# Key Advances in Pervasive Edge Computing for Industrial Internet of Things in 5G and Beyond

**ARUN NARAYANAN**[1], (Member, IEEE),
**ARTHUR SOUSA DE SENA**[1], (Student Member, IEEE),
**DANIEL GUTIERREZ-ROJAS**[1], (Graduate Student Member, IEEE),
**DICK CARRILLO MELGAREJO**[1], (Graduate Student Member, IEEE),
**HAFIZ MAJID HUSSAIN**[1], (Member, IEEE),
**MEHAR ULLAH**[1], (Graduate Student Member, IEEE),
**SUZAN BAYHAN**[2], **AND PEDRO H. J. NARDELLI**[1], (Senior Member, IEEE)

[1]Department of Electrical Engineering, School of Energy Systems, LUT University, 53850 Lappeenranta, Finland
[2]Design and Analysis of Communication Systems (DACS), Faculty of EEMCS, University of Twente, 7500 AE Enschede, The Netherlands

Corresponding author: Arun Narayanan (arun.narayanan@lut.fi)

**ABSTRACT** This article surveys emerging technologies related to pervasive edge computing (PEC) for industrial internet-of-things (IIoT) enabled by fifth-generation (5G) and beyond communication networks. PEC encompasses all devices that are capable of performing computational tasks locally, including those at the edge of the core network (edge servers co-located with 5G base stations) and in the radio access network (sensors, actuators, etc.). The main advantages of this paradigm are core network offloading (and benefits therefrom) and low latency for delay-sensitive applications (e.g., automatic control). We have reviewed the state-of-the-art in the PEC paradigm and its applications to the IIoT domain, which have been enabled by the recent developments in 5G technology. We have classified and described three important research areas related to PEC—distributed artificial intelligence methods, energy efficiency, and cyber security. We have also identified the main open challenges that must be solved to have a scalable PEC-based IIoT network that operates efficiently under different conditions. By explaining the applications, challenges, and opportunities, our paper reinforces the perspective that the PEC paradigm is an extremely suitable and important deployment model for industrial communication networks, considering the modern trend toward private industrial 5G networks with local operations and flexible management.

**INDEX TERMS** Edge computing, industrial Internet of Things, 5G network, energy efficiency, artificial intelligence, cyber security.

## I. INTRODUCTION

Today, the Internet of Things (IoT) is a well-established paradigm in modern wireless telecommunications with numerous applications to society and the industry. IoT systems have radically evolