolutions could be finished with only addition and subtraction, which is different with [253]. Hence, the training is speed up 2 \times . With XNOR-net, in addition to weights, the input to convolutional layers are approximately binarised.

TABLE VIII
THE COMPARISON AMONGST STANDARD CONVOLUTION, BINARY-WEIGHT AND XNOR-NET.

	Input	Weight	Convolution operation	Memory saving	Computation saving	Accuracy (imageNet)
Standard Convolution	Real value	Real value	$\times, +, -$	$1\times$	$1\times$	56.7%
Binary-Weight	Real value	Binary value	$+, -$	$\sim 32\times$	$\sim 2\times$	56.8%
XNOR-Net	Binary value	Binary value	XNOR, bitcount	$\sim 32\times$	$\sim 58\times$	44.2%

Moreover, they further simplify the convolution with XNOR operations, which achieves a speed up of $58\times$. The comparison amongst standard convolution, Binary-Weight and XNOR-net is presented as Table. VIII.

Lin *et al.* propose to use binary connect to reduce multiplications in DNN [255]. In the forward pass, the authors stochastically binarise weights by binary connect. Afterwards, they quantise the representations at each layer to replace the remaining multiply operations into bit-shifts. Their results show that there is no loss in accuracy in training and sometimes this approach surprisingly achieves even better performance than standard stochastic gradient descent training.

Soudry *et al.* prove that binary weights and activations could be used in Expectation Backpropagation (EBP) and achieves high performance [262]. This is based on a variational Bayesian approach. The authors test eight binary text classification tasks with EBP-trained multilayer neural networks (MNN). The results show that binary weights always achieve better performance than continuous weights. Esser *et al.* further develop a fully binary network with the same approach to EBP to improve the energy efficiency on neuromorphic chips [263]. They perform the experimentation on the MNIST dataset, and the results show that the method achieves 99.42% accuracy at $108 \mu J$ per image.

6) *Applications*: Some efforts try to use these compression techniques on practical applications and prototypes at the edge, including image analysis [264]–[266], compression service [269], and automotive [267], [268].

Mathur *et al.* develop a wearable camera, called DeepEye, that runs multiple cloud-scale deep learning models at edge provide real-time analysis on the captured images [264]. DeepEye enables the creation of five state-of-the-art image recognition models. After camera captures an image, the image pre-processing component deals with the image according to the adopted deep model. There is a model compression component inside the inference engine, which applies available compression techniques to reduce energy consumption and the running time. Finally, DeepEye use the optimised BLAS library to optimise the numeric operations on hardware.

To correctly identify prescription pills for patients based on their visual appearance, Zeng *et al.* develop MobileDeepPill, a pill image recognition system [265]. The pill image recognition model is based on ImageNet [417]. Fig. 30 presents the architecture of MobileDeepPill. In the training phase, the system first localises and splits the pill image in consumer and pill references. The system then enrich samples through running data augmentation module. Finally, the system imports CNNs as the teacher model to supervise the student model. In the inference phase, the system first processes the pill

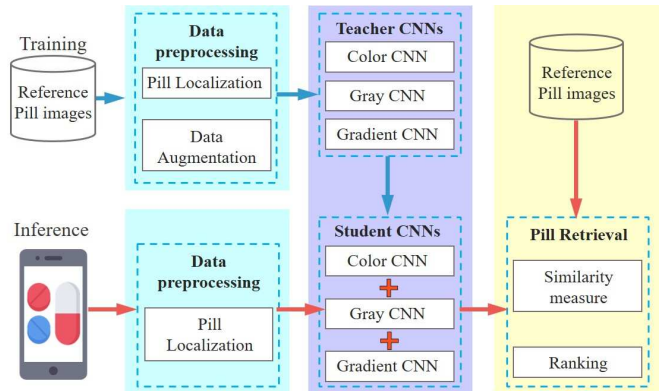


Fig. 30. The architecture of MobileDeepPill. The blue arrows indicates the flow of the training phase, whilst the red arrows indicate the inference phase.

photo and extracts features to perform the student CNNs. As a last step, the system ranks the results according to their possibilities.

Wang *et al.* propose a fast image search framework to implement the content-based image retrieval (CBIR) service from cloud servers to edge devices [266]. Traditional CBIR services are based on the cloud, which suffers from high latency and privacy concerns. The authors propose to reduce the resource requirements of the model and to deploy it on edge devices. For the two components consuming most resources, i.e., object detection and feature extraction, the authors use low-rank approximation to compress these two parts. The compressed model achieves $6.1\times$ speedup for inference.

Liu *et al.* develop an on-demand customised compression system, named AdaDeep [269]. Various kinds of compression approaches could be jointly used in AdaDeep to balance the performance and resource constraints. Specifically, the authors propose a reinforcement learning based optimiser to automatically select the combination of compression approaches to achieve appropriate trade-offs among multiple metrics such as accuracy, storage, and energy consumption.

With growing interests from the automotive industry, various large deep learning models with high accuracy have been implemented in smart vehicles with the assistance of compression techniques. Kim *et al.* develop a DL based object recognition system to recognise vehicles [267]. The vehicle recognition system is based on faster-RCNN. To deploy the system on vehicles, the authors apply network pruning and parameter quantisation to compress the network. Evaluations show that these two compression techniques reduce the network size to 16% and reduce runtime to 64%. Xu *et al.* propose an RNN based driving behaviour analysis system on

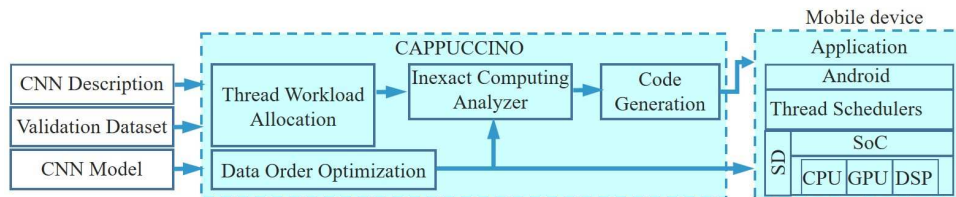


Fig. 31. The architecture of Cappuccino. Thread workload allocation component optimises the workload of each thread. Data order optimisation component converts data format. Inexact computing analyser determines the tradeoff amongst multiple metrics.

vehicles [268]. The system uses the raw data collected by a variety of sensors on vehicles to predict the driving patterns. To deploy the system on automobiles, the authors apply parameter quantisation to reduce the energy consumption and model size. After compression, the system size is reduced to 44 KB and the power overhead is 7.7 mW.

C. Inference Acceleration

The computing capacities of edge devices have been increased and some embedded devices, such as NVIDIA Jetson TX2 [418] could directly perform CNN. However, it is still difficult for most edge devices to directly run large models. Model compression techniques reduce the required resources to create neural network models and facilitate the performance of these models on edge devices. Model acceleration techniques further speed up the performance of the compressed model on edge devices. The main idea of model acceleration in inference is to reduce the run-time of inference on edge devices and realise real-time responses for specific neural network based applications without changing the structure of the trained model. According to acceleration approaches, research works on inference acceleration could be divided into two categories: hardware acceleration and software acceleration. Hardware acceleration methods focuses on parallelising inference tasks to available hardware, such as CPU, GPU, and DSP. Software acceleration method focuses on optimising resource management, pipeline design, and compiler.

1) *Hardware Acceleration*: Recently, mobile devices are becoming increasingly powerful. More and more mobile platforms are equipped with GPUs. Since mobile CPUs are not suitable for the computing of deep neural networks, the embedded GPU could be used to share the computing tasks and accelerate the inference. Table IX summaries existing literature on hardware acceleration.

Alzantot *et al.* evaluate the performance of CNNs and RNNs only on CPU, and compares against the execution in parallel on all available computing resources, e.g., CPU, GPU, DSP, etc. [270]. Results show that the parallel computing paradigm is much faster. Loukadakis *et al.* propose two parallel implementations of VGG-16 network on ODRROID-XU4 board: OpenMP version and OpenCL version [271]. The former parallelises the inference within the CPU, whilst the latter one parallelises within the Mali GPU. These two approaches achieve $2.8\times$ and $11.6\times$ speedup, respectively. Oskouei *et al.* design a mobile GPU-based accelerator for using deep CNN on mobile platforms, which executes inference in parallel

on both CPU and GPU [272]. The accelerator achieves $60\times$ speedup. The authors further develop a GPU-based accelerated library for Android devices, called CNNdroid, which could achieve up to $60\times$ speedup and $130\times$ energy reduction Android platforms [273].

With the consideration that the memory on edge devices are usually not sufficient for neural networks, Tsung *et al.* propose to optimise the flow to accelerate inference [274]. They use a matrix multiplication function to improve the cache hit rate in memory, which indirectly speeds up the execution of the model.

Nvidia has developed a parallelisation framework, named Compute Unified Device Architecture (CUDA) for desktop GPUs to reduce the complexity of neural networks and improve inference speed. For example, in [419], CUDA significantly improves the execution efficiency of RNN on desktop GPUs. Some efforts implement the CUDA framework onto mobile platforms. Rizvi *et al.* propose an approach for image classification on embedded devices based on the CUDA framework [275]. The approach features the most common layers in CNN models, e.g., convolutions, max-pooling, batch-normalisation, and activation functions. General Purpose Computing GPU (GPGPU) is used to speed up the most computation-intensive operations in each layer. The approach is also used to implement an Italian license plate detection and recognition system on tablets [276]. They subsequently introduce matrix multiplication to reduce the computational complexity of convolution in a similar system to achieve real-time object classification on mobile devices [277]. They also apply the approach in a robotic controller system [278].

However, the experiments in [279] show that directly applying CUDA on mobile GPUs may be ineffective, or even deteriorates the performance. Cao *et al.* propose to accelerate RNN on mobile devices based on a parallelisation framework, called RenderScript [279]. RenderScript [280] is a component of the Android platform, which provides an API for hardware acceleration. RenderScript could automatically parallelise the data structure across available GPU cores. The proposed framework could reduce latency by $4\times$.

Motamedi *et al.* implement SqueezeNet on mobile and evaluates the performance on three different Android devices based on RenderScript [281]. Results show that it achieves $310.74\times$ speedup on a Nexus 5. They further develop a general framework, called Cappuccino, for automatic synthesis of efficient inference on edge devices [282]. The structure of Cappuccino is shown as in Fig. 31. There are three inputs for the framework: basic information of the model, model file, and

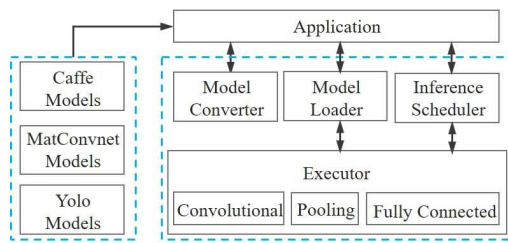


Fig. 32. The structure of Deepsense. The model converter converts the format of the input model, then the model loader loads the model into memory. Inference scheduler is responsible for task scheduling for GPU. The executor runs the allocated tasks on a GPU.

dataset. There are three kinds of parallelisation: kernel-level, filter bank-level, and output-level parallelisation. The thread workload allocation component allocates tasks by using these three kinds of parallelisation. They specially investigate the optimal degree of concurrency for each task, i.e., the number of threads in [283]. The data order optimisation component is used to convert the data format. Cappuccino enables imprecise computing in exchange for high speed and energy efficiency. The inexact computing analyser component is used to analyse the effect of imprecise computing and determine the tradeoff amongst accuracy, speed and energy efficiency.

Huynh *et al.* propose Deepsense, a GPU-based CNN framework to run various CNN models in soft real-time on mobile devices with GPUs [284]. To minimise the latency, Deepsense applies various optimisation strategies, including branch divergence elimination and memory vectorisation. The structure of Deepsense is shown as in Fig. 32. The model converter first converts pre-trained models with different representations into a pre-defined format. Then, the model loader component loads the converted model into memory. When inference starts, the inference scheduler allocates tasks to the GPU sequentially. The executor takes inputted data and the model for executing. During the execution pipeline, CPU is only responsible for padding and intermediate memory allocation, whilst most computing tasks are done by the GPU. The authors further present a demo of the framework in [96] for continuous vision sensing applications on mobile devices.

The heterogeneous multi-core architecture, including CPU and GPU on mobile enables the application of neural networks. By reasonably mapping tasks to cores could improve energy efficiency and inference speed. Taylor *et al.* propose a machine learning based approach to map OpenCL kernels onto proper heterogeneous multi-cores to achieve given objectives, e.g., speedup, energy-efficiency or a tradeoff [285]. The framework first trains the mapping model with the optimisation setting for each objective, then it uses the learned model to schedule OpenCL kernels based on the information of the application.

Rallapalli *et al.* find that the memory of GPUs severely limits the operation of deep CNNs on mobile devices, and proposes to properly allocate part of computation of the fully-connected layers to the CPU [286]. The fully-connected layers are split into several parts, which are executed sequentially. Meanwhile, part of these tasks are loaded into the memory of the CPU for processing. They evaluate the method with an

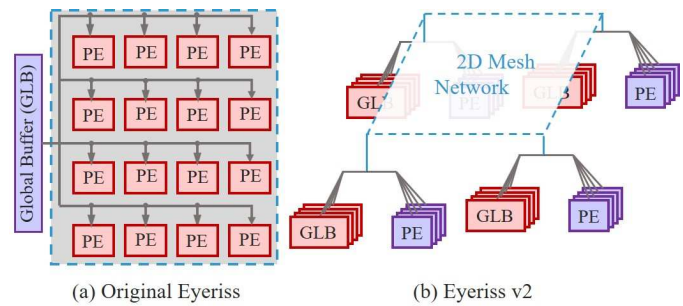


Fig. 33. The comparison between Eyeriss and Eyeriss v2. Both of them are composed of GLB and PE. Eyeriss v2 adopts a hierarchical structure to reduce communication cost.

object detection model, YOLO [287] on Jetson TK1 board and achieve $60\times$ speedup.

In addition to commonly used hardware, i.e., CPUs, mobile GPUs, GPGPU, and DSP, field-programmable gate arrays (FPGAs) could also be used for acceleration. Different from CPUs and GPUs, which run software code, FPGA uses hardware level programming, which means that FPGA is much faster than CPU and GPU. Bettoni *et al.* implement an object recognition CNN model on FPGA via Tiling and Pipelining parallelisation [288]. Ma *et al.* exploit the data reuse and data movement in a convolution loop and proposes to use loop optimisation (including loop unrolling, tiling, and interchange) to accelerate the inference of CNN models in FPGA [289]. A similar approach is also adopted in [290].

Lots of literature focus on developing energy-efficiency DNNs. However, the diversity of DNNs makes them inflexible for hardware [291]. Hence, some researchers attempt to design special accelerating chips to flexibly use DNNs. Chen *et al.* develop an energy-efficient hardware accelerator, called Eyeriss [292]. Eyeriss uses two methods to accelerate the performance of DNNs. The first method is to exploit data reuse to minimise data movement, whilst the second method is to exploit data statistics to avoid unnecessary reads and computations, which improves energy efficiency. Subsequently, they change the structure of the accelerator and propose a new version, Eyeriss v2, to run compact and sparse DNNs [293]. Fig. 33 shows the comparison between Eyeriss and Eyeriss v2. Both of them consist of an array of processing elements (PE) and global buffers (GLB). The main difference is the structure. Eyeriss v2 is hierarchical, in which PEs and GLBs are grouped to reduce communication cost.

2) *Software Acceleration*: Different from hardware acceleration, which depends on the parallelisation of tasks on available hardware, software acceleration mainly focuses on optimising resource management, pipeline design, and compilers. Hardware acceleration methods speed up inference through increasing available computing powers, which usually does not affect the accuracy, whilst software acceleration methods maximise the performance of limited resources for speedup, which may lead to a drop in accuracy with some cases. For example, in [295], the authors sacrifice accuracy for real-time response. Table X summarises existing literature on software acceleration.

TABLE IX
LITERATURE SUMMARY OF HARDWARE ACCELERATION.

Ref.	Model	Executor	Strategy	Object	Performance
[270]	CNN, RNN	CPU, GPU	RenderScript	Feasibility check	3× faster
[271]	VGG-16	CPU, GPU	SIMD	Speed up inference	11.6× faster
[272]	CNN	GPU	SIMD	Speed up inference	60× faster
[273]	CNN	GPU	SIMD	Speed up inference	60× faster 130× energy-saving
[274]	DNN	GPU	Flow optimisation	Enable DNN on mobile device	58× faster 104× energy-saving
[275]	CNN	GPGPU	CUDA	Maximise throughput	50× faster
[276]	DNN	GPU	CUDA	Real-time character detection	250ms per time
[277]	DNN	GPU	Matrix multiplication	Real-time character detection	3× faster
[279]	LSTM	GPU	RenderScript	Rnn RNN on mobile platform	4× reduction on latency
[281]	SqueezeNet	GPU	RenderScript	Acceleration, energy-efficiency	310.74× faster 249.47× energy-saving
[283]	CNN	CPU, GPU, DSP	RenderScript	Optimise thread number	2.37× faster
[282]	CNN	CPU, GPU, DSP	RenderScript	Automatic speedup	272.03× faster at most
[284]	CNN	GPU	Memory vectorisation	Real-time response	VGG-F 361ms
[96]	YOLO	GPU	Tucker decomposition	Real-time response	36% faster
[285]	OpenCL	CPU, GPU	Kernel mapping	Adaptive optimisation	1.2× faster 1.6× energy saving
[286]	YOLO	CPU, GPU	Memory optimisation	Enable CNN on mobile device	0.42s for YOLO
[288]	CNN	FPGA	Tiling, Pipelining	Enable CNN in FPGA	15× faster
[289]	CNN	FPGA	Loop optimisation	Memory and data movement	3.2× faster 23% faster
[290]	CNN	FPGA	Loop optimisation	Improve energy efficiency	9.05× energy-saving
[292]	DNN	Eyeriss	Data reuse	Improve energy efficiency	45% power saving
[293]	DNN	Eyeriss v2	Hierarchical mesh	Hardware processing efficiency	12.6× faster 2.5× energy-saving
[294]	CNN	TPU	Systolic tensor array	Improve systolic array	3.14× faster

Georgiev *et al.* investigate the tradeoff between performance and energy consumption of an audio sensing model on edge devices [296]. Work items need to access different kinds of memory, i.e., global, shared, and private memory. Global memory has the maximum size but minimum speed, whilst private memory is fastest and smallest but exclusive to each work item. Shared memory is between global and private memory. Typical audio sensing models have the maximum parameters, which surpasses the capacity of memory. They use memory access optimisation techniques to speed up the inference, including vectorisation, shared memory sliding window, and tiling.

Lane *et al.* propose DeepX to reduce the resource usage on mobile devices based on two resource management algorithms [297]. The first resource management algorithm is for runtime layer compression. The model compression method discussed in Section V-B could also be used to remove redundancy from original layers. Specifically, they use a SVD-based layer compression technique to simplify the model. The second algorithm is for architecture decomposition, which decomposes the model into blocks that could be performed in parallel. The workflow of DeepX is shown in Fig. 34. They further develop a prototype of DeepX on wearable devices [298]. Subsequently, they develop the DeepX toolkit (DXTK) [299]. A number of pre-trained and compressed deep neural network models are packaged in DXTK. Users could directly use DXTK for specific applications.

Yang *et al.* propose an adaptive software accelerator, Netadpt, which could dynamically speed up the model according to specific metrics [300]. They use empirical measurements on

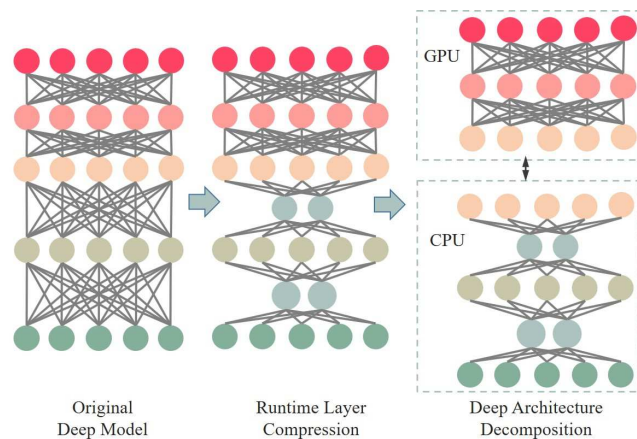


Fig. 34. The workflow of DeepX. Layer compression could reduce the requirement on resource, whilst the architecture decomposition divides the model into multiple blocks that could be performed in parallel.

practical devices to evaluate the performance of the accelerator. Fig. 35 shows the structure of Netadpat. Netadpat adjusts the network according to the given budget, i.e., latency, energy, etc. During iteration, the framework generates many network proposals. Then, Netadpat evaluates these proposals according to direct empirical measurements, and selects one with maximum accuracy. The framework is similar to [269], which caches multiple model compression systems, and compresses the input model according to users' demand.

Ma *et al.* introduce the concept of quality of service (QoS) in model acceleration and develop an accelerator, DeepRT

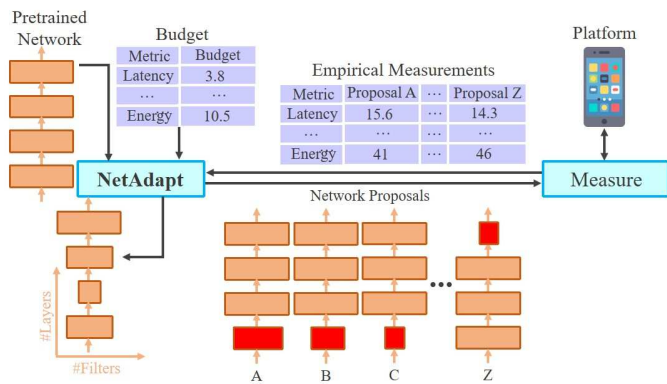


Fig. 35. The structure of Netadapt. Netadapt caches multiple pre-trained models. When requests arrive, Netadapt selects a specific model and adjusts its structure according to the given budget. Then it chooses the best proposal as the accelerating scheme according to empirical measurement.

[301]. The QoS of an accelerated model is defined as a tuple $Q = (d, C)$, where d is a desired response time and C denotes model compression bound. There is a QoS manager component in DeepRT, which controls the system resources to support the QoS during the acceleration.

Liu *et al.* find that fast Fourier transform (FFT) could effectively speed up convolution operation [420]. Abtahi *et al.* apply FFT-based convolution ResNet-20 on NVIDIA Jetson TX1 and evaluates the performance [302]. Results show the inference speed is improved several times. However, FFT-based convolution only works when the convolution kernel is big, e.g., bigger than $9 \times 9 \times 9$. Most models adopt smaller kernels in practice. Hence, there are few applications of FFT-based convolution in practice.

In continuous mobile vision applications, mobile devices are required to deal with continuous videos or images for classification, object recognition, text translation, etc. These continuous videos or images contain large amounts of repeated frames, which are computed through the model again and again during the inference. In such applications, caching mechanisms are promising for acceleration. Xu *et al.* propose CNNCache, a cache-based software accelerator for mobile continuous vision applications, which reuses the computation of similar image regions to avoid unnecessary computation and saves resources on mobile devices [39]. Cavigelli *et al.* present a similar framework, named CBinfer [95]. The difference is that CBinfer considers the threshold of the pixel size when matching frames. However, CBinfer only matches frames of the same position, which may be ineffective in mobile scenarios. [96] also considers reusing the result of the similar input in inference. Different from [39] and [95], the authors extract histogram-based features to match frames, instead of comparing pixels.

VI. EDGE OFFLOADING

Computation is of utmost importance for supporting edge intelligence, which powers the other three components. Most edge devices and edge servers are not as powerful as central servers or computing clusters. Hence, there are two approaches

to enable computation-intensive intelligent applications at the edge: reducing the computational complexity of applications and improving the computing power of edge devices and edge servers. The former approach has been discussed in previous sections. In this section, we focus on the latter approach.

Considering the hardware limitation of edge devices, computation offloading [16], [359], [423]–[425] offers promising approaches to increase computation capability. Literature of this area mainly focuses on designing an optimal offloading strategy to achieve a particular objective, such as latency minimisation, energy-efficiency, and privacy preservation. According to their realisation approaches, these strategies could be divided into five categories: device-to-cloud (D2C), device-to-edge (D2E), device-to-device (D2D), hybrid architecture, and caching.

A. D2C offloading strategy

It consumes considerable computing resources and energy to deal with streamed AI tasks, such as video analysis and continuous speech translation. Most applications, such as Apple Siri and Google Assistant, adopt pure cloud based offloading strategy, in which devices upload input data, e.g., speech or image to cloud server through cellular or WiFi networks. The inference through a giant neural model with high accuracy is done by powerful computers and the results are transmitted back through the same way. There are three main disadvantages in this procedure: (1) mobile devices are required to upload enormous volumes of data to the cloud, which has proved to be the bottleneck of the whole procedure [306]. Such a bottleneck increases users' waiting time; (2) the execution depends on the Internet connectivity; once the device is offline, relative applications could not be used; (3) the uploaded data from mobile devices may contain private information of users, e.g., personal photos, which might be attacked by malicious hackers during the inference on cloud server [426]. There are some efforts trying to solve these problems, which will be discussed next. Table XI summarises existing literature on D2C offloading strategy.

There are usually many layers in a typical deep neural network, which processes the input data layer by layer. The size of intermediate data could be scaled down through layers. Li *et al.* propose a deep neural network layer schedule scheme for the edge environment, leveraging this characteristic of deep neural networks [303]. Fig. 36 shows the structure of neural network layer scheduling-based offloading scheme. Edge devices lacking computing resources, such as IoT devices, first upload the collected data to nearby edge server, which would process the original input data through few low network layers. The generated intermediate data would be uploaded to the cloud server for further processing and eventually output the classification results. The framework is also adopted in [304]. The authors use edge server to pre-process raw data and extract key features.

The model partitioning and layer scheduling could be designed from multiple perspectives, e.g., energy-efficiency, latency, and privacy. Eshratifar *et al.* propose a layer scheduling algorithm from the perspective of energy-efficiency in a

TABLE X
LITERATURE SUMMARY OF SOFTWARE ACCELERATION.

Ref.	Model	Strategy	Object	Performance	Accuracy
[295]	DNN	memory access optimisation	Performance-energy tradeoff	42ms, 83% accuracy	Lossy
[296]	DNN	Resource management	Accelerate inference	6.5× faster, 3 – 4× less energy	Lossless
[297]	DNN	Compression, decomposition	Reduce resource use	5.8× faster	4.9% loss
[298]	DNN	Compression, decomposition	Reduce resource use	5.8× faster	4.9% loss
[299]	DNN	Caching	Reduce resource use	5.8× faster	4.9% Lossless
[300]	NN	Caching	Adaptive speedup	1.7× speedup	4.9% Lossless
[421]	DNN	Caching, model selection	Optimizing DL inference	1.8× speedup	7.52% improvement
[301]	DNN	QoS control	Improve QoS	N/A	N/A
[302]	CNN	FFT-based convolution	Accelerate convolution	10916× faster at most	N/A
[39]	CNN	Cache mechanism	Accelerate inference	20.2% faster	3.51% drop
[95]	CNN	Caching, pixel matching	Accelerate inference	9.1× faster	0.1% drop
[96]	YOLO	Caching, feature extraction	Real-time response	36% faster	3.8%-6.2% drop
[422]	NN	Optimized computing library	Ultra-low-power computing	Up to 63× faster	Negligible loss

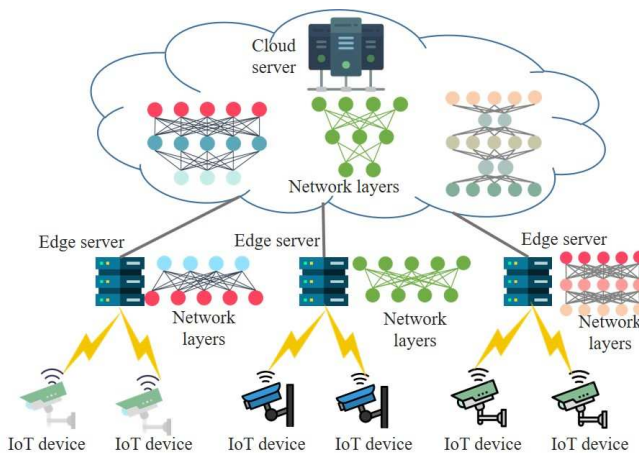


Fig. 36. The structure of neural network layer scheduling-based offloading. IoT devices first upload collected data to edge server, where few neural network layers are deployed. The raw data is first pre-processed on edge servers. Then the intermediate results are uploaded to cloud server for further processing.

similar offloading framework [305]. Kang *et al.* investigate this problem between edge and cloud side [306]. They propose to partition the computing tasks of DNN between local mobile devices and cloud server and design a system, called Neurosurgeon, to intelligently partition DNN based on the prediction of system dynamics. Osia *et al.* consider the layer scheduling from the perspective of privacy preservation [307]. They add a feature extractor module to identify private features from raw data, which will be sent to the cloud for further processing. Analogous approaches are also adopted in [308], [309].

In continuous computer vision analysis, video streams need to be uploaded to the cloud server, which requires a large amount of network resources and consumes battery energy. Ananthanarayanan *et al.* propose a geographically distributed architecture of clouds and edge servers for real-time video analysis [310]. Fixed (e.g., traffic light) and mobile cameras (e.g., event data recorder) upload video streams to available edge servers for pre-processing. The pre-processed data would be further transmitted to a central cloud server in a geographic location for inference. Similarly, Ali *et al.* leverage all available edge resources to pre-process data for large-scale video

stream analytics [311]. Deep learning based video analytic applications contain four stages, including motion detection, frame enhancement, object detection based on shallow networks, and object detection based on deep networks. With the traditional cloud-based approach, these four stages are executed on a cloud server. The authors propose to execute the first two stages locally, which does not require much computation capacity. The output is then transmitted to edge servers for further processing (the third stage). The output is then uploaded to the cloud for final recognition.

Some efforts [427], [428] propose to upload only ‘interesting’ frames to the cloud, which significantly reduces the amount of uploaded data. However, detecting these ‘interesting’ frames also requires intensive computation. Naderiparizi *et al.* develop a framework, Glimpse, to select valid frames by performing coarse visual processing with low energy consumption [312]. Glimpse adopts gating sensors and redesigns the processing pipeline to save energy.

The easiest offloading strategy is to offload the inference task to the cloud when the network condition is good, otherwise perform model compression locally. For example, [313] only considers the network condition when offloading healthcare inference tasks. Cui *et al.* characterise the resource requirements of data processing applications on an edge gateway and cloud server [429]. Hanhirova *et al.* explore and characterise the performance of some CNNs, e.g., object detection and recognition, both on smartphones and cloud server [314]. They find that the latency and throughput are conflicting metrics, in turn, are difficult to be jointly optimised on both mobile devices and cloud server. Some efforts focus on designing an offloading scheme to decide when to use a neural network locally and when to offload the task to the cloud server, instead of always executing on cloud. Considering the capacities of local devices and the network conditions, Qi *et al.* design an adaptive decision scheme to dynamically perform tasks [315]. To enable the local execution, mobile devices adopt compressed models, which achieve lower accuracy than the complete model on the cloud server. If the network condition could not guarantee a real-time response, the inference task would be executed locally. [316] also proposes an offloading decision scheme for the same problem with extra consideration to energy consumption. Ran *et al.* subsequently

TABLE XI
LITERATURE SUMMARY OF D2C OFFLOADING STRATEGY.

Ref.	Model	Execution platform	Focus and problem	Latency	Energy efficiency
[303]	DNN	Edge and cloud	Layer partitioning to reduce uploaded data	0.2s	N/A
[304]	DNN	Edge and cloud	Framework design	$3.23\times$ faster	N/A
[305]	DNN	Edge and cloud	Layer partitioning for energy-efficiency	$3.07\times$ faster	$4.26\times$ higher
[306]	DNN	Edge and cloud	Layer partitioning for latency, energy	$3.1\times$ faster	140.5% higher
[307]	DNN	Local and cloud	Layer partitioning for privacy	N/A	N/A
[308]	DNN	Local and cloud	Feature obfuscation for sensitive data	N/A	N/A
[309]	DNN	Local and cloud	Feature obfuscation for privacy protection	N/A	N/A
[310]	DNN	Edge and cloud	Resource-accuracy tradeoff for real-time performance	N/A	N/A
[311]	CNN	Edge and cloud	Task allocation for QoS	$3.1\times$ faster	140.5% higher
[312]	CV	Edge and cloud	Hardware-based energy and computation efficiency	$10 - 20\times$ faster	$7 - 25\times$ higher
[313]	DNN	Edge or cloud	Offloading decision for acceleration	N/A	N/A
[314]	CNN	Edge or cloud	Performance characterisation and measurement	N/A	N/A
[315]	CNN	Edge or cloud	Latency-accuracy tradeoff for computation-efficiency	N/A	N/A
[316]	NN	Edge or cloud	Multi-objective tradeoff for real-time performance	N/A	N/A
[317]	NN	Edge or cloud	Multi-objective tradeoff for real-time performance	N/A	N/A
[318]	N/A	Local and cloud	Optimal schedule for energy efficiency	Real-time	$1.6 - 3\times$ higher

extend the work with a measurement-driven mathematical framework for achieving a tradeoff between data compression, network condition, energy consumption, latency, and accuracy [317].

Georgiev *et al.* consider a collective offloading scheme for heterogeneous mobile processors and cloud for sensor based applications, which makes best possible use out of different kinds of computing resources on mobile devices, e.g., CPU, GPU, and DSP [318]. They designed a task scheduler running on low-power co-processor unit (LPU) to dynamically restructure and allocate tasks from applications across heterogeneous computing resources based on fluctuations in device and network.

B. D2E offloading strategy

Three main disadvantages with the D2C offloading strategy have been discussed, i.e., latency, wireless network dependency, and privacy concern. Although various solutions have been proposed to alleviate these problems, they do not address these fundamental challenges. Users still need to wait for a long time. Congested wireless networks lead to failed inference. Moreover, the potential risk of private information leakage still exists. Hence, some works try to explore the potential of D2E offloading, which may effectively address these three problems. Edge server refers to the powerful servers (more powerful than ordinary edge devices) that is physically near mobile devices. For example, wearable devices could offload the inference tasks to their connected smartphones. Smartphones could offload computing tasks to cloudlets deployed at roadside. Table XII summarises the existing literature on D2E offloading strategy.

First, we review the works that offload inference tasks to specialised edge servers, e.g., cloudlets and surrogates [430], which refer to infrastructure deployed at edge of the network. There are two general problems that need to be considered in this scenario, including (1) which component of the model could be offloaded to the edge; and (2) which edge server should be selected to offload to. Ra *et al.* develop a framework, named Odessa for interactive perception applications, which enables parallel execution of the inference on local devices

and edge server [319]. They propose a greedy algorithm to partition the model based on the interactive deadlines. The edge servers and edge devices in Odessa are assumed to be fixed, meaning they do not consider problem (2). Streiffer *et al.* appoint an edge server for mobile devices that requests video frame analytics [320]. They evaluate the impact of distance between edge server and mobile devices on latency and packet loss and find that offloading inference tasks to an edge server at a city-scale distance could achieve the similar performance with execution locally on each mobile devices.

Similar to D2C offloading, where the partitioned model layers could be simultaneously deployed on both cloud server and local edge device, the partitioned model layers could also be deployed on edge servers and edge devices. This strategy reduces the transmitted data, which further reduces latency and preserve privacy. Li *et al.* propose Edgent, a device-edge co-inference framework to realise this strategy [19]. The core idea of Edgen is to run computation-intensive layers on powerful edge servers and run the rest layers on device. They also adopt model compression techniques to reduce the model size and further reduce the latency. Similarly, Ko *et al.* propose a model partitioning approach with the consideration of energy efficiency [321]. Due to the difference of available resources between edge devices and edge servers, partitioning the network at a deeper layer would reduce the energy efficiency. Hence, they propose to partition the network at the end of the convolution layers. The output features through the layers on edge device would be compressed before transmitted to edge server to minimise the bandwidth usage.

Some efforts [431]–[433] attempt to encrypt the sensitive data locally before uploading. On cloud side, non-linear layers of a model are converted into linear layers, and then they use homomorphic encryption to execute inference over encrypted input data. This offloading paradigm could also be adopted on edge servers. However, the encryption operation is also computation-intensive. Tian *et al.* propose a private CNN inference framework, LEP-CNN, to offload most inference tasks to edge servers and to avoid privacy breaches [322]. The authors propose an online/offline encryption method to speed up the encryption, which trades offline computation and

TABLE XII
LITERATURE SUMMARY OF D2E OFFLOADING STRATEGY.

Ref.	Model	Problem	Object	Latency	Energy consumption
[319]	Object recognition	Model partitioning	Responsiveness and accuracy	3×	N/A
[320]	Object recognition	Model partitioning	Responsiveness and accuracy	3×	N/A
[19]	DNN	Model partitioning	Reduce latency	100-1000ms	N/A
[321]	DNN	Model partitioning	Energy-efficiency	N/A	4.5× enhanced
[322]	CNN	online/offline encryption	Privacy and latency	35×	95.56% saved
[323]	N/A	Edge server selection	Optimal task migration	N/A	N/A
[324]	DNN	Execution state migration	Computation resource	N/A	N/A

storage for online computation speedup. The execution of the inference over encrypted input data on edge server addresses privacy issues.

Mobility of devices introduces a challenge during the offloading, e.g., in autonomous driving scenarios. Mobile devices may lose the connection with edge server before the inference is done. Hence, selecting an appropriate edge server according to users' mobility pattern is crucial. Zhang *et al.* use reinforcement learning to decide when and which edge server to offload to [323]. A deep Q-network (DQN) based approach is used to automatically learn the optimal offloading scheme from previous experiences. If one mobile device moves away before the edge server finishes the execution of the task, the edge server must drop the task, which wastes the computing resources. Jeong *et al.* propose to move the execution state of the task back to the mobile device from the edge server before the mobile device moves away in the context of web apps [324]. The mobile device continues the execution of the task in this way.

Since the number of edge servers and computation capacity of edge servers are limited, edge devices may compete for resources on edge servers. Hence, proper task scheduling and resource management schemes should be proposed to provide better services at edge. Yi *et al.* propose a latency-aware video edge analytic (LAVEA) system to schedule the tasks from edge devices [328]. For a single edge server, they adopt Johnson's rule [434] to partition the inference task into a two-stage job and prioritise all received offloading requests from edge devices. LAVEA also enables cooperation among different edge servers. They propose three inter-server task scheduling algorithms based on transmission time, scheduling time, and queue length, respectively.

C. D2D offloading strategy

Most neural networks could be executed on mobile devices after compression and achieve a compatible performance. For example, the width-halved GoogLeNet on unmanned aerial vehicles achieves 99% accuracy [435]. Some works consider a more static scenario, where edge devices, such as smart watches are linked to smartphones or home gateways. Wearable devices could offload their model inference tasks to connected powerful devices. There are two kinds of offloading paradigms in this scenario, including binary decision-based offloading and partial offloading. Binary decision offloading refers to executing the task either on a local device or through offloading. This paradigm is similar to D2C offloading. Partial offloading means dividing the inference task into multiple

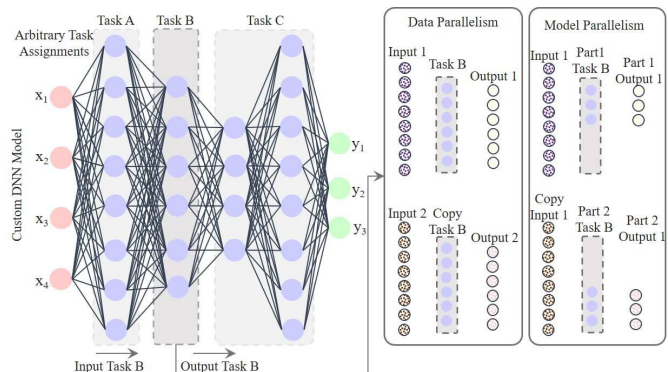


Fig. 37. The parallelism structure of Musical Chair. Task B is a layer-level task, which are further partitioned into two sub-tasks on two devices. These two devices adopt different input to double the system throughput.

sub-tasks and offloading some of them to associated devices. In fact, although the associated devices are more powerful, the performance of the complete offloading is not necessarily better than partial offloading. Because complete offloading is required to transmit complete input data to the associated device, which increases the latency. Table XIII summarises the existing literature for D2D offloading strategy.

Xu *et al.* present CoINF, an offloading framework for wearable devices, which offloads partial inference tasks to associated smartphones [325]. CoINF partitions the model into two sub-models, in which the first sub-model could be executed on the wearable devices, while the second sub-model could be performed on smartphones. They find that the performance of partial offloading outperforms the complete offloading in some scenarios. They further develop a library and provide API for developers. Liu *et al.* also propose EdgeEye, an open source edge computing framework to provide real-time service of video analysis, which provides a task-specific API for developers [326]. Such APIs help developers focus on application logic. Similar methods are also adopted in [436].

If one edge device is not powerful enough to provide real-time response for model inference, a cluster of edge devices could cooperate and help each other to provide enough computation resources. For example, if a camera needs to perform image recognition task, it could partition the CNN model by layers, and transmit the partitioned tasks to other devices nearby. In this scenario, a cluster of edge devices could be organised as a virtual edge server, which could execute inference tasks from both inside and outside of the cluster. Hadidi *et al.* propose Musical Chair, an offloading framework

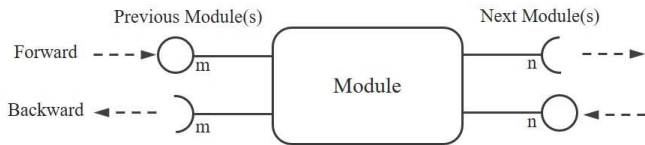


Fig. 38. The illustration of a DIANNE module. Each module has references to its predecessors and successors for feedforward and back-propagation during training.

that harvests available computing resources in an IoT network for cooperation [329]. In Musical Chair, the authors develop data parallelism and model parallelism scheme to speed up the inference. Data parallelism refers to duplicating devices that performs the same task whilst model parallelism is about performing different sub-tasks of a task on different devices. Fig. 37 shows the parallel structure for a layer-level task. Talagala *et al.* use a graph-based overlay network to specify the pipeline dependencies in neural networks and propose a server/agent architecture to schedule computing tasks amongst edge devices in similar scenarios [330].

Coninck *et al.* develop DIANNE, a modular distributed framework, which treats neural network layers as directed graphs [331]. As shown in Fig. 38, each module provides forward and backward methods, corresponding to feedforward and back-propagation, respectively. Each module is a unit deployed on an edge device. There is a job scheduler component, which assigns learning jobs to devices with spare resources. Fukushima *et al.* propose the MicroDeep framework, which assigns neurons of CNN to wireless sensors for image recognition model training [332]. The structure of MicroDeep is similar to DIANNE. Each CNN unit is allocated to a wireless sensor, which executes the training of the unit parameters. In feedforward phase, sensors exchange their output data. Once a sensor receives the necessary input, it executes its unit and broadcasts its output for subsequent layer. If a sensor with an output layer unit obtains its input and ground-truth, it starts the back-propagation phase. They adopt a 2D coordinate based approach to approximately allocate a CNN unit to sensors.

Distributed solo learning enables edge devices or edge servers to train models with local data. Consequently, each model may become local experts that are good at predicting local phenomena. For example, RSUs use local trained models to predict local traffic condition. However, users are interested in the traffic condition of places they plan to visit, in addition to the local traffic condition. Bach *et al.* propose a routing strategy to forward the queries to devices that have the specific knowledge [333]. The strategy is similar to the routing strategy in TCP/IP networks. Each device maintains a routing table to guide the forwarding. The strategy achieve 95% accuracy in their experiments. However, latency is a big problem in such frameworks.

D. Hybrid offloading

The hybrid offloading architecture, also named osmotic computing [437], refers to the computing paradigm that is supported by the seamless collaboration between edge and

cloud computing resources, along with the assistance of data transfer protocols. The hybrid computing architecture takes advantage of cloud, edge, and mobile devices in a holistic manner. There are some efforts focusing on distributing deep learning models in such environments. Table XIV presents a summary of these efforts.

Morshed *et al.* investigate ‘deep osmosis’ and analyses the challenges involved with the holistic distributed deep learning architecture, as well as the data and resource architecture [334]. Teerapittayanon *et al.* propose distributed deep neural networks (DDNNs) based on the holistic computing architecture, which maps sections of a DNN onto a distributed computing hierarchy [335]. All sections are jointly trained in the cloud to minimise communication and resource usage for edge devices. During inference, each edge device performs local computation and then all outputs are aggregated to output the final results.

There is always the risk that the physical nodes i.e., edge devices and edge servers may fail, which results in the failure of DNN units deployed on these physical nodes. Yousefpour *et al.* introduce ‘deepFogGuard’ to make the distributed DNN inference failure-resilient [336]. Similar to residual connections [225], which skips DNN layers to reduce the runtime, ‘deepFogGuard’ skips physical nodes to minimise the impact of failed DNN units. The authors also verify the resilience of ‘deepFogGuard’ on sensing and vision applications.

E. Applications

There exists some works applying the above mentioned offloading strategy to practical applications, such as intelligent transportation [337], smart industry [338], smart city [339], and healthcare [340] [341]. Specifically, in [337], the authors design an edge-centric architecture for intelligent transportation, where roadside smart sensors and vehicles could work as edge servers to provide low latency deep learning based services. [338] proposes a deep learning based classification model to detect defective products in assembly lines, which leverages an edge server to provide computing resources. Tang *et al.* develop a hierarchical distributed framework to support data intensive analytics in smart cities [339]. They develop a pipeline monitoring system for anomaly detection. Edge devices and servers provide computing resources for the execution of these detection models. Liu *et al.* design an edge based food recognition system for dietary assessment, which splits the recognition tasks between nearby edge devices and cloud server to solve the latency and energy consumption problem [340]. Muhammed *et al.* develop a ubiquitous healthcare framework, called UbeHealth, which makes full use of deep learning, big data, and computing resources [341]. They use big data to predict the network traffic, which in turn is used to assist the edge server to make the optimal offloading decision.

VII. FUTURE DIRECTIONS AND OPEN CHALLENGES

We present a thorough and comprehensive survey on the literature surrounding edge intelligence. The benefits of edge intelligence are obvious - it paves the way for the last mile of AI and to provide high-efficient intelligent services for people,

TABLE XIII
LITERATURE SUMMARY OF D2D OFFLOADING STRATEGY.

Ref.	Model	Problem	Object	Latency	Energy consumption
[325]	DNN	Model partition	Acceleration, save energy	23×	85.5% reduction
[326]	DNN	Open source framework	Enable edge inference	N/A	N/A
[327]	AlexNet, VGGNet	Incremental training	Improve accuracy	1.4 – 3.3×	30%-70% saving
[328]	DNN	Task scheduling	Minimise latency	1.2 – 1.7×	N/A
[329]	DNN	Data and task parallelism	Computing power	90×	200× reduction
[330]	N/A	Execution management	ML deployments at edge	N/A	N/A
[331]	AlexNet	Data, model parallelism	Modular architecture	N/A	N/A
[332]	CNN	Neuron assignment	Enable training/inference	N/A	N/A
[333]	Bayesian	Knowledge retrieval	Routing strategy	N/A	N/A

TABLE XIV
LITERATURE SUMMARY OF HYBRID OFFLOADING STRATEGY.

Ref.	Contribution	Solution	Performance
[334]	Challenge analysis in deep osmosis	N/A	N/A
[335]	DDNN frame	Joint training of DNN sections	20× cost reduction
[336]	deepFogGuard	Skip hyperconnections	16% improvement on accuracy

which significantly lessens the dependency on central cloud servers, and can effectively protect data privacy. It is worth recapping that there are still some unsolved open challenges in realising edge intelligence. It is crucial to identify and analyze these challenges and seek for novel theoretical and technical solutions. In this view, we discuss some prominent challenges in edge intelligence along some possible solutions. These challenges include data scarcity at edge, data consistency on edge devices, bad adaptability of statically trained model, privacy and security issues, and Incentive mechanism.

A. Data scarcity at edge

For most machine learning algorithms, especially supervised machine learning, the high performance depends on sufficiently high-quality training instances. However, it often does not work in edge intelligence scenarios, where the collected data is sparse and unlabelled, e.g., in HAR and speech recognition applications. Different from traditional cloud based intelligent services, where the training instances are all gathered in a central database, edge devices use the self-generated data or the data captured from surrounding environments to generate models. High-quality training instances, e.g., good image features are lacking in such datasets. Most existing works ignore this challenge, assuming that the training instances are of good quality. Moreover, the training dataset is often unlabelled. Some works [123], [128] propose to use active learning to solve the problem of unlabelled training instances, which requires manual intervention for annotation. Such an approach could only be used in scenarios with few instances and classifications. Federated learning approaches leverage the decentralised characteristic of data to effectively solve the problem. However, federated learning is only suitable for collaboration training, instead of the solo training needed for personalised models.

We discuss several possible solutions for this problem as follows.

- Adopt shallow models, which could be trained with only a small dataset. Generally, the simpler the machine learning

algorithm is, the better the algorithm will learn from the small datasets. A simple model, e.g., Naive Bayes, linear model, and decision tree, are enough to deal with the problem in some scenarios, compared with complicated models, e.g., neural network, since they are essentially trying to learn less. Hence, adopting an appropriate model should be taken into consideration when dealing with practical problems.

- Incremental learning based methods. Edge devices could re-train a commonly-used pre-trained model in an incremental fashion to accommodate their new data. In such a manner, only few training instances are required to generate a customised model.
- Transfer learning based methods, e.g., few-shot learning. Transfer learning uses the learned knowledge from other models to enhance the performance of a related model, typically avoiding the cold-start problem and reducing the amount of required training data. Hence, transfer learning could be a possible solution, when there is not enough target training data, and the source and target domains have some similarities.
- Data augmentation based methods. Data augmentation enables a model to be more robust by enriching data during the training phase [438]. For example, increasing the number of images without changing the semantic meaning of the labels through flipping, rotation, scaling, translation, cropping, etc. Through the training on augmented data, the network would be invariant to these deformations and have better performance to unseen data.

B. Data consistency on edge devices

Edge intelligence based applications, e.g., speech recognition, activity recognition, emotion recognition, etc., usually collect data from large amounts of sensors distributed at the edge network. Nevertheless, the collected data may not be consistent. Two factors contribute to this problem: different sensing environments, and sensor heterogeneity. The environment (e.g., street and library) and its conditions (e.g.,

raining, windy) add background noise to the collected sensor data, which may have an impact on the model accuracy. The heterogeneity of sensors (e.g., hardware and software) may also cause the unexpected variation in their collected data. For example, different sensors have different sensitivities, sampling rates, and sensing efficiencies. Even the sensor data collected from the same source may vary on different sensors. Consequently, the variation of the data would result in the variation on the model training, e.g., the parameters of features [400], [439], [440]. Such variation is still a challenge for existing sensing applications.

This problem could be solved easily if the model is trained in a centralised manner. The centralised large training set guarantees that the invariant features to the variations could be learned. However, this is not the scope of edge intelligence. Future efforts of this problem should focus on how to block the negative effect of the variation on model accuracy. To this end, two possible research directions maybe considered: data augmentation, and representation learning. Data augmentation could enrich the data during the model training process to enable the model to be more robust to noise. For example, adding various kinds of background noises to block the variation caused by the environments in speech recognition applications on mobile devices. Similarly, the noise caused by the hardware of sensors could also be added to deal with the inconsistency problem. Through the training of the augmented data, the models are more robust with these variations.

Data representation heavily affects the performance of models. Representation learning focuses on learning the representation of data to extract more effective features when building models [441], which could also be used to hide the differences between different hardware. For this problem, if we could make a ‘translation’ on the representations between two sensors which are working on the same data source, the performance of the model would be improved significantly. Hence, representation learning is a promising solution to diminish the impact of data inconsistency. Future efforts could be made on this direction, e.g., design more effective processing pipelines and data transformations.

C. *Bad adaptability of statically trained model*

In most edge intelligence based AI applications, the model is first trained on a central server, then deployed on edge devices. The trained model will not be retrained, once the training procedure is finished. These statically trained models cannot effectively deal with the unknown new data and tasks in unfamiliar environments, which results in low performance and bad user experience. On the other hands, for models trained with a decentralised learning manner, only the local experience is used. Consequently, such models may become experts only in their small local areas. When the serving area broadens, the quality of service decreases.

To cope with this problem, two possible solutions may be considered: lifelong machine learning and knowledge sharing. Lifelong machine learning (LML) [442] is an advanced learning paradigm, which enables continuous knowledge accumulation and self-learning on new tasks. Machines are taught

to learn new knowledge by themselves based on previously learned knowledge, instead of being trained by humans. LML is slightly different from meta learning [443], which enables machines to automatically learn new models. Edge devices with a series of learned tasks could use LML to adapt to changing environments and to deal with unknown data. It is worth recapping that the LML is not primarily designed for edge devices, which means that the machines are expected to be computationally powerful. Accordingly, model design, model compression, and offloading strategies should be also considered if LML is applied.

Knowledge sharing [444] enables the knowledge communication between different edge servers. When there is a task submitted to an edge server that does not have enough knowledge to provide a good service, the server could send knowledge queries to other edge servers. Since the knowledge is allocated on different edge servers, the server with the required knowledge responds to the query and performs the task for users. A knowledge assessment method and knowledge query system are required in such a knowledge sharing paradigm.

D. *Privacy and security issues*

To realise edge intelligence, heterogeneous edge devices and edge servers are required to work collaboratively to provide computing powers. In this procedure, the locally cached data and computing tasks (either training or inference tasks) might be sent to unfamiliar devices for further processing. The data may contain users’ private information, e.g. photos and tokens, which leads to the risk of privacy leakage and attacks from malicious users. If the data is not encrypted, malicious users could easily obtain private information from the data. Some efforts [303], [305], [306], [426] propose to do some preliminary processing locally, which could hide private information and reduce the amount of transmitted data. However, it is still possible to extract private information from the processed data [148]. Moreover, malicious users could also attack and control a device that provides computing power through inserting a virus in the computing tasks. The key challenge is the lack of relevant privacy preserving and security protocols or mechanisms to protect users’ privacy and security from being attacked.

Credit system maybe a possible solution. This is similar with the credit system used in banks, which authenticates each user participated in the system and checks their credit information. Users with bad credit records would be deleted from the system. Consequently, all devices that provide computing powers are credible and all users are safe.

Encryption could be used to protect privacy, which is already applied in some works [55], [132]. However, the encrypted data need to be decrypted before the training or inference tasks are executed, which requires an increase in the amount of computation needed. To cope with the problem, future efforts could pay more attention to homomorphic encryption [151]. Homomorphic encryption refers to an encryption method that allows direct computation on ciphertexts and generate encrypted results. After decryption,

the result is the same as the result achieved by computation on the unencrypted data. Hence, by applying homomorphic encryption, the training or inference task could be directly executed on encrypted data.

E. Incentive mechanism

Data collection and model training/inference are two utmost important steps for edge intelligence. For data collection, it is challenging to ensure the quality and usability of information of the collected data. Data collectors consume their own resources, e.g., battery, bandwidth, and the time to sense and collect data. It is not realistic to assume that all data collectors are willing to contribute, let alone for preprocessing data cleaning, feature extraction and encryption, which further consumes more resources. For model training/inference in a collaborative manner, all participants are required to unselfishly work together for a given task. For example, the architecture proposed in [121] consists of one master and multi workers. Workers recognise objects in a particular mobile visual domain and provides training instances for masters through pipelines. Such architecture works in private scenarios, e.g., at home, where all devices are inherently motivated to collaboratively create a better intelligent model for their master, i.e., their owner. However, it would not work well in public scenarios, where the master initialises a task and allocates sub-tasks to unfamiliar participants. In this context, additional incentive issue arises, which is not typically considered in smart environments where all devices are not under the ownership of a single master. Participants need to be incentivised to perform data collection and task execution.

Reasonable incentive mechanisms should be considered for future efforts. On one hand, participants have different missions, e.g., data collection, data processing, and data analysis, which have different resource consumptions. All participants hope to get as much reward as possible. On the other hand, the operator hopes to achieve the best model accuracy with as a low cost as possible. The challenges of designing the optimal incentive mechanism are how to quantify the workloads of different missions to match corresponding rewards and how to jointly optimise these two conflicting objectives. Future efforts could focus on addressing these challenges.

VIII. CONCLUSIONS

In this paper, we present a thorough and comprehensive survey on the literature surrounding edge intelligence. Specifically, we identify critical components of edge intelligence: edge caching, edge training, edge inference, and edge off-loading. Based on this, we provide a systematic classification of literature by reviewing research achievements for each component and present a systematic taxonomy according to practical problems, adopted techniques, application goals, etc. We compare, discuss and analyse literature in the taxonomy from multi-dimensions, i.e., adopted techniques, objectives, performance, advantages and drawbacks, etc. Moreover, we also discuss important open issues and present possible theoretical research directions. Concerning the era of edge intelligence, We believe that this is only the tip of iceberg.

Along with the explosive development trend of IoT and AI, we expect that more and more research efforts would be carried out to completely realize edge intelligence in the following decades.

REFERENCES

- [1] O. M. Parkhi, A. Vedaldi, and A. Zisserman, "Deep face recognition," 2015.
- [2] Y. Sun, D. Liang, X. Wang, and X. Tang, "Deepid3: Face recognition with very deep neural networks," *arXiv preprint arXiv:1502.00873*, 2015.
- [3] Y. Bengio, R. Ducharme, P. Vincent, and C. Jauvin, "A neural probabilistic language model," *Journal of machine learning research*, vol. 3, no. Feb, pp. 1137–1155, 2003.
- [4] R. Collobert and J. Weston, "A unified architecture for natural language processing: Deep neural networks with multitask learning," in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 160–167.
- [5] A. Kendall and Y. Gal, "What uncertainties do we need in bayesian deep learning for computer vision?" in *Advances in neural information processing systems*, 2017, pp. 5574–5584.
- [6] H. A. Alhajja, S. K. Mustikovela, L. Mescheder, A. Geiger, and C. Rother, "Augmented reality meets deep learning for car instance segmentation in urban scenes," in *British machine vision conference*, vol. 1, 2017, p. 2.
- [7] W. Huang, G. Song, H. Hong, and K. Xie, "Deep architecture for traffic flow prediction: deep belief networks with multitask learning," *IEEE Transactions on Intelligent Transportation Systems*, vol. 15, no. 5, pp. 2191–2201, 2014.
- [8] Y. Lv, Y. Duan, W. Kang, Z. Li, and F.-Y. Wang, "Traffic flow prediction with big data: a deep learning approach," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 2, pp. 865–873, 2014.
- [9] Z. Fang, F. Fei, Y. Fang, C. Lee, N. Xiong, L. Shu, and S. Chen, "Abnormal event detection in crowded scenes based on deep learning," *Multimedia Tools and Applications*, vol. 75, no. 22, pp. 14 617–14 639, 2016.
- [10] C. Potes, S. Parvaneh, A. Rahman, and B. Conroy, "Ensemble of feature-based and deep learning-based classifiers for detection of abnormal heart sounds," in *2016 Computing in Cardiology Conference (CinC)*. IEEE, 2016, pp. 621–624.
- [11] "Cisco visual networking index: Global mobile data traffic forecast update (2017–2022)," <http://www.cisco.com/c/en/us/solutions/collateral/service-provider/visual-networking-index-vni/mobile-white-paper-c11-520862.html>, 2016.
- [12] P. Voigt and A. Von dem Bussche, "The eu general data protection regulation (gdpr)," *A Practical Guide, 1st Ed.*, Cham: Springer International Publishing, 2017.
- [13] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," *IEEE Internet of Things Journal*, vol. 3, no. 5, pp. 637–646, Oct 2016.
- [14] T. X. Tran, A. Hajisami, P. Pandey, and D. Pompili, "Collaborative mobile edge computing in 5g networks: New paradigms, scenarios, and challenges," *IEEE Communications Magazine*, vol. 55, no. 4, pp. 54–61, 2017.
- [15] P. Garcia Lopez, A. Montesor, D. Epema, A. Datta, T. Higashino, A. Iamnitchi, M. Barcellos, P. Felber, and E. Riviere, "Edge-centric computing: Vision and challenges," *ACM SIGCOMM Computer Communication Review*, vol. 45, no. 5, pp. 37–42, 2015.
- [16] Y. C. Hu, M. Patel, D. Sabella, N. Sprecher, and V. Young, "Mobile edge computing a key technology towards 5g," *ETSI white paper*, vol. 11, no. 11, pp. 1–16, 2015.
- [17] F. Bonomi, R. Milito, J. Zhu, and S. Addepalli, "Fog computing and its role in the internet of things," in *Proceedings of the first edition of the MCC workshop on Mobile cloud computing*, 2012, pp. 13–16.
- [18] X. Wang, Y. Han, C. Wang, Q. Zhao, X. Chen, and M. Chen, "In-edge ai: Intelligentizing mobile edge computing, caching and communication by federated learning," *IEEE Network*, vol. 33, no. 5, pp. 156–165, 2019.
- [19] E. Li, Z. Zhou, and X. Chen, "Edge intelligence: On-demand deep learning model co-inference with device-edge synergy," in *Proceedings of the 2018 Workshop on Mobile Edge Communications*. ACM, 2018, pp. 31–36.
- [20] Z. Wang, Y. Cui, and Z. Lai, "A first look at mobile intelligence: Architecture, experimentation and challenges," *IEEE Network*, 2019.

- [21] H. Khelifi, S. Luo, B. Nour, A. Sellami, H. Moun gla, S. H. Ahmed, and M. Guizani, "Bringing deep learning at the edge of information-centric internet of things," *IEEE Communications Letters*, vol. 23, no. 1, pp. 52–55, 2018.
- [22] N. D. Lane and P. Warden, "The deep (learning) transformation of mobile and embedded computing," *Computer*, vol. 51, no. 5, pp. 12–16, 2018.
- [23] F. Chen, Z. Dong, Z. Li, and X. He, "Federated meta-learning for recommendation," *arXiv preprint arXiv:1802.07876*, 2018.
- [24] Y. Chen, J. Wang, C. Yu, W. Gao, and X. Qin, "Fedhealth: A federated transfer learning framework for wearable healthcare," *arXiv preprint arXiv:1907.09173*, 2019.
- [25] E. Peltonen, M. Bennis, M. Capobianco, M. Debbah, A. Ding, F. Gil-Castifeira, M. Jurmu, T. Karvonen, M. Kelanti, A. Kliks *et al.*, "6g white paper on edge intelligence," *arXiv preprint arXiv:2004.14850*, 2020.
- [26] X. Wang, Y. Han, V. C. Leung, D. Niyato, X. Yan, and X. Chen, "Convergence of edge computing and deep learning: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, 2020.
- [27] G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, "Toward an intelligent edge: Wireless communication meets machine learning," *IEEE Communications Magazine*, vol. 58, no. 1, pp. 19–25, 2020.
- [28] S. Yi, Z. Hao, Z. Qin, and Q. Li, "Fog computing: Platform and applications," in *2015 Third IEEE Workshop on Hot Topics in Web Systems and Technologies (HotWeb)*. IEEE, 2015, pp. 73–78.
- [29] K. Ha, Z. Chen, W. Hu, W. Richter, P. Pillai, and M. Satyanarayanan, "Towards wearable cognitive assistance," in *Proceedings of the 12th annual international conference on Mobile systems, applications, and services*, 2014, pp. 68–81.
- [30] 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.
- [31] V. Radu, C. Tong, S. Bhattacharya, N. D. Lane, C. Mascolo, M. K. Marina, and F. Kawsar, "Multimodal deep learning for activity and context recognition," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 1, no. 4, p. 157, 2018.
- [32] N. D. Lane, S. Bhattacharya, P. Georgiev, C. Forlivesi, and F. Kawsar, "An early resource characterization of deep learning on wearables, smartphones and internet-of-things devices," in *Proceedings of the 2015 international workshop on internet of things towards applications*. ACM, 2015, pp. 7–12.
- [33] B. McMahan and D. Ramage, "Federated learning: Collaborative machine learning without centralized training data," *Google Research Blog*, vol. 3, 2017.
- [34] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Transactions on Intelligent Systems and Technology (TIST)*, vol. 10, no. 2, p. 12, 2019.
- [35] 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.
- [36] M. Mohammadi, A. Al-Fuqaha, S. Sorour, and M. Guizani, "Deep learning for iot big data and streaming analytics: A survey," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 4, pp. 2923–2960, 2018.
- [37] C. Zhang, P. Patras, and H. Haddadi, "Deep learning in mobile and wireless networking: A survey," *IEEE Communications Surveys & Tutorials*, 2019.
- [38] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proceedings of the IEEE*, vol. 107, no. 8, pp. 1738–1762, 2019.
- [39] M. Xu, X. Liu, Y. Liu, and F. X. Lin, "Accelerating convolutional neural networks for continuous mobile vision via cache reuse," *arXiv preprint arXiv:1712.01670*, 2017.
- [40] D. Liu and C. Yang, "A learning-based approach to joint content caching and recommendation at base stations," in *2018 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2018, pp. 1–7.
- [41] R. Crane and D. Sornette, "Robust dynamic classes revealed by measuring the response function of a social system," *Proceedings of the National Academy of Sciences*, vol. 105, no. 41, pp. 15 649–15 653, 2008.
- [42] E. Adar, J. Teevan, and S. T. Dumais, "Large scale analysis of web revisitation patterns," in *Proceedings of the SIGCHI conference on Human Factors in Computing Systems*. ACM, 2008, pp. 1197–1206.
- [43] S. Traverso, M. Ahmed, M. Garetto, P. Giaccone, E. Leonardi, and S. Niccolini, "Temporal locality in today's content caching: why it matters and how to model it," *ACM SIGCOMM Computer Communication Review*, vol. 43, no. 5, pp. 5–12, 2013.
- [44] S. Dernbach, N. Taft, J. Kurose, U. Weinsberg, C. Diot, and A. Ashkan, "Cache content-selection policies for streaming video services," in *IEEE INFOCOM 2016-The 35th Annual IEEE International Conference on Computer Communications*. IEEE, 2016, pp. 1–9.
- [45] P. Guo, B. Hu, R. Li, and W. Hu, "Foggycache: Cross-device approximate computation reuse," in *Proceedings of the 24th Annual International Conference on Mobile Computing and Networking*. ACM, 2018, pp. 19–34.
- [46] Google Street View Image API, <https://developers.google.com/maps/documentation/streetview/intro>, 2019.
- [47] R. Huitl, G. Schroth, S. Hilsenbeck, F. Schweiger, and E. Steinbach, "Tumindoor: An extensive image and point cloud dataset for visual indoor localization and mapping," in *2012 19th IEEE International Conference on Image Processing*. IEEE, 2012, pp. 1773–1776.
- [48] D. Liu, B. Chen, C. Yang, and A. F. Molisch, "Caching at the wireless edge: design aspects, challenges, and future directions," *IEEE Communications Magazine*, vol. 54, no. 9, pp. 22–28, 2016.
- [49] T. Li, Z. Xiao, H. M. Georges, Z. Luo, and D. Wang, "Performance analysis of co-and cross-tier device-to-device communication underlying macro-small cell wireless networks," *KSI Transactions on Internet & Information Systems*, vol. 10, no. 4, 2016.
- [50] B. Blaszczyszyn and A. Giovanidis, "Optimal geographic caching in cellular networks," in *2015 IEEE International Conference on Communications (ICC)*. IEEE, 2015, pp. 3358–3363.
- [51] J. G. Andrews, F. Baccelli, and R. K. Ganti, "A tractable approach to coverage and rate in cellular networks," *IEEE Transactions on communications*, vol. 59, no. 11, pp. 3122–3134, 2011.
- [52] C. Roadknight, I. Marshall, and D. Veareer, "File popularity characterisation," *ACM Sigmetrics Performance Evaluation Review*, vol. 27, no. 4, pp. 45–50, 2000.
- [53] H. Ahlehagh and S. Dey, "Video-aware scheduling and caching in the radio access network," *IEEE/ACM Transactions on Networking (TON)*, vol. 22, no. 5, pp. 1444–1462, 2014.
- [54] L. E. Chatzieftheriou, M. Karaliopoulos, and I. Koutsopoulos, "Caching-aware recommendations: Nudging user preferences towards better caching performance," in *IEEE INFOCOM 2017-IEEE Conference on Computer Communications*. IEEE, 2017, pp. 1–9.
- [55] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv preprint arXiv:1610.05492*, 2016.
- [56] S. Li, Y. Cheng, Y. Liu, W. Wang, and T. Chen, "Abnormal client behavior detection in federated learning," *arXiv preprint arXiv:1910.09933*, 2019.
- [57] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *arXiv preprint arXiv:1908.07873*, 2019.
- [58] M. S. H. Abad, E. Ozfatura, D. Gunduz, and O. Ercetin, "Hierarchical federated learning across heterogeneous cellular networks," *arXiv preprint arXiv:1909.02362*, 2019.
- [59] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, "Federated learning in mobile edge networks: A comprehensive survey," *arXiv preprint arXiv:1909.11875*, 2019.
- [60] X. Wang, M. Chen, T. Taleb, A. Ksentini, and V. C. Leung, "Cache in the air: Exploiting content caching and delivery techniques for 5g systems," *IEEE Communications Magazine*, vol. 52, no. 2, pp. 131–139, 2014.
- [61] X. Peng, J.-C. Shen, J. Zhang, and K. B. Letaief, "Backhaul-aware caching placement for wireless networks," in *2015 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2015, pp. 1–6.
- [62] F. Xu, Y. Li, H. Wang, P. Zhang, and D. Jin, "Understanding mobile traffic patterns of large scale cellular towers in urban environment," *IEEE/ACM transactions on networking (TON)*, vol. 25, no. 2, pp. 1147–1161, 2017.
- [63] Z. Xiao, T. Li, W. Ding, D. Wang, and J. Zhang, "Dynamic pci allocation on avoiding handover confusion via cell status prediction in lte heterogeneous small cell networks," *Wireless Communications and Mobile Computing*, vol. 16, no. 14, pp. 1972–1986, 2016.
- [64] Z. Xiao, H. Liu, V. Havvarimana, T. Li, and D. Wang, "Analytical study on multi-tier 5g heterogeneous small cell networks: Coverage performance and energy efficiency," *Sensors*, vol. 16, no. 11, p. 1854, 2016.
- [65] W. K. Lai, C.-S. Shieh, C.-S. Ho, and Y.-R. Chen, "A clustering-based energy saving scheme for dense small cell networks," *IEEE Access*, vol. 7, pp. 2880–2893, 2019.

- [66] Z. Xiao, J. Yu, T. Li, Z. Xiang, D. Wang, and W. Chen, "Resource allocation via hierarchical clustering in dense small cell networks: a correlated equilibrium approach," in *2016 IEEE 27th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*. IEEE, 2016, pp. 1–5.
- [67] K. Hamidouche, W. Saad, M. Debbah, and H. V. Poor, "Mean-field games for distributed caching in ultra-dense small cell networks," in *2016 American Control Conference (ACC)*. IEEE, 2016, pp. 4699–4704.
- [68] N. Zhao, X. Liu, F. R. Yu, M. Li, and V. C. Leung, "Communications, caching, and computing oriented small cell networks with interference alignment," *IEEE Communications Magazine*, vol. 54, no. 9, pp. 29–35, 2016.
- [69] D. Liu and C. Yang, "Cache-enabled heterogeneous cellular networks: Comparison and tradeoffs," in *2016 IEEE International Conference on Communications (ICC)*. IEEE, 2016, pp. 1–6.
- [70] F. Pantisano, M. Bennis, W. Saad, and M. Debbah, "Cache-aware user association in backhaul-constrained small cell networks," in *2014 12th International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt)*. IEEE, 2014, pp. 37–42.
- [71] K. Shanmugam, N. Golrezaei, A. G. Dimakis, A. F. Molisch, and G. Caire, "Femtocaching: Wireless content delivery through distributed caching helpers," *IEEE Transactions on Information Theory*, vol. 59, no. 12, pp. 8402–8413, 2013.
- [72] J. Liu and S. Sun, "Energy efficiency analysis of cache-enabled cooperative dense small cell networks," *IET Communications*, vol. 11, no. 4, pp. 477–482, 2017.
- [73] W. C. Ao and K. Psounis, "Distributed caching and small cell cooperation for fast content delivery," in *Proceedings of the 16th ACM International Symposium on Mobile Ad Hoc Networking and Computing*. ACM, 2015, pp. 127–136.
- [74] —, "Fast content delivery via distributed caching and small cell cooperation," *IEEE Transactions on Mobile Computing*, vol. 17, no. 5, pp. 1048–1061, 2018.
- [75] Z. Chen, J. Lee, T. Q. Quek, and M. Kountouris, "Cooperative caching and transmission design in cluster-centric small cell networks," *IEEE Transactions on Wireless Communications*, vol. 16, no. 5, pp. 3401–3415, 2017.
- [76] S. Krishnan, M. Afshang, and H. S. Dhillon, "Effect of retransmissions on optimal caching in cache-enabled small cell networks," *IEEE Transactions on Vehicular Technology*, vol. 66, no. 12, pp. 11383–11387, 2017.
- [77] Y. Guan, Y. Xiao, H. Feng, C.-C. Shen, and L. J. Cimini, "Mobicacher: Mobility-aware content caching in small-cell networks," in *2014 IEEE Global Communications Conference*. IEEE, 2014, pp. 4537–4542.
- [78] K. Poularakis and L. Tassiulas, "Code, cache and deliver on the move: A novel caching paradigm in hyper-dense small-cell networks," *IEEE Transactions on Mobile Computing*, vol. 16, no. 3, pp. 675–687, 2017.
- [79] E. Ozfatura and D. Gndz, "Mobility and popularity-aware coded small-cell caching," *IEEE Communications Letters*, vol. 22, no. 2, pp. 288–291, 2018.
- [80] Z. Chang, L. Lei, Z. Zhou, S. Mao, and T. Ristaniemi, "Learn to cache: Machine learning for network edge caching in the big data era," *IEEE Wireless Communications*, vol. 25, no. 3, pp. 28–35, 2018.
- [81] M. A. Kader, E. Bastug, M. Bennis, E. Zeydan, A. Karatepe, A. S. Er, and M. Debbah, "Leveraging big data analytics for cache-enabled wireless networks," in *2015 IEEE Globecom Workshops (GC Wkshps)*. IEEE, 2015, pp. 1–6.
- [82] F. Pantisano, M. Bennis, W. Saad, and M. Debbah, "Match to cache: Joint user association and backhaul allocation in cache-aware small cell networks," in *2015 IEEE International Conference on Communications (ICC)*. IEEE, 2015, pp. 3082–3087.
- [83] P. Cheng, C. Ma, M. Ding, Y. Hu, Z. Lin, Y. Li, and B. Vucetic, "Localized small cell caching: A machine learning approach based on rating data," *IEEE Transactions on Communications*, vol. 67, no. 2, pp. 1663–1676, 2019.
- [84] E. Baştuğ, M. Bennis, and M. Debbah, "A transfer learning approach for cache-enabled wireless networks," in *2015 13th International Symposium on Modeling and Optimization in Mobile, Ad Hoc, and Wireless Networks (WiOpt)*. IEEE, 2015, pp. 161–166.
- [85] H. El-Sayed, S. Sankar, M. Prasad, D. Puthal, A. Gupta, M. Mohanty, and C.-T. Lin, "Edge of things: The big picture on the integration of edge, iot and the cloud in a distributed computing environment," *IEEE Access*, vol. 6, pp. 1706–1717, 2018.
- [86] J. Quevedo, D. Corujo, and R. Aguiar, "A case for icn usage in iot environments," in *2014 IEEE Global Communications Conference*. IEEE, 2014, pp. 2770–2775.
- [87] S. K. Sharma and X. Wang, "Live data analytics with collaborative edge and cloud processing in wireless iot networks," *IEEE Access*, vol. 5, pp. 4621–4635, 2017.
- [88] U. Drolia, K. Guo, J. Tan, R. Gandhi, and P. Narasimhan, "Cachier: Edge-caching for recognition applications," in *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2017, pp. 276–286.
- [89] Q. Lv, W. Josephson, Z. Wang, M. Charikar, and K. Li, "Multi-probe lsh: efficient indexing for high-dimensional similarity search," in *Proceedings of the 33rd international conference on Very large data bases*. VLDB Endowment, 2007, pp. 950–961.
- [90] U. Drolia, K. Guo, and P. Narasimhan, "Precog: prefetching for image recognition applications at the edge," in *Proceedings of the Second ACM/IEEE Symposium on Edge Computing*. ACM, 2017, p. 17.
- [91] S. Venugopal, M. Gazzetti, Y. Gkoufas, and K. Katrinis, "Shadow puppets: Cloud-level accurate {AI} inference at the speed and economy of edge," in *{USENIX} Workshop on Hot Topics in Edge Computing (HotEdge 18)*, 2018.
- [92] B. Taylor, V. S. Marco, W. Wolff, Y. Elkhatib, and Z. Wang, "Adaptive deep learning model selection on embedded systems," in *ACM SIGPLAN Notices*, vol. 53, no. 6. ACM, 2018, pp. 31–43.
- [93] J. Zhao, R. Mortier, J. Crowcroft, and L. Wang, "Privacy-preserving machine learning based data analytics on edge devices," in *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*. ACM, 2018, pp. 341–346.
- [94] S. S. Ogden and T. Guo, "{MODI}: Mobile deep inference made efficient by edge computing," in *{USENIX} Workshop on Hot Topics in Edge Computing (HotEdge 18)*, 2018.
- [95] L. Cavigelli and L. Benini, "Cbinfer: Exploiting frame-to-frame locality for faster convolutional network inference on video streams," *IEEE Transactions on Circuits and Systems for Video Technology*, 2019.
- [96] N. L. HUYNH, R. K. Balan, and Y. Lee, "Deepmon-building mobile gpu deep learning models for continuous vision applications," 2017.
- [97] T. Y.-H. Chen, L. Ravindranath, S. Deng, P. Bahl, and H. Balakrishnan, "Glimpse: Continuous, real-time object recognition on mobile devices," in *Proceedings of the 13th ACM Conference on Embedded Networked Sensor Systems*. ACM, 2015, pp. 155–168.
- [98] B. Chen, C. Yang, and Z. Xiong, "Optimal caching and scheduling for cache-enabled d2d communications," *IEEE Communications Letters*, vol. 21, no. 5, pp. 1155–1158, 2017.
- [99] N. Giatsoglou, K. Ntontin, E. Kartsakli, A. Antonopoulos, and C. Verikoukis, "D2d-aware device caching in mmwave-cellular networks," *IEEE Journal on Selected Areas in Communications*, vol. 35, no. 9, pp. 2025–2037, 2017.
- [100] L. Qiu and G. Cao, "Popularity-aware caching increases the capacity of wireless networks," *IEEE Transactions on Mobile Computing*, 2019.
- [101] D. Malak and M. Al-Shalash, "Optimal caching for device-to-device content distribution in 5g networks," in *2014 IEEE Globecom Workshops (GC Wkshps)*. IEEE, 2014, pp. 863–868.
- [102] S. Peng, L. Li, X. Tan, G. Zhao, and Z. Chen, "Optimal caching strategy in device-to-device wireless networks," in *2018 IEEE Wireless Communications and Networking Conference Workshops (WCNCW)*. IEEE, 2018, pp. 78–82.
- [103] M. Ji, G. Caire, and A. F. Molisch, "Optimal throughput-outage trade-off in wireless one-hop caching networks," in *2013 IEEE International Symposium on Information Theory*. IEEE, 2013, pp. 1461–1465.
- [104] —, "The throughput-outage tradeoff of wireless one-hop caching networks," *IEEE Transactions on Information Theory*, vol. 61, no. 12, pp. 6833–6859, 2015.
- [105] B. Chen, C. Yang, and G. Wang, "Cooperative device-to-device communications with caching," in *2016 IEEE 83rd Vehicular Technology Conference (VTC Spring)*. IEEE, 2016, pp. 1–5.
- [106] —, "High-throughput opportunistic cooperative device-to-device communications with caching," *IEEE transactions on vehicular technology*, vol. 66, no. 8, pp. 7527–7539, 2017.
- [107] M. Afshang, H. S. Dhillon, and P. H. J. Chong, "Fundamentals of cluster-centric content placement in cache-enabled device-to-device networks," *IEEE Transactions on Communications*, vol. 64, no. 6, pp. 2511–2526, 2016.
- [108] N. Golrezaei, A. G. Dimakis, and A. F. Molisch, "Wireless device-to-device communications with distributed caching," in *2012 IEEE International Symposium on Information Theory Proceedings*. IEEE, 2012, pp. 2781–2785.
- [109] N. Naderializadeh, D. T. Kao, and A. S. Avestimehr, "How to utilize caching to improve spectral efficiency in device-to-device wireless networks," in *2014 52nd Annual Allerton Conference on Communication, Control, and Computing (Allerton)*. IEEE, 2014, pp. 415–422.

- [110] C. Jarray and A. Giovanidis, "The effects of mobility on the hit performance of cached d2d networks," in *2016 14th international symposium on modeling and optimization in mobile, ad hoc, and wireless networks (WiOpt)*. IEEE, 2016, pp. 1–8.
- [111] E. Bastug, M. Bennis, and M. Debbah, "Living on the edge: The role of proactive caching in 5g wireless networks," *IEEE Communications Magazine*, vol. 52, no. 8, pp. 82–89, 2014.
- [112] M. Newman, *Networks: an introduction*. Oxford university press, 2010.
- [113] B. Bai, L. Wang, Z. Han, W. Chen, and T. Svensson, "Caching based socially-aware d2d communications in wireless content delivery networks: A hypergraph framework," *IEEE Wireless Communications*, vol. 23, no. 4, pp. 74–81, 2016.
- [114] Z. Chen, Y. Liu, B. Zhou, and M. Tao, "Caching incentive design in wireless d2d networks: A stackelberg game approach," in *2016 IEEE International Conference on Communications (ICC)*. IEEE, 2016, pp. 1–6.
- [115] M. Taghizadeh, K. Micinski, S. Biswas, C. Ofria, and E. Torng, "Distributed cooperative caching in social wireless networks," *IEEE Transactions on Mobile Computing*, vol. 12, no. 6, pp. 1037–1053, 2013.
- [116] P. Blasco and D. Gündüz, "Learning-based optimization of cache content in a small cell base station," in *2014 IEEE International Conference on Communications (ICC)*. IEEE, 2014, pp. 1897–1903.
- [117] E. Baştuğ, J.-L. Guénelo, and M. Debbah, "Proactive small cell networks," in *ICT 2013*. IEEE, 2013, pp. 1–5.
- [118] X. Wang, Y. Han, C. Wang, Q. Zhao, X. Chen, and M. Chen, "In-edge ai: Intelligentizing mobile edge computing, caching and communication by federated learning," *IEEE Network*, vol. 33, no. 5, pp. 156–165, 2019.
- [119] Y. Chen, S. Biokaghazadeh, and M. Zhao, "Exploring the capabilities of mobile devices supporting deep learning," in *Proceedings of the 27th International Symposium on High-Performance Parallel and Distributed Computing*. ACM, 2018, pp. 17–18.
- [120] D. Li, T. Salonidis, N. V. Desai, and M. C. Chuah, "Deepcham: Collaborative edge-mediated adaptive deep learning for mobile object recognition," in *2016 IEEE/ACM Symposium on Edge Computing (SEC)*. IEEE, 2016, pp. 64–76.
- [121] Y. Huang, Y. Zhu, X. Fan, X. Ma, F. Wang, J. Liu, Z. Wang, and Y. Cui, "Task scheduling with optimized transmission time in collaborative cloud-edge learning," in *2018 27th International Conference on Computer Communication and Networks (ICCCN)*. IEEE, 2018, pp. 1–9.
- [122] L. Valerio, A. Passarella, and M. Conti, "A communication efficient distributed learning framework for smart environments," *Pervasive and Mobile Computing*, vol. 41, pp. 46–68, 2017.
- [123] T. Xing, S. S. Sandha, B. Balaji, S. Chakraborty, and M. Srivastava, "Enabling edge devices that learn from each other: Cross modal training for activity recognition," in *Proceedings of the 1st International Workshop on Edge Systems, Analytics and Networking*. ACM, 2018, pp. 37–42.
- [124] O. Valery, P. Liu, and J.-J. Wu, "Cpu/gpu collaboration techniques for transfer learning on mobile devices," in *2017 IEEE 23rd International Conference on Parallel and Distributed Systems (ICPADS)*. IEEE, 2017, pp. 477–484.
- [125] —, "Low precision deep learning training on mobile heterogeneous platform," in *2018 26th Euromicro International Conference on Parallel, Distributed and Network-based Processing (PDP)*. IEEE, 2018, pp. 109–117.
- [126] T. Miu, P. Missier, and T. Plötz, "Bootstrapping personalised human activity recognition models using online active learning," in *2015 IEEE International Conference on Computer and Information Technology; Ubiquitous Computing and Communications; Dependable, Autonomic and Secure Computing; Pervasive Intelligence and Computing*. IEEE, 2015, pp. 1138–1147.
- [127] S. Ambrogio, P. Narayanan, H. Tsai, C. Mackin, K. Spoon, A. Chen, A. Fasoli, A. Friz, and G. W. Burr, "Accelerating deep neural networks with analog memory devices," in *2020 2nd IEEE International Conference on Artificial Intelligence Circuits and Systems (AICAS)*, 2020, pp. 149–152.
- [128] F. Shahmohammadi, A. Hosseini, C. E. King, and M. Sarrafzadeh, "Smartwatch based activity recognition using active learning," in *Proceedings of the Second IEEE/ACM International Conference on Connected Health: Applications, Systems and Engineering Technologies*. IEEE Press, 2017, pp. 321–329.
- [129] S. Flutura, A. Seiderer, I. Aslan, C.-T. Dang, R. Schwarz, D. Schiller, and E. André, "Drinkwatch: A mobile wellbeing application based on interactive and cooperative machine learning," in *Proceedings of the 2018 International Conference on Digital Health*. ACM, 2018, pp. 65–74.
- [130] V. Smith, C.-K. Chiang, M. Sanjabi, and A. S. Talwalkar, "Federated multi-task learning," in *Advances in Neural Information Processing Systems*, 2017, pp. 4424–4434.
- [131] H. Zeng and V. Prasanna, "Graphact: Accelerating gcn training on cpu-fpga heterogeneous platforms," in *The 2020 ACM/SIGDA International Symposium on Field-Programmable Gate Arrays*, ser. FPGA 20. New York, NY, USA: Association for Computing Machinery, 2020, p. 255265. [Online]. Available: <https://doi.org/10.1145/3373087.3375312>
- [132] J. Konečný, B. McMahan, and D. Ramage, "Federated optimization: Distributed optimization beyond the datacenter," *arXiv preprint arXiv:1511.03575*, 2015.
- [133] H. B. McMahan, E. Moore, D. Ramage, S. Hampson *et al.*, "Communication-efficient learning of deep networks from decentralized data," *arXiv preprint arXiv:1602.05629*, 2016.
- [134] Y. Zhao, M. Li, L. Lai, N. Suda, D. Civin, and V. Chandra, "Federated learning with non-iid data," *arXiv preprint arXiv:1806.00582*, 2018.
- [135] S. R. Pandey, N. H. Tran, M. Bennis, Y. K. Tun, A. Manzoor, and C. S. Hong, "A crowdsourcing framework for on-device federated learning," *IEEE Transactions on Wireless Communications*, vol. 19, no. 5, pp. 3241–3256, 2020.
- [136] Y. Zhan, P. Li, Z. Qu, D. Zeng, and S. Guo, "A learning-based incentive mechanism for federated learning," *IEEE Internet of Things Journal*, pp. 1–1, 2020.
- [137] S. Wang, T. Tuor, T. Salonidis, K. K. Leung, C. Makaya, T. He, and K. Chan, "Adaptive federated learning in resource constrained edge computing systems," *IEEE Journal on Selected Areas in Communications*, vol. 37, no. 6, pp. 1205–1221, 2019.
- [138] Y. Wang, "Co-op: Cooperative machine learning from mobile devices," 2017.
- [139] K. Yang, T. Jiang, Y. Shi, and Z. Ding, "Federated learning via over-the-air computation," *IEEE Transactions on Wireless Communications*, vol. 19, no. 3, pp. 2022–2035, 2020.
- [140] F. Ang, L. Chen, N. Zhao, Y. Chen, W. Wang, and F. R. Yu, "Robust federated learning with noisy communication," *IEEE Transactions on Communications*, pp. 1–1, 2020.
- [141] M. M. Amiri and D. Gndz, "Federated learning over wireless fading channels," *IEEE Transactions on Wireless Communications*, vol. 19, no. 5, pp. 3546–3557, 2020.
- [142] S. Savazzi, M. Nicoli, and V. Rampa, "Federated learning with cooperating devices: A consensus approach for massive iot networks," *IEEE Internet of Things Journal*, vol. 7, no. 5, pp. 4641–4654, 2020.
- [143] J. Chen, R. Monga, S. Bengio, and R. Jozefowicz, "Revisiting distributed synchronous sgd," 04 2016.
- [144] J. Konečný, "Stochastic, distributed and federated optimization for machine learning," *arXiv preprint arXiv:1707.01155*, 2017.
- [145] Y. Lin, S. Han, H. Mao, Y. Wang, and W. J. Dally, "Deep gradient compression: Reducing the communication bandwidth for distributed training," *arXiv preprint arXiv:1712.01887*, 2017.
- [146] C. Hardy, E. Le Merrer, and B. Sericola, "Distributed deep learning on edge-devices: feasibility via adaptive compression," in *2017 IEEE 16th International Symposium on Network Computing and Applications (NCA)*. IEEE, 2017, pp. 1–8.
- [147] S. Caldas, J. Konečný, H. B. McMahan, and A. Talwalkar, "Expanding the reach of federated learning by reducing client resource requirements," *arXiv preprint arXiv:1812.07210*, 2018.
- [148] W. Yang, S. Wang, J. Hu, G. Zheng, J. Yang, and C. Valli, "Securing deep learning based edge finger-vein biometrics with binary decision diagram," *IEEE Transactions on Industrial Informatics*, 2019.
- [149] K. Bonawitz, V. Ivanov, B. Kreuter, A. Marcedone, H. B. McMahan, S. Patel, D. Ramage, A. Segal, and K. Seth, "Practical secure aggregation for privacy-preserving machine learning," in *Proceedings of the 2017 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2017, pp. 1175–1191.
- [150] Y. Liu, T. Chen, and Q. Yang, "Secure federated transfer learning," *arXiv preprint arXiv:1812.03337*, 2018.
- [151] C. Gentry *et al.*, "Fully homomorphic encryption using ideal lattices," in *Stoc*, vol. 9, no. 2009, 2009, pp. 169–178.
- [152] R. C. Geyer, T. Klein, and M. Nabi, "Differentially private federated learning: A client level perspective," *arXiv preprint arXiv:1712.07557*, 2017.
- [153] K. Wei, J. Li, M. Ding, C. Ma, H. H. Yang, F. Farokhi, S. Jin, T. Q. S. Quek, and H. V. Poor, "Federated learning with differential privacy: Algorithms and performance analysis," *IEEE Transactions on Information Forensics and Security*, pp. 1–1, 2020.

- [154] C. Dwork, "Differential privacy," *Encyclopedia of Cryptography and Security*, pp. 338–340, 2011.
- [155] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, "Deep learning with differential privacy," in *Proceedings of the 2016 ACM SIGSAC Conference on Computer and Communications Security*. ACM, 2016, pp. 308–318.
- [156] H. B. McMahan, D. Ramage, K. Talwar, and L. Zhang, "Learning differentially private recurrent language models," *arXiv preprint arXiv:1710.06963*, 2017.
- [157] H. H. Zhuo, W. Feng, Q. Xu, Q. Yang, and Y. Lin, "Federated reinforcement learning," *arXiv preprint arXiv:1901.08277*, 2019.
- [158] K. Cheng, T. Fan, Y. Jin, Y. Liu, T. Chen, and Q. Yang, "Secureboost: A lossless federated learning framework," *arXiv preprint arXiv:1901.08755*, 2019.
- [159] B. Biggio, B. Nelson, and P. Laskov, "Poisoning attacks against support vector machines," in *Proceedings of the 29th International Conference on Machine Learning*, ser. ICML12. Madison, WI, USA: Omnipress, 2012, p. 14671474.
- [160] J. Steinhardt, P. W. W. Koh, and P. S. Liang, "Certified defenses for data poisoning attacks," in *Advances in neural information processing systems*, 2017, pp. 3517–3529.
- [161] C. Fung, C. J. Yoon, and I. Beschastnikh, "Mitigating sybils in federated learning poisoning," *arXiv preprint arXiv:1808.04866*, 2018.
- [162] J. R. Douceur, "The sybil attack," in *International workshop on peer-to-peer systems*. Springer, 2002, pp. 251–260.
- [163] P. Blanchard, R. Guerraoui, J. Stainer *et al.*, "Machine learning with adversaries: Byzantine tolerant gradient descent," in *Advances in Neural Information Processing Systems*, 2017, pp. 119–129.
- [164] Y. Chen, L. Su, and J. Xu, "Distributed statistical machine learning in adversarial settings: Byzantine gradient descent," *Proceedings of the ACM on Measurement and Analysis of Computing Systems*, vol. 1, no. 2, p. 44, 2017.
- [165] D. Yin, Y. Chen, R. Kannan, and P. Bartlett, "Byzantine-robust distributed learning: Towards optimal statistical rates," in *Proceedings of the 35th International Conference on Machine Learning*, ser. Proceedings of Machine Learning Research, J. Dy and A. Krause, Eds., vol. 80. Stockholm, Sweden: PMLR, 10–15 Jul 2018, pp. 5650–5659.
- [166] X. Chen, C. Liu, B. Li, K. Lu, and D. Song, "Targeted backdoor attacks on deep learning systems using data poisoning," *arXiv preprint arXiv:1712.05526*, 2017.
- [167] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," *arXiv preprint arXiv:1807.00459*, 2018.
- [168] A. N. Bhagoji, S. Chakraborty, P. Mittal, and S. Calo, "Analyzing federated learning through an adversarial lens," *arXiv preprint arXiv:1811.12470*, 2018.
- [169] Y. Gal and Z. Ghahramani, "Dropout as a bayesian approximation: Representing model uncertainty in deep learning," in *international conference on machine learning*, 2016, pp. 1050–1059.
- [170] S. Yao, Y. Zhao, H. Shao, A. Zhang, C. Zhang, S. Li, and T. Abdelzaher, "Rdeepsense: Reliable deep mobile computing models with uncertainty estimations," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 1, no. 4, p. 173, 2018.
- [171] K. Bonawitz, H. Eichner, W. Grieskamp, D. Huba, A. Ingerman, V. Ivanov, C. Kiddon, J. Konecny, S. Mazzocchi, H. B. McMahan *et al.*, "Towards federated learning at scale: System design," *arXiv preprint arXiv:1902.01046*, 2019.
- [172] T. Nishio and R. Yonetani, "Client selection for federated learning with heterogeneous resources in mobile edge," in *ICC 2019-2019 IEEE International Conference on Communications (ICC)*. IEEE, 2019, pp. 1–7.
- [173] A. Hard, K. Rao, R. Mathews, F. Beaufays, S. Augenstein, H. Eichner, C. Kiddon, and D. Ramage, "Federated learning for mobile keyboard prediction," *arXiv preprint arXiv:1811.03604*, 2018.
- [174] M. Chen, R. Mathews, T. Ouyang, and F. Beaufays, "Federated learning of out-of-vocabulary words," *arXiv preprint arXiv:1903.10635*, 2019.
- [175] T. Yang, G. Andrew, H. Eichner, H. Sun, W. Li, N. Kong, D. Ramage, and F. Beaufays, "Applied federated learning: Improving google keyboard query suggestions," *arXiv preprint arXiv:1812.02903*, 2018.
- [176] S. Ramaswamy, R. Mathews, K. Rao, and F. Beaufays, "Federated learning for emoji prediction in a mobile keyboard," *arXiv preprint arXiv:1906.04329*, 2019.
- [177] M. J. Sheller, G. A. Reina, B. Edwards, J. Martin, and S. Bakas, "Multi-institutional deep learning modeling without sharing patient data: A feasibility study on brain tumor segmentation," in *International MICCAI Brainlesion Workshop*. Springer, 2018, pp. 92–104.
- [178] A. G. Roy, S. Siddiqui, S. Pölsterl, N. Navab, and C. Wachinger, "Braitentor: A peer-to-peer environment for decentralized federated learning," *arXiv preprint arXiv:1905.06731*, 2019.
- [179] S. Samarakoon, M. Bennis, W. Saad, and M. Debbah, "Distributed federated learning for ultra-reliable low-latency vehicular communications," *IEEE Transactions on Communications*, 2019.
- [180] T. D. Nguyen, S. Marchal, M. Miettinen, N. Asokan, and A. Sadeghi, "Diot: A self-learning system for detecting compromised iot devices," *CoRR*, vol. abs/1804.07474, 2018.
- [181] S. Chen, Y. Liu, X. Gao, and Z. Han, "Mobilefacenet: Efficient cnns for accurate real-time face verification on mobile devices," in *Chinese Conference on Biometric Recognition*. Springer, 2018, pp. 428–438.
- [182] C. N. Duong, K. G. Quach, N. Le, N. Nguyen, and K. Luu, "Mobiface: A lightweight deep learning face recognition on mobile devices," *arXiv preprint arXiv:1811.11080*, 2018.
- [183] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, "Learning transferable architectures for scalable image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 8697–8710.
- [184] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, "Regularized evolution for image classifier architecture search," in *Proceedings of the aaai conference on artificial intelligence*, vol. 33, 2019, pp. 4780–4789.
- [185] C. Cortes, X. Gonzalvo, V. Kuznetsov, M. Mohri, and S. Yang, "Adanet: Adaptive structural learning of artificial neural networks," in *Proceedings of the 34th International Conference on Machine Learning-Volume 70*. JMLR. org, 2017, pp. 874–883.
- [186] 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 *The IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2019.
- [187] H. Liu, K. Simonyan, and Y. Yang, "Darts: Differentiable architecture search," *arXiv preprint arXiv:1806.09055*, 2018.
- [188] 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.
- [189] Y. Zhang, N. Suda, L. Lai, and V. Chandra, "Hello edge: Keyword spotting on microcontrollers," *arXiv preprint arXiv:1711.07128*, 2017.
- [190] D. Wofk, F. Ma, T.-J. Yang, S. Karaman, and V. Sze, "Fastdepth: Fast monocular depth estimation on embedded systems," in *2019 International Conference on Robotics and Automation (ICRA)*. IEEE, 2019, pp. 6101–6108.
- [191] 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.
- [192] Z. Qin, Z. Zhang, S. Zhang, H. Yu, and Y. Peng, "Merging-and-evolution networks for mobile vision applications," *IEEE Access*, vol. 6, pp. 31 294–31 306, 2018.
- [193] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, and L.-C. Chen, "Mobilenetv2: Inverted residuals and linear bottlenecks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 4510–4520.
- [194] S. Bhattacharya and N. D. Lane, "From smart to deep: Robust activity recognition on smartwatches using deep learning," in *2016 IEEE International Conference on Pervasive Computing and Communication Workshops (PerCom Workshops)*. IEEE, 2016, pp. 1–6.
- [195] B. Almaslakh, J. Al Muhtadi, and A. M. Artoli, "A robust convolutional neural network for online smartphone-based human activity recognition," *Journal of Intelligent & Fuzzy Systems*, no. Preprint, pp. 1–12, 2018.
- [196] B. Almaslakh, A. 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.
- [197] 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*. ACM, 2018, pp. 31–36.
- [198] 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*. ACM, 2016, pp. 185–188.
- [199] 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.

- [200] X. Bo, C. Poellabauer, M. K. O'Brien, C. K. Mummisetti, and A. Jayaraman, "Detecting label errors in crowd-sourced smartphone sensor data," in *2018 International Workshop on Social Sensing (SocialSens)*. IEEE, 2018, pp. 20–25.
- [201] S. Yao, S. Hu, Y. Zhao, A. Zhang, and T. Abdelzaher, "Deepsense: A unified deep learning framework for time-series mobile sensing data processing," in *Proceedings of the 26th International Conference on World Wide Web*. International World Wide Web Conferences Steering Committee, 2017, pp. 351–360.
- [202] S. Yao, Y. Zhao, S. Hu, and T. Abdelzaher, "Qualitydeepsense: Quality-aware deep learning framework for internet of things applications with sensor-temporal attention," in *Proceedings of the 2nd International Workshop on Embedded and Mobile Deep Learning*. ACM, 2018, pp. 42–47.
- [203] C. Streiffer, R. Raghavendra, T. Benson, and M. Srivatsa, "Darnet: a deep learning solution for distracted driving detection," in *Proceedings of the 18th ACM/IFIP/USENIX Middleware Conference: Industrial Track*. ACM, 2017, pp. 22–28.
- [204] L. Liu, C. Karatas, H. Li, S. Tan, M. Gruteser, J. Yang, Y. Chen, and R. P. Martin, "Toward detection of unsafe driving with wearables," in *Proceedings of the 2015 workshop on Wearable Systems and Applications*. ACM, 2015, pp. 27–32.
- [205] C. Bo, X. Jian, X.-Y. Li, X. Mao, Y. Wang, and F. Li, "You're driving and texting: detecting drivers using personal smart phones by leveraging inertial sensors," in *Proceedings of the 19th annual international conference on Mobile computing & networking*. ACM, 2013, pp. 199–202.
- [206] J. Yang, S. Sidhom, G. Chandrasekaran, T. Vu, H. Liu, N. Cekan, Y. Chen, M. Gruteser, and R. P. Martin, "Detecting driver phone use leveraging car speakers," in *Proceedings of the 17th annual international conference on Mobile computing and networking*. ACM, 2011, pp. 97–108.
- [207] 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*. ACM, 2015, pp. 283–294.
- [208] P. Georgiev, S. Bhattacharya, N. D. Lane, and C. Mascolo, "Low-resource multi-task audio sensing for mobile and embedded devices via shared deep neural network representations," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 1, no. 3, p. 50, 2017.
- [209] M. Jaderberg, A. Vedaldi, and A. Zisserman, "Speeding up convolutional neural networks with low rank expansions," *arXiv preprint arXiv:1405.3866*, 2014.
- [210] E. L. Denton, W. Zaremba, J. Bruna, Y. LeCun, and R. Fergus, "Exploiting linear structure within convolutional networks for efficient evaluation," in *Advances in neural information processing systems*, 2014, pp. 1269–1277.
- [211] P. Maji, D. Bates, A. Chadwick, and R. Mullins, "Adapt: optimizing cnn inference on iot and mobile devices using approximately separable 1-d kernels," in *Proceedings of the 1st International Conference on Internet of Things and Machine Learning*. ACM, 2017, p. 43.
- [212] Y.-D. Kim, E. Park, S. Yoo, T. Choi, L. Yang, and D. Shin, "Compression of deep convolutional neural networks for fast and low power mobile applications," *arXiv preprint arXiv:1511.06530*, 2015.
- [213] P. Wang and J. Cheng, "Accelerating convolutional neural networks for mobile applications," in *Proceedings of the 24th ACM international conference on Multimedia*. ACM, 2016, pp. 541–545.
- [214] S. Bhattacharya and N. D. Lane, "Sparsification and separation of deep learning layers for constrained resource inference on wearables," in *Proceedings of the 14th ACM Conference on Embedded Network Sensor Systems CD-ROM*. ACM, 2016, pp. 176–189.
- [215] C. Bucilu, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2006, pp. 535–541.
- [216] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [217] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "Fitnets: Hints for thin deep nets," *arXiv preprint arXiv:1412.6550*, 2014.
- [218] N. Komodakis and S. Zagoruyko, "Paying more attention to attention: improving the performance of convolutional neural networks via attention transfer," in *ICLR*, Paris, France, Jun. 2017. [Online]. Available: <https://hal-enpc.archives-ouvertes.fr/hal-01832769>
- [219] B. B. Sau and V. N. Balasubramanian, "Deep model compression: Distilling knowledge from noisy teachers," *arXiv preprint arXiv:1610.09650*, 2016.
- [220] E. J. Crowley, G. Gray, and A. J. Storkey, "Moonshine: Distilling with cheap convolutions," in *Advances in Neural Information Processing Systems*, 2018, pp. 2888–2898.
- [221] 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.
- [222] G. Zhou, Y. Fan, R. Cui, W. Bian, X. Zhu, and K. Gai, "Rocket launching: A universal and efficient framework for training well-performing light net," in *Thirty-Second AAAI Conference on Artificial Intelligence*, 2018.
- [223] R. G. Lopes, S. Fenu, and T. Starner, "Data-free knowledge distillation for deep neural networks," *arXiv preprint arXiv:1710.07535*, 2017.
- [224] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2015, pp. 1–9.
- [225] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 770–778.
- [226] A. Wong, M. Famuori, M. J. Shafiee, F. Li, B. Chwyl, and J. Chung, "Yolo nano: a highly compact you only look once convolutional neural network for object detection," 2019.
- [227] 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.
- [228] C. Szegedy, S. Ioffe, V. Vanhoucke, and A. A. Alemi, "Inception-v4, inception-resnet and the impact of residual connections on learning," in *Thirty-First AAAI Conference on Artificial Intelligence*, 2017.
- [229] M. J. Shafiee, F. Li, B. Chwyl, and A. Wong, "Squishednets: Squishing squeezenet further for edge device scenarios via deep evolutionary synthesis," *arXiv preprint arXiv:1711.07459*, 2017.
- [230] 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*, 2019.
- [231] J. Zhang, X. Wang, D. Li, and Y. Wang, "Dynamically hierarchy revolution: dirnet for compressing recurrent neural network on mobile devices," *arXiv preprint arXiv:1806.01248*, 2018.
- [232] Y. Shen, T. Han, Q. Yang, X. Yang, Y. Wang, F. Li, and H. Wen, "Cs-cnn: Enabling robust and efficient convolutional neural networks inference for internet-of-things applications," *IEEE Access*, vol. 6, pp. 13 439–13 448, 2018.
- [233] J. Guo and M. Potkonjak, "Pruning filters and classes: Towards on-device customization of convolutional neural networks," in *Proceedings of the 1st International Workshop on Deep Learning for Mobile Systems and Applications*. ACM, 2017, pp. 13–17.
- [234] S. Han, J. Pool, J. Tran, and W. Dally, "Learning both weights and connections for efficient neural network," in *Advances in neural information processing systems*, 2015, pp. 1135–1143.
- [235] A. Gordon, E. Eban, O. Nachum, B. Chen, H. Wu, T.-J. Yang, and E. Choi, "Morphnet: Fast & simple resource-constrained structure learning of deep networks," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 1586–1595.
- [236] F. Manessi, A. Rozza, S. Bianco, P. Napolitano, and R. Schettini, "Automated pruning for deep neural network compression," in *2018 24th International Conference on Pattern Recognition (ICPR)*. IEEE, 2018, pp. 657–664.
- [237] P. Molchanov, S. Tyree, T. Karras, T. Aila, and J. Kautz, "Pruning convolutional neural networks for resource efficient transfer learning," *arXiv preprint arXiv:1611.06440*, vol. 3, 2016.
- [238] Z. You, K. Yan, J. Ye, M. Ma, and P. Wang, "Gate decorator: Global filter pruning method for accelerating deep convolutional neural networks," in *Advances in Neural Information Processing Systems 32*. Curran Associates, Inc., 2019, pp. 2133–2144.
- [239] T.-J. Yang, Y.-H. Chen, and V. Sze, "Designing energy-efficient convolutional neural networks using energy-aware pruning," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 5687–5695.
- [240] 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*. ACM, 2017, p. 4.
- [241] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, and C. Zhang, "Learning efficient convolutional networks through network slimming," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2736–2744.

- [242] C.-F. Chen, G. G. Lee, V. Sritapan, and C.-Y. Lin, "Deep convolutional neural network on ios mobile devices," in *2016 IEEE International Workshop on Signal Processing Systems (SiPS)*. IEEE, 2016, pp. 130–135.
- [243] J.-H. Luo and J. Wu, "Autopruner: An end-to-end trainable filter pruning method for efficient deep model inference," *Pattern Recognition*, p. 107461, 2020. [Online]. Available: <http://www.sciencedirect.com/science/article/pii/S0031320320302648>
- [244] J. Guo, W. Zhang, W. Ouyang, and D. Xu, "Model compression using progressive channel pruning," *IEEE Transactions on Circuits and Systems for Video Technology*, pp. 1–1, 2020.
- [245] O. Oyedotun, D. Aouada, and B. Ottersten, "Structured compression of deep neural networks with debiased elastic group lasso," in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, March 2020.
- [246] P. Singh, V. K. Verma, P. Rai, and V. Nambodiri, "Leveraging filter correlations for deep model compression," in *The IEEE Winter Conference on Applications of Computer Vision (WACV)*, March 2020.
- [247] J. Wang, H. Bai, J. Wu, and J. Cheng, "Bayesian automatic model compression," *IEEE Journal of Selected Topics in Signal Processing*, pp. 1–1, 2020.
- [248] Y. Gong, L. Liu, M. Yang, and L. Bourdev, "Compressing deep convolutional networks using vector quantization," *arXiv preprint arXiv:1412.6115*, 2014.
- [249] W. Chen, J. Wilson, S. Tyree, K. Weinberger, and Y. Chen, "Compressing neural networks with the hashing trick," in *International Conference on Machine Learning*, 2015, pp. 2285–2294.
- [250] 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.
- [251] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng, "Quantized convolutional neural networks for mobile devices," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 4820–4828.
- [252] M. Courbariaux, Y. Bengio, and J.-P. David, "Binaryconnect: Training deep neural networks with binary weights during propagations," in *Advances in neural information processing systems*, 2015, pp. 3123–3131.
- [253] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio, "Binarized neural networks: Training deep neural networks with weights and activations constrained to +1 or -1," *arXiv preprint arXiv:1602.02830*, 2016.
- [254] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "Xnor-net: Imagenet classification using binary convolutional neural networks," in *European Conference on Computer Vision*. Springer, 2016, pp. 525–542.
- [255] Z. Lin, M. Courbariaux, R. Memisevic, and Y. Bengio, "Neural networks with few multiplications," *arXiv preprint arXiv:1510.03009*, 2015.
- [256] V. Vanhoucke, A. Senior, and M. Z. Mao, "Improving the speed of deep neural networks on cpus," 2011.
- [257] R. Alvarez, R. Prabhavalkar, and A. Bakhtin, "On the efficient representation and execution of deep acoustic models," *arXiv preprint arXiv:1607.04683*, 2016.
- [258] M. A. Nasution, D. Chahyati, and M. I. Fanany, "Faster r-cnn with structured sparsity learning and ristretto for mobile environment," in *2017 International Conference on Advanced Computer Science and Information Systems (ICACSIS)*. IEEE, 2017, pp. 309–314.
- [259] P. Peng, Y. Mingyu, and X. Weisheng, "Running 8-bit dynamic fixed-point convolutional neural network on low-cost arm platforms," in *2017 Chinese Automation Congress (CAC)*. IEEE, 2017, pp. 4564–4568.
- [260] S. Anwar, K. Hwang, and W. Sung, "Fixed point optimization of deep convolutional neural networks for object recognition," in *2015 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2015, pp. 1131–1135.
- [261] S. H. F. Langroudi, T. Pandit, and D. Kudithipudi, "Deep learning inference on embedded devices: Fixed-point vs posit," in *2018 1st Workshop on Energy Efficient Machine Learning and Cognitive Computing for Embedded Applications (EMC2)*. IEEE, 2018, pp. 19–23.
- [262] D. Soudry, I. Hubara, and R. Meir, "Expectation backpropagation: Parameter-free training of multilayer neural networks with continuous or discrete weights," in *Advances in Neural Information Processing Systems*, 2014, pp. 963–971.
- [263] S. K. Esser, R. Appuswamy, P. Merolla, J. V. Arthur, and D. S. Modha, "Backpropagation for energy-efficient neuromorphic computing," in *Advances in Neural Information Processing Systems*, 2015, pp. 1117–1125.
- [264] A. Mathur, N. D. Lane, S. Bhattacharya, A. Boran, C. Forlivesi, and F. Kawsar, "Deepeye: Resource efficient local execution of multiple deep vision models using wearable commodity hardware," in *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*. ACM, 2017, pp. 68–81.
- [265] X. Zeng, K. Cao, and M. Zhang, "Mobiledeppill: A small-footprint mobile deep learning system for recognizing unconstrained pill images," in *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*. ACM, 2017, pp. 56–67.
- [266] P. Wang, Q. Hu, Z. Fang, C. Zhao, and J. Cheng, "Deepsearch: A fast image search framework for mobile devices," *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, vol. 14, no. 1, p. 6, 2018.
- [267] B. Kim, Y. Jeon, H. Park, D. Han, and Y. Baek, "Design and implementation of the vehicular camera system using deep neural network compression," in *Proceedings of the 1st International Workshop on Deep Learning for Mobile Systems and Applications*. ACM, 2017, pp. 25–30.
- [268] X. Xu, S. Yin, and P. Ouyang, "Fast and low-power behavior analysis on vehicles using smartphones," in *2017 6th International Symposium on Next Generation Electronics (ISNE)*. IEEE, 2017, pp. 1–4.
- [269] S. Liu, Y. Lin, Z. Zhou, K. Nan, H. Liu, and J. Du, "On-demand deep model compression for mobile devices: A usage-driven model selection framework," in *Proceedings of the 16th Annual International Conference on Mobile Systems, Applications, and Services*. ACM, 2018, pp. 389–400.
- [270] 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*. ACM, 2017, pp. 7–12.
- [271] M. Loukidakis, J. Cano, and M. OBoyle, "Accelerating deep neural networks on low power heterogeneous architectures," 2018.
- [272] S. S. L. Oskouei, H. Golestani, M. Kachuee, M. Hashemi, H. Mohammadzade, and S. Ghiasi, "Gpu-based acceleration of deep convolutional neural networks on mobile platforms," *Distrib. Parallel Clust. Comput.*, 2015.
- [273] S. S. Latifi Oskouei, H. Golestani, M. Hashemi, and S. Ghiasi, "Cndroid: Gpu-accelerated execution of trained deep convolutional neural networks on android," in *Proceedings of the 24th ACM international conference on Multimedia*. ACM, 2016, pp. 1201–1205.
- [274] P.-K. Tsung, S.-F. Tsai, A. Pai, S.-J. Lai, and C. Lu, "High performance deep neural network on low cost mobile gpu," in *2016 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 2016, pp. 69–70.
- [275] S. Rizvi, G. Cabodi, D. Patti, and G. Francini, "Gpgpu accelerated deep object classification on a heterogeneous mobile platform," *Electronics*, vol. 5, no. 4, p. 88, 2016.
- [276] S. Rizvi, D. Patti, T. Björklund, G. Cabodi, and G. Francini, "Deep classifiers-based license plate detection, localization and recognition on gpu-powered mobile platform," *Future Internet*, vol. 9, no. 4, p. 66, 2017.
- [277] S. Rizvi, G. Cabodi, and G. Francini, "Optimized deep neural networks for real-time object classification on embedded gpus," *Applied Sciences*, vol. 7, no. 8, p. 826, 2017.
- [278] S. Rizvi, G. Cabodi, D. Patti, and M. Gulzar, "A general-purpose graphics processing unit (gpgpu)-accelerated robotic controller using a low power mobile platform," *Journal of Low Power Electronics and Applications*, vol. 7, no. 2, p. 10, 2017.
- [279] Q. Cao, N. Balasubramanian, and A. Balasubramanian, "Mobirnn: Efficient recurrent neural network execution on mobile gpu," in *Proceedings of the 1st International Workshop on Deep Learning for Mobile Systems and Applications*. ACM, 2017, pp. 1–6.
- [280] H. Guihot, "Renderscript," in *Pro Android Apps Performance Optimization*. Springer, 2012, pp. 231–263.
- [281] M. Motamedi, D. Fong, and S. Ghiasi, "Fast and energy-efficient cnn inference on iot devices," *arXiv preprint arXiv:1611.07151*, 2016.
- [282] —, "Cappuccino: Efficient cnn inference software synthesis for mobile system-on-chips," *IEEE Embedded Systems Letters*, vol. 11, no. 1, pp. 9–12, 2018.
- [283] —, "Machine intelligence on resource-constrained iot devices: The case of thread granularity optimization for cnn inference," *ACM Transactions on Embedded Computing Systems (TECS)*, vol. 16, no. 5s, p. 151, 2017.
- [284] L. N. Huynh, R. K. Balan, and Y. Lee, "Deepense: A gpu-based deep convolutional neural network framework on commodity mobile

- devices,” in *Proceedings of the 2016 Workshop on Wearable Systems and Applications*. ACM, 2016, pp. 25–30.
- [285] B. Taylor, V. S. Marco, and Z. Wang, “Adaptive optimization for opencv programs on embedded heterogeneous systems,” in *ACM SIGPLAN Notices*, vol. 52, no. 5. ACM, 2017, pp. 11–20.
- [286] S. Rallapalli, H. Qiu, A. Bency, S. Karthikeyan, R. Govindan, B. Manjunath, and R. Uргаonkar, “Are very deep neural networks feasible on mobile devices,” *IEEE Trans. Circ. Syst. Video Technol.*, 2016.
- [287] 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.
- [288] M. Bettoni, G. Urgese, Y. Kobayashi, E. Macii, and A. Acquaviva, “A convolutional neural network fully implemented on fpga for embedded platforms,” in *2017 New Generation of CAS (NGCAS)*. IEEE, 2017, pp. 49–52.
- [289] 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*. ACM, 2017, pp. 45–54.
- [290] S.-S. Park, K.-B. Park, and K.-S. Chung, “Implementation of a cnn accelerator on an embedded soc platform using sdsoc,” in *Proceedings of the 2nd International Conference on Digital Signal Processing*. ACM, 2018, pp. 161–165.
- [291] Y.-H. Chen, T.-J. Yang, J. Emer, and V. Sze, “Understanding the limitations of existing energy-efficient design approaches for deep neural networks,” *Energy*, vol. 2, no. L1, p. L3, 2018.
- [292] 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.
- [293] Y.-H. Chen, T.-J. Yang, J. Emer, and V. Sze, “Eyeriss v2: A flexible accelerator for emerging deep neural networks on mobile devices,” *IEEE Journal on Emerging and Selected Topics in Circuits and Systems*, 2019.
- [294] Z. Liu, P. N. Whatmough, and M. Mattina, “Systolic tensor array: An efficient structured-sparse gemm accelerator for mobile cnn inference,” *IEEE Computer Architecture Letters*, vol. 19, no. 1, pp. 34–37, 2020.
- [295] S. Han, H. Shen, M. Philipose, S. Agarwal, A. Wolman, and A. Krishnamurthy, “Mcdnn: An approximation-based execution framework for deep stream processing under resource constraints,” in *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services*. ACM, 2016, pp. 123–136.
- [296] P. Georgiev, N. D. Lane, C. Mascolo, and D. Chu, “Accelerating mobile audio sensing algorithms through on-chip gpu offloading,” in *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*. ACM, 2017, pp. 306–318.
- [297] 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 *Proceedings of the 15th International Conference on Information Processing in Sensor Networks*. IEEE Press, 2016, p. 23.
- [298] N. D. Lane, S. Bhattacharya, P. Georgiev, C. Forlivesi, and F. Kawsar, “Accelerated deep learning inference for embedded and wearable devices using deepx,” in *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services Companion*. ACM, 2016, pp. 109–109.
- [299] 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.
- [300] T.-J. Yang, A. Howard, B. Chen, X. Zhang, A. Go, M. Sandler, V. Sze, and H. Adam, “Netadapt: Platform-aware neural network adaptation for mobile applications,” in *Proceedings of the European Conference on Computer Vision (ECCV)*, 2018, pp. 285–300.
- [301] C. Ma, Z. Zhu, J. Ye, J. Yang, J. Pei, S. Xu, R. Zhou, C. Yu, F. Mo, B. Wen *et al.*, “Deeppt: deep learning for peptide retention time prediction in proteomics,” *arXiv preprint arXiv:1705.05368*, 2017.
- [302] T. Abtahi, C. Shea, A. Kulkarni, and T. Mohsenin, “Accelerating convolutional neural network with fft on embedded hardware,” *IEEE Transactions on Very Large Scale Integration (VLSI) Systems*, vol. 26, no. 9, pp. 1737–1749, 2018.
- [303] H. Li, K. Ota, and M. Dong, “Learning iot in edge: Deep learning for the internet of things with edge computing,” *IEEE Network*, vol. 32, no. 1, pp. 96–101, 2018.
- [304] Y. Huang, X. Ma, X. Fan, J. Liu, and W. Gong, “When deep learning meets edge computing,” in *2017 IEEE 25th international conference on network protocols (ICNP)*. IEEE, 2017, pp. 1–2.
- [305] A. E. Eshratifar and M. Pedram, “Energy and performance efficient computation offloading for deep neural networks in a mobile cloud computing environment,” in *Proceedings of the 2018 on Great Lakes Symposium on VLSI*. ACM, 2018, pp. 111–116.
- [306] Y. Kang, J. Hauswald, C. Gao, A. Rovinski, T. Mudge, J. Mars, and L. Tang, “Neurosurgeon: Collaborative intelligence between the cloud and mobile edge,” in *ACM SIGARCH Computer Architecture News*, vol. 45, no. 1. ACM, 2017, pp. 615–629.
- [307] S. A. Osia, A. S. Shamsabadi, S. Sajadmanesh, A. Taheri, K. Katevas, H. R. Rabiee, N. D. Lane, and H. Haddadi, “A hybrid deep learning architecture for privacy-preserving mobile analytics,” *IEEE Internet of Things Journal*, 2020.
- [308] C. Liu, S. Chakraborty, and P. Mittal, “Deepprotect: Enabling inference-based access control on mobile sensing applications,” *arXiv preprint arXiv:1702.06159*, 2017.
- [309] C. Xu, J. Ren, L. She, Y. Zhang, Z. Qin, and K. Ren, “Edgesanitizer: Locally differentially private deep inference at the edge for mobile data analytics,” *IEEE Internet of Things Journal*, 2019.
- [310] G. Ananthanarayanan, P. Bahl, P. Bodík, K. Chintalapudi, M. Philipose, L. Ravindranath, and S. Sinha, “Real-time video analytics: The killer app for edge computing,” *computer*, vol. 50, no. 10, pp. 58–67, 2017.
- [311] M. Ali, A. Anjum, M. U. Yaseen, A. R. Zamani, D. Balouek-Thomert, O. Rana, and M. Parashar, “Edge enhanced deep learning system for large-scale video stream analytics,” in *2018 IEEE 2nd International Conference on Fog and Edge Computing (ICFEC)*. IEEE, 2018, pp. 1–10.
- [312] S. Naderiparizi, P. Zhang, M. Philipose, B. Priyantha, J. Liu, and D. Ganesan, “Glimpse: A programmable early-discard camera architecture for continuous mobile vision,” in *Proceedings of the 15th Annual International Conference on Mobile Systems, Applications, and Services*. ACM, 2017, pp. 292–305.
- [313] P. Sanabria, J. I. Benedetto, A. Neyem, J. Navon, and C. Poellabauer, “Code offloading solutions for audio processing in mobile healthcare applications: a case study,” in *2018 IEEE/ACM 5th International Conference on Mobile Software Engineering and Systems (MOBILESoft)*. IEEE, 2018, pp. 117–121.
- [314] J. Hanhiova, T. Kämäräinen, S. Seppälä, M. Siekkinen, V. Hirvisalo, and A. Ylä-Jääski, “Latency and throughput characterization of convolutional neural networks for mobile computer vision,” in *Proceedings of the 9th ACM Multimedia Systems Conference*. ACM, 2018, pp. 204–215.
- [315] B. Qi, M. Wu, and L. Zhang, “A dnn-based object detection system on mobile cloud computing,” in *2017 17th International Symposium on Communications and Information Technologies (ISCIT)*. IEEE, 2017, pp. 1–6.
- [316] X. Ran, H. Chen, Z. Liu, and J. Chen, “Delivering deep learning to mobile devices via offloading,” in *Proceedings of the Workshop on Virtual Reality and Augmented Reality Network*. ACM, 2017, pp. 42–47.
- [317] X. Ran, H. Chen, X. Zhu, Z. Liu, and J. Chen, “Deepdecision: A mobile deep learning framework for edge video analytics,” in *IEEE INFOCOM 2018-IEEE Conference on Computer Communications*. IEEE, 2018, pp. 1421–1429.
- [318] P. Georgiev, N. D. Lane, K. K. Rachuri, and C. Mascolo, “Leo: Scheduling sensor inference algorithms across heterogeneous mobile processors and network resources,” in *Proceedings of the 22nd Annual International Conference on Mobile Computing and Networking*. ACM, 2016, pp. 320–333.
- [319] M.-R. Ra, A. Sheth, L. Mummert, P. Pillai, D. Wetherall, and R. Govindan, “Odessa: enabling interactive perception applications on mobile devices,” in *Proceedings of the 9th international conference on Mobile systems, applications, and services*. ACM, 2011, pp. 43–56.
- [320] C. Streiffer, A. Srivastava, V. Orlikowski, Y. Velasco, V. Martin, N. Raval, A. Machanavajjhala, and L. P. Cox, “epivateeye: To the edge and beyond!” in *Proceedings of the Second ACM/IEEE Symposium on Edge Computing*. ACM, 2017, p. 18.
- [321] J. H. Ko, T. Na, M. F. Amir, and S. Mukhopadhyay, “Edge-host partitioning of deep neural networks with feature space encoding for resource-constrained internet-of-things platforms,” in *2018 15th IEEE International Conference on Advanced Video and Signal Based Surveillance (AVSS)*. IEEE, 2018, pp. 1–6.
- [322] Y. Tian, J. Yuan, S. Yu, and Y. Hou, “Lep-cnn: A lightweight edge device assisted privacy-preserving cnn inference solution for iot,” *arXiv preprint arXiv:1901.04100*, 2019.

- [323] C. Zhang and Z. Zheng, "Task migration for mobile edge computing using deep reinforcement learning," *Future Generation Computer Systems*, vol. 96, pp. 111–118, 2019.
- [324] H.-J. Jeong, I. Jeong, H.-J. Lee, and S.-M. Moon, "Computation offloading for machine learning web apps in the edge server environment," in *2018 IEEE 38th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2018, pp. 1492–1499.
- [325] M. Xu, F. Qian, and S. Pushp, "Enabling cooperative inference of deep learning on wearables and smartphones," *arXiv preprint arXiv:1712.03073*, 2017.
- [326] P. Liu, B. Qi, and S. Banerjee, "Edgeeye: An edge service framework for real-time intelligent video analytics," in *Proceedings of the 1st International Workshop on Edge Systems, Analytics and Networking*. ACM, 2018, pp. 1–6.
- [327] 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.
- [328] S. Yi, Z. Hao, Q. Zhang, Q. Zhang, W. Shi, and Q. Li, "Lavea: Latency-aware video analytics on edge computing platform," in *Proceedings of the Second ACM/IEEE Symposium on Edge Computing*. ACM, 2017, p. 15.
- [329] R. Hadidi, J. Cao, M. Woodward, M. S. Ryoo, and H. Kim, "Musical chair: Efficient real-time recognition using collaborative iot devices," *arXiv preprint arXiv:1802.02138*, 2018.
- [330] N. Talagala, S. Sundararaman, V. Sridhar, D. Arteaga, Q. Luo, S. Subramanian, S. Ghanta, L. Khermosh, and D. Roselli, "{ECO}: Harmonizing edge and cloud with ml/dl orchestration," in *{USENIX} Workshop on Hot Topics in Edge Computing (HotEdge 18)*, 2018.
- [331] E. De Coninck, S. Bohez, S. Leroux, T. Verbelen, B. Vankeirsbilck, P. Simoens, and B. Dhoedt, "Dianne: a modular framework for designing, training and deploying deep neural networks on heterogeneous distributed infrastructure," *Journal of Systems and Software*, vol. 141, pp. 52–65, 2018.
- [332] Y. Fukushima, D. Miura, T. Hamatani, H. Yamaguchi, and T. Higashino, "Microdeep: In-network deep learning by micro-sensor coordination for pervasive computing," in *2018 IEEE International Conference on Smart Computing (SMARTCOMP)*. IEEE, 2018, pp. 163–170.
- [333] T. Bach, M. A. Tariq, R. Mayer, and K. Rothermel, "Knowledge is at the edge! how to search in distributed machine learning models," in *OTM Confederated International Conferences "On the Move to Meaningful Internet Systems"*. Springer, 2017, pp. 410–428.
- [334] A. Morshed, P. P. Jayaraman, T. Sellis, D. Georgakopoulos, M. Villari, and R. Ranjan, "Deep osmosis: Holistic distributed deep learning in osmotic computing," *IEEE Cloud Computing*, vol. 4, no. 6, pp. 22–32, 2017.
- [335] S. Teerapittayanon, B. McDanel, and H.-T. Kung, "Distributed deep neural networks over the cloud, the edge and end devices," in *2017 IEEE 37th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2017, pp. 328–339.
- [336] A. Yousefpour, S. Devic, B. Q. Nguyen, A. Kreidieh, A. Liao, A. M. Bayen, and J. P. Jue, "Guardians of the deep fog: Failure-resilient dnn inference from edge to cloud," in *Proceedings of the First International Workshop on Challenges in Artificial Intelligence and Machine Learning for Internet of Things*, 2019, pp. 25–31.
- [337] A. Ferdowsi, U. Challita, and W. Saad, "Deep learning for reliable mobile edge analytics in intelligent transportation systems: An overview," *IEEE vehicular technology magazine*, vol. 14, no. 1, pp. 62–70, 2019.
- [338] L. Li, K. Ota, and M. Dong, "Deep learning for smart industry: Efficient manufacture inspection system with fog computing," *IEEE Transactions on Industrial Informatics*, vol. 14, no. 10, pp. 4665–4673, 2018.
- [339] B. Tang, Z. Chen, G. Hefferman, S. Pei, T. Wei, H. He, and Q. Yang, "Incorporating intelligence in fog computing for big data analysis in smart cities," *IEEE Transactions on Industrial Informatics*, vol. 13, no. 5, pp. 2140–2150, 2017.
- [340] C. Liu, Y. Cao, Y. Luo, G. Chen, V. Vokkarane, M. Yunsheng, S. Chen, and P. Hou, "A new deep learning-based food recognition system for dietary assessment on an edge computing service infrastructure," *IEEE Transactions on Services Computing*, vol. 11, no. 2, pp. 249–261, 2017.
- [341] T. Muhammed, R. Mehmood, A. Albeshri, and I. Katib, "Ubehealth: a personalized ubiquitous cloud and edge-enabled networked healthcare system for smart cities," *IEEE Access*, vol. 6, pp. 32 258–32 285, 2018.
- [342] M. Schwabacher and K. Goebel, "A survey of artificial intelligence for prognostics," in *AAAI Fall Symposium: Artificial Intelligence for Prognostics*, 2007, pp. 108–115.
- [343] A. He, K. K. Bae, T. R. Newman, J. Gaeddert, K. Kim, R. Menon, L. Morales-Tirado, Y. Zhao, J. H. Reed, W. H. Tranter *et al.*, "A survey of artificial intelligence for cognitive radios," *IEEE Transactions on Vehicular Technology*, vol. 59, no. 4, pp. 1578–1592, 2010.
- [344] A. Bahrammirzaee, "A comparative survey of artificial intelligence applications in finance: artificial neural networks, expert system and hybrid intelligent systems," *Neural Computing and Applications*, vol. 19, no. 8, pp. 1165–1195, 2010.
- [345] Y. Zhang, J. Ren, J. Liu, C. Xu, H. Guo, and Y. Liu, "A survey on emerging computing paradigms for big data," *Chinese Journal of Electronics*, vol. 26, no. 1, pp. 1–12, 2017.
- [346] D. Singh and C. K. Reddy, "A survey on platforms for big data analytics," *Journal of big data*, vol. 2, no. 1, p. 8, 2015.
- [347] S. Yi, C. Li, and Q. Li, "A survey of fog computing: concepts, applications and issues," in *Proceedings of the 2015 workshop on mobile big data*, 2015, pp. 37–42.
- [348] T.-h. Kim, C. Ramos, and S. Mohammed, "Smart city and iot," 2017.
- [349] K. Su, J. Li, and H. Fu, "Smart city and the applications," in *2011 international conference on electronics, communications and control (ICECC)*. IEEE, 2011, pp. 1028–1031.
- [350] W. Zhang, B. Han, and P. Hui, "Jaguar: Low latency mobile augmented reality with flexible tracking," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 355–363.
- [351] W. Zhang, S. Lin, F. H. Bijarbooneh, H. F. Cheng, and P. Hui, "Cloudar: A cloud-based framework for mobile augmented reality," in *Proceedings of the on Thematic Workshops of ACM Multimedia 2017*, 2017, pp. 194–200.
- [352] S. Lin, H. F. Cheng, W. Li, Z. Huang, P. Hui, and C. Peylo, "Ubii: Physical world interaction through augmented reality," *IEEE Transactions on Mobile Computing*, vol. 16, no. 3, pp. 872–885, 2016.
- [353] J. L. Hennessy and D. A. Patterson, *Computer architecture: a quantitative approach*. Elsevier, 2011.
- [354] G. S. Paschos, G. Iosifidis, M. Tao, D. Towsley, and G. Caire, "The role of caching in future communication systems and networks," *IEEE Journal on Selected Areas in Communications*, vol. 36, no. 6, pp. 1111–1125, 2018.
- [355] L. Li, G. Zhao, and R. S. Blum, "A survey of caching techniques in cellular networks: Research issues and challenges in content placement and delivery strategies," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 3, pp. 1710–1732, 2018.
- [356] S. Wang, X. Zhang, Y. Zhang, L. Wang, J. Yang, and W. Wang, "A survey on mobile edge networks: Convergence of computing, caching and communications," *IEEE Access*, vol. 5, pp. 6757–6779, 2017.
- [357] Google Glass, https://en.wikipedia.org/wiki/Google_Glass, 2019.
- [358] Microsoft Hololens, https://en.wikipedia.org/wiki/Microsoft_HoloLens, 2019.
- [359] T. Braud, P. Zhou, J. Kangasharju, and P. Hui, "Multipath computation offloading for mobile augmented reality," in *In Proceedings of the IEEE International Conference on Pervasive Computing and Communications (PerCom 2020)*, Austin USA, 2020.
- [360] L. Lovagnini, W. Zhang, F. H. Bijarbooneh, and P. Hui, "Circe: Real-time caching for instance recognition on cloud environments and multi-core architectures," in *Proceedings of the 26th ACM international conference on Multimedia*, 2018, pp. 346–354.
- [361] Z. Xiao, T. Li, W. Cheng, and D. Wang, "Apollonius circles based outbound handover in macro-small wireless cellular networks," in *2016 IEEE Global Communications Conference (GLOBECOM)*. IEEE, 2016, pp. 1–6.
- [362] E. Bacstug, M. Bennis, M. Kountouris, and M. Debbah, "Cache-enabled small cell networks: Modeling and tradeoffs," *EURASIP Journal on Wireless Communications and Networking*, vol. 2015, no. 1, p. 41, 2015.
- [363] Z. Chen and M. Kountouris, "Cache-enabled small cell networks with local user interest correlation," in *2015 IEEE 16th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2015, pp. 680–684.
- [364] J. Liao, K.-K. Wong, M. R. Khandaker, and Z. Zheng, "Optimizing cache placement for heterogeneous small cell networks," *IEEE Communications Letters*, vol. 21, no. 1, pp. 120–123, 2017.
- [365] K. Poularakis, G. Iosifidis, V. Sourlas, and L. Tassiulas, "Multicast-aware caching for small cell networks," in *2014 IEEE Wireless Communications and Networking Conference (WCNC)*. IEEE, 2014, pp. 2300–2305.
- [366] H. Dahrouj and W. Yu, "Coordinated beamforming for the multicell multi-antenna wireless system," *IEEE transactions on wireless communications*, vol. 9, no. 5, pp. 1748–1759, 2010.

- [367] P. Marsch and G. P. Fettweis, *Coordinated Multi-Point in Mobile Communications: from theory to practice*. Cambridge University Press, 2011.
- [368] M. Ji, G. Caire, and A. F. Molisch, "Wireless device-to-device caching networks: Basic principles and system performance," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 1, pp. 176–189, 2016.
- [369] W. Chen, T. Li, Z. Xiao, and D. Wang, "On mitigating interference under device-to-device communication in macro-small cell networks," in *2016 International Conference on Computer, Information and Telecommunication Systems (CITS)*. IEEE, 2016, pp. 1–5.
- [370] Z. Chen and M. Kountouris, "D2d caching vs. small cell caching: Where to cache content in a wireless network?" in *2016 IEEE 17th International Workshop on Signal Processing Advances in Wireless Communications (SPAWC)*. IEEE, 2016, pp. 1–6.
- [371] M. Gregori, J. Gómez-Vilardebó, J. Matamoros, and D. Gündüz, "Wireless content caching for small cell and d2d networks," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 5, pp. 1222–1234, 2016.
- [372] D. Liu and C. Yang, "Will caching at base station improve energy efficiency of downlink transmission?" in *2014 IEEE Global Conference on Signal and Information Processing (GlobalSIP)*. IEEE, 2014, pp. 173–177.
- [373] —, "Energy efficiency of downlink networks with caching at base stations," *IEEE Journal on Selected Areas in Communications*, vol. 34, no. 4, pp. 907–922, 2016.
- [374] J. Gu, W. Wang, A. Huang, and H. Shan, "Proactive storage at caching-enable base stations in cellular networks," in *2013 IEEE 24th Annual International Symposium on Personal, Indoor, and Mobile Radio Communications (PIMRC)*. IEEE, 2013, pp. 1543–1547.
- [375] A. Khreishah and J. Chakareski, "Collaborative caching for multicell-coordinated systems," in *2015 IEEE Conference on Computer Communications Workshops (INFOCOM WKSHPs)*. IEEE, 2015, pp. 257–262.
- [376] P. Ostovari, J. Wu, and A. Khreishah, "Efficient online collaborative caching in cellular networks with multiple base stations," in *2016 IEEE 13th International Conference on Mobile Ad Hoc and Sensor Systems (MASS)*. IEEE, 2016, pp. 136–144.
- [377] R. Wang, X. Peng, J. Zhang, and K. B. Letaief, "Mobility-aware caching for content-centric wireless networks: Modeling and methodology," *IEEE Communications Magazine*, vol. 54, no. 8, pp. 77–83, 2016.
- [378] J. Li, C. Shunfeng, F. Shu, J. Wu, and D. N. K. Jayakody, "Contract-based small-cell caching for data disseminations in ultra-dense cellular networks," *IEEE Transactions on Mobile Computing*, 2018.
- [379] K. Poularakis, V. Sourlas, P. Flegkas, and L. Tassiulas, "On exploiting network coding in cache-capable small-cell networks," in *2014 IEEE Symposium on Computers and Communications (ISCC)*. IEEE, 2014, pp. 1–5.
- [380] S. Krishnan and H. S. Dhillon, "Effect of user mobility on the performance of device-to-device networks with distributed caching," *IEEE Wireless Communications Letters*, vol. 6, no. 2, pp. 194–197, 2017.
- [381] A. Ioannou and S. Weber, "A survey of caching policies and forwarding mechanisms in information-centric networking," *IEEE Communications Surveys & Tutorials*, vol. 18, no. 4, pp. 2847–2886, 2016.
- [382] S. Amershi, M. Cakmak, W. B. Knox, and T. Kulesza, "Power to the people: The role of humans in interactive machine learning," *Ai Magazine*, vol. 35, no. 4, pp. 105–120, 2014.
- [383] M. Ware, E. Frank, G. Holmes, M. Hall, and I. H. Witten, "Interactive machine learning: letting users build classifiers," *International Journal of Human-Computer Studies*, vol. 55, no. 3, pp. 281–292, 2001.
- [384] J. B. Predd, S. B. Kulkarni, and H. V. Poor, "Distributed learning in wireless sensor networks," *IEEE Signal Processing Magazine*, vol. 23, no. 4, pp. 56–69, 2006.
- [385] M. Lavassani, S. Forsström, U. Jennehag, and T. Zhang, "Combining fog computing with sensor mote machine learning for industrial iot," *Sensors*, vol. 18, no. 5, p. 1532, 2018.
- [386] C. Ma, J. Li, M. Ding, H. H. Yang, F. Shu, T. Q. S. Quek, and H. V. Poor, "On safeguarding privacy and security in the framework of federated learning," *IEEE Network*, pp. 1–7, 2020.
- [387] S. Yao, Y. Zhao, H. Shao, C. Zhang, A. Zhang, D. Liu, S. Liu, L. Su, and T. Abdelzaker, "Apdeepsense: Deep learning uncertainty estimation without the pain for iot applications," in *2018 IEEE 38th International Conference on Distributed Computing Systems (ICDCS)*. IEEE, 2018, pp. 334–343.
- [388] J. Hu, L. Shen, and G. Sun, "Squeeze-and-excitation networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2018, pp. 7132–7141.
- [389] D. Kong, "Science driven innovations powering mobile product: Cloud ai vs. device ai solutions on smart device," *arXiv preprint arXiv:1711.07580*, 2017.
- [390] T. Guo, "Cloud-based or on-device: An empirical study of mobile deep inference," in *2018 IEEE International Conference on Cloud Engineering (IC2E)*. IEEE, 2018, pp. 184–190.
- [391] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *arXiv preprint arXiv:1808.05377*, 2018.
- [392] M. Wistuba, A. Rawat, and T. Pedapati, "A survey on neural architecture search," *arXiv preprint arXiv:1905.01392*, 2019.
- [393] F. Chollet, "Xception: Deep learning with depthwise separable convolutions," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1251–1258.
- [394] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He, "Aggregated residual transformations for deep neural networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 1492–1500.
- [395] B. Fan, X. Liu, X. Su, J. Niu, and P. Hui, "Emgauth: An emg-based smartphone unlocking system using siamese network," in *In Proceedings of the IEEE International Conference on Pervasive Computing and Communications (PerCom 2020)*, Austin USA. IEEE, 2020.
- [396] D. Wen, H. Han, and A. K. Jain, "Face spoof detection with image distortion analysis," *IEEE Transactions on Information Forensics and Security*, vol. 10, no. 4, pp. 746–761, 2015.
- [397] T. Plötz and Y. Guan, "Deep learning for human activity recognition in mobile computing," *Computer*, vol. 51, no. 5, pp. 50–59, 2018.
- [398] Y. Guan and T. Plötz, "Ensembles of deep lstm learners for activity recognition using wearables," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 1, no. 2, p. 11, 2017.
- [399] M. Shoaib, O. D. Incel, H. Scolten, and P. Havinga, "Resource consumption analysis of online activity recognition on mobile phones and smartwatches," in *2017 IEEE 36th International Performance Computing and Communications Conference (IPCCC)*. IEEE, 2017, pp. 1–6.
- [400] A. Stisen, H. Blunck, S. Bhattacharya, T. S. Prentow, M. B. Kjørgaard, A. Dey, T. Sonne, and M. M. Jensen, "Smart devices are different: Assessing and mitigating mobile sensing heterogeneities for activity recognition," in *Proceedings of the 13th ACM Conference on Embedded Networked Sensor Systems*. ACM, 2015, pp. 127–140.
- [401] A. Rosenfeld and J. K. Tsotsos, "Incremental learning through deep adaptation," *IEEE transactions on pattern analysis and machine intelligence*, 2018.
- [402] W. T. Ang, P. K. Khosla, and C. N. Riviere, "Nonlinear regression model of a low-g mems accelerometer," *IEEE Sensors Journal*, vol. 7, no. 1, pp. 81–88, 2007.
- [403] S. G. Klauer, F. Guo, B. G. Simons-Morton, M. C. Ouimet, S. E. Lee, and T. A. Dingus, "Distracted driving and risk of road crashes among novice and experienced drivers," *New England journal of medicine*, vol. 370, no. 1, pp. 54–59, 2014.
- [404] A. Graves, A.-r. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in *2013 IEEE international conference on acoustics, speech and signal processing*. IEEE, 2013, pp. 6645–6649.
- [405] M. Rabbi, S. Ali, T. Choudhury, and E. Berke, "Passive and in-situ assessment of mental and physical well-being using mobile sensors," in *Proceedings of the 13th international conference on Ubiquitous computing*. ACM, 2011, pp. 385–394.
- [406] A. Gebhart, "Google home to the amazon echo: 'anything you can do...'," *cnet, May*, vol. 18, p. 7, 2017.
- [407] X. Xu, Y. Ding, S. X. Hu, M. Niemier, J. Cong, Y. Hu, and Y. Shi, "Scaling for edge inference of deep neural networks," *Nature Electronics*, vol. 1, no. 4, p. 216, 2018.
- [408] M. Denil, B. Shakibi, L. Dinh, N. De Freitas *et al.*, "Predicting parameters in deep learning," in *Advances in neural information processing systems*, 2013, pp. 2148–2156.
- [409] C. Bucilu, R. Caruana, and A. Niculescu-Mizil, "Model compression," in *Proceedings of the 12th ACM SIGKDD international conference on Knowledge discovery and data mining*. ACM, 2006, pp. 535–541.
- [410] R. G. Baraniuk, "Compressive sensing," *IEEE signal processing magazine*, vol. 24, no. 4, 2007.
- [411] Y. LeCun, J. S. Denker, and S. A. Solla, "Optimal brain damage," in *Advances in neural information processing systems*, 1990, pp. 598–605.

- [412] B. Hassibi and D. G. Stork, "Second order derivatives for network pruning: Optimal brain surgeon," in *Advances in neural information processing systems*, 1993, pp. 164–171.
- [413] J. Van Leeuwen, "On the construction of huffman trees." in *ICALP*, 1976, pp. 382–410.
- [414] S. Malki and L. Spaanenburg, "Cnn image processing on a xilinx virtex-ii 6000," in *Proceedings ECCTD*, vol. 3, 2003, pp. 261–264.
- [415] J. L. Gustafson and I. T. Yonemoto, "Beating floating point at its own game: Posit arithmetic," *Supercomputing Frontiers and Innovations*, vol. 4, no. 2, pp. 71–86, 2017.
- [416] R. Morris, "Tapered floating point: A new floating-point representation," *IEEE Transactions on Computers*, vol. 100, no. 12, pp. 1578–1579, 1971.
- [417] 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.
- [418] 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*, 2019.
- [419] J. Appleyard, T. Kocisky, and P. Blunsom, "Optimizing performance of recurrent neural networks on gpus," *arXiv preprint arXiv:1604.01946*, 2016.
- [420] S. Liu, Q. Wang, and G. Liu, "A versatile method of discrete convolution and fft (dc-fft) for contact analyses," *Wear*, vol. 243, no. 1-2, pp. 101–111, 2000.
- [421] V. S. Marco, B. Taylor, Z. Wang, and Y. Elkhatib, "Optimizing deep learning inference on embedded systems through adaptive model selection," *ACM Trans. Embed. Comput. Syst.*, vol. 19, no. 1, Feb. 2020. [Online]. Available: <https://doi.org/10.1145/3371154>
- [422] A. Garofalo, M. Rusci, F. Conti, D. Rossi, and L. Benini, "Pulp-nn: accelerating quantized neural networks on parallel ultra-low-power risc-v processors," *Philosophical Transactions of the Royal Society A*, vol. 378, no. 2164, p. 20190155, 2020.
- [423] N. Abbas, Y. Zhang, A. Taherkordi, and T. Skeie, "Mobile edge computing: A survey," *IEEE Internet of Things Journal*, vol. 5, no. 1, pp. 450–465, 2017.
- [424] K. A. Shatilov, D. Chatzopoulos, A. W. T. Hang, and P. Hui, "Using deep learning and mobile offloading to control a 3d-printed prosthetic hand," *Proceedings of the ACM on Interactive, Mobile, Wearable and Ubiquitous Technologies*, vol. 3, no. 3, pp. 1–19, 2019.
- [425] S. Kosta, A. Aucinas, P. Hui, R. Mortier, and X. Zhang, "Thinkair: Dynamic resource allocation and parallel execution in the cloud for mobile code offloading," in *2012 proceedings IEEE Infocom*. IEEE, 2012, pp. 945–953.
- [426] N. Raval, A. Srivastava, A. Razeen, K. Lebeck, A. Machanavajjhala, and L. P. Cox, "What you mark is what apps see," in *Proceedings of the 14th Annual International Conference on Mobile Systems, Applications, and Services*. ACM, 2016, pp. 249–261.
- [427] J. Hoisko, "Context triggered visual episodic memory prosthesis," in *Digest of Papers. Fourth International Symposium on Wearable Computers*. IEEE, 2000, pp. 185–186.
- [428] S. Hodges, L. Williams, E. Berry, S. Izadi, J. Srinivasan, A. Butler, G. Smyth, N. Kapur, and K. Wood, "Sensecam: A retrospective memory aid," in *International Conference on Ubiquitous Computing*. Springer, 2006, pp. 177–193.
- [429] W. Cui, Y. Kim, and T. S. Rosing, "Cross-platform machine learning characterization for task allocation in iot ecosystems," in *2017 IEEE 7th Annual Computing and Communication Workshop and Conference (CCWC)*. IEEE, 2017, pp. 1–7.
- [430] D. Xu, Y. Li, X. Chen, J. Li, P. Hui, S. Chen, and J. Crowcroft, "A survey of opportunistic offloading," *IEEE Communications Surveys & Tutorials*, vol. 20, no. 3, pp. 2198–2236, 2018.
- [431] R. Gilad-Bachrach, N. Dowlin, K. Laine, K. Lauter, M. Naehrig, and J. Wernsing, "Cryptonets: Applying neural networks to encrypted data with high throughput and accuracy," in *International Conference on Machine Learning*, 2016, pp. 201–210.
- [432] H. Chabanne, A. de Wargny, J. Milgram, C. Morel, and E. Prouff, "Privacy-preserving classification on deep neural network." *IACR Cryptology ePrint Archive*, vol. 2017, p. 35, 2017.
- [433] E. Hesamifard, H. Takabi, and M. Ghasemi, "Cryptodl: Deep neural networks over encrypted data," *arXiv preprint arXiv:1711.05189*, 2017.
- [434] S. M. Johnson, "Optimal two-and three-stage production schedules with setup times included," *Naval research logistics quarterly*, vol. 1, no. 1, pp. 61–68, 1954.
- [435] W. Lee, S. Kim, Y.-T. Lee, H.-W. Lee, and M. Choi, "Deep neural networks for wild fire detection with unmanned aerial vehicle," in *2017 IEEE international conference on consumer electronics (ICCE)*. IEEE, 2017, pp. 252–253.
- [436] A. Thomas, Y. Guo, Y. Kim, B. Aksanli, A. Kumar, and T. S. Rosing, "Pushing down machine learning inference to the edge in heterogeneous internet of things applications," 2018.
- [437] M. Villari, M. Fazio, S. Dustdar, O. Rana, and R. Ranjan, "Osmotic computing: A new paradigm for edge/cloud integration," *IEEE Cloud Computing*, vol. 3, no. 6, pp. 76–83, 2016.
- [438] A. Mathur, T. Zhang, S. Bhattacharya, P. Velickovic, L. Joffe, N. D. Lane, F. Kawsar, and P. Lió, "Using deep data augmentation training to address software and hardware heterogeneities in wearable and smartphone sensing devices," in *2018 17th ACM/IEEE International Conference on Information Processing in Sensor Networks (IPSN)*. IEEE, 2018, pp. 200–211.
- [439] A. Das, N. Borisov, and M. Caesar, "Fingerprinting smart devices through embedded acoustic components," *arXiv preprint arXiv:1403.3366*, 2014.
- [440] A. Mathur, A. Isopoussu, F. Kawsar, R. Smith, N. D. Lane, and N. Berthouze, "On robustness of cloud speech apis: An early characterization," in *Proceedings of the 2018 ACM International Joint Conference and 2018 International Symposium on Pervasive and Ubiquitous Computing and Wearable Computers*. ACM, 2018, pp. 1409–1413.
- [441] Y. Bengio, A. Courville, and P. Vincent, "Representation learning: A review and new perspectives," *IEEE transactions on pattern analysis and machine intelligence*, vol. 35, no. 8, pp. 1798–1828, 2013.
- [442] Z. Chen and B. Liu, "Lifelong machine learning," *Synthesis Lectures on Artificial Intelligence and Machine Learning*, vol. 10, no. 3, pp. 1–145, 2016.
- [443] R. Vilalta and Y. Drissi, "A perspective view and survey of meta-learning," *Artificial intelligence review*, vol. 18, no. 2, pp. 77–95, 2002.
- [444] A. Soller, J. Wiebe, and A. Lesgold, "A machine learning approach to assessing knowledge sharing during collaborative learning activities," in *Proceedings of the Conference on Computer Support for Collaborative Learning: Foundations for a CSCL Community*. International Society of the Learning Sciences, 2002, pp. 128–137.



Dianlei Xu received the B.S. degree from Anhui University, Hefei, China, and is currently a joint doctoral student in the Department of Computer Science, Helsinki, Finland and Beijing National Research Center for Information Science and Technology (BNRist), Department of Electronic Engineering, Tsinghua University, Beijing, China.

His research interests include edge/fog computing and AIoT.

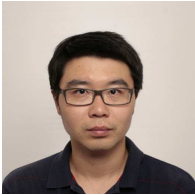


Tong Li received the B.S. degree and M.S. degree in communication engineering from Hunan University, China, in 2014 and 2017. At present, he is a dual Ph.D. student at the Hong Kong University of Science and Technology and the University of Helsinki. His research interests include distributed systems, edge network, and data-driven network. He is an IEEE student member.



Yong Li (M'09-SM'16) received the B.S. degree in electronics and information engineering from Huazhong University of Science and Technology, Wuhan, China, in 2007 and the Ph.D. degree in electronic engineering from Tsinghua University, Beijing, China, in 2012. He is currently a Faculty Member of the Department of Electronic Engineering, Tsinghua University.

Dr. Li has served as General Chair, TPC Chair, SPC/TPC Member for several international workshops and conferences, and he is on the editorial board of two IEEE journals. His papers have total citations more than 6900. Among them, ten are ESI Highly Cited Papers in Computer Science, and four receive conference Best Paper (run-up) Awards. He received IEEE 2016 ComSoc Asia-Pacific Outstanding Young Researchers, Young Talent Program of China Association for Science and Technology, and the National Youth Talent Support Program.



Xiang Su received his Ph.D. in technology from the University of Oulu in 2016. He is currently an Academy of Finland postdoc fellow and a senior postdoctoral researcher in computer science in the University of Helsinki. Dr. Su has extensive expertise on Internet of Things, edge computing, mobile augmented reality, knowledge representations, and context modeling and reasoning.



Sasu Tarkoma received the MSc and PhD degrees in computer science from the Department of Computer Science, University of Helsinki. He is a full professor at the Department of Computer Science, University of Helsinki, and the deputy head of the department. He has managed and participated in national and international research projects at the University of Helsinki, Aalto University, and the Helsinki Institute for Information Technology. His research interests include mobile computing, Internet technologies, and middleware. He is a senior member

of the IEEE.



Tao Jiang (M'06-SM'10-F'19) received the Ph.D. degree in information and communication engineering from the Huazhong University of Science and Technology, Wuhan, China, in April 2004. From August 2004 to December 2007, he worked in some universities, such as Brunel University and University of Michigan-Dearborn, respectively. He is currently a Distinguished Professor with the Wuhan National Laboratory for Optoelectronics and School of Electronics Information and Communications, Huazhong University of Science and Technology.

He has authored or coauthored more than 300 technical articles in major journals and conferences and nine books/chapters in the areas of communications and networks. He served or is serving as symposium technical program committee membership of some major IEEE conferences, including INFOCOM, GLOBECOM, and ICC. He was invited to serve as a TPC Symposium Chair for the IEEE GLOBECOM 2013, the IEEE WCNC 2013, and ICC 2013. He is served or serving as an Associate Editor of some technical journals in communications, including in the IEEE NETWORK, the IEEE TRANSACTIONS ON SIGNAL PROCESSING, the IEEE COMMUNICATIONS SURVEYS AND TUTORIALS, the IEEE TRANSACTIONS ON VEHICULAR TECHNOLOGY, and the IEEE INTERNET OF THINGS JOURNAL. He is the Associate Editor-in-Chief of China Communications.



Jon Crowcroft (SM'95-F'04) graduated in physics from Trinity College, Cambridge University, United Kingdom, in 1979, and received the MSc degree in computing in 1981 and the PhD degree in 1993 from University College London (UCL), United Kingdom. He is currently the Marconi Professor of Communications Systems in the Computer Lab at the University of Cambridge, United Kingdom. Professor Crowcroft is a fellow of the United Kingdom Royal Academy of Engineering, a fellow of the ACM, and a fellow of IET. He was a recipient of

the ACM Sigcomm Award in 2009.



Pan Hui (SM'14-F'18) received his PhD from the Computer Laboratory at University of Cambridge, and both his Bachelor and MPhil degrees from the University of Hong Kong.

He is the Nokia Chair Professor in Data Science and Professor of Computer Science at the University of Helsinki. He is also the director of the HKUST-DT System and Media Lab at the Hong Kong University of Science and Technology. He was an adjunct Professor of social computing and networking at Aalto University from 2012 to 2017.

He was a senior research scientist and then a Distinguished Scientist for Telekom Innovation Laboratories (T-labs) Germany from 2008 to 2015. His industrial profile also includes his research at Intel Research Cambridge and Thomson Research Paris from 2004 to 2006. His research has been generously sponsored by Nokia, Deutsche Telekom, Microsoft Research, and China Mobile. He has published more than 300 research papers and with over 17,500 citations. He has 30 granted and filed European and US patents in the areas of augmented reality, data science, and mobile computing.

He has founded and chaired several IEEE/ACM conferences/workshops, and has served as track chair, senior program committee member, organising committee member, and program committee member of numerous top conferences including ACM WWW, ACM SIGCOMM, ACM Mobisys, ACM MobiCom, ACM CoNext, IEEE Infocom, IEEE ICNP, IEEE ICDCS, IJCAI, AAAI, and ICWSM. He is an associate editor for the leading journals IEEE Transactions on Mobile Computing and IEEE Transactions on Cloud Computing. He is an IEEE Fellow, an ACM Distinguished Scientist, and a member of Academia Europaea.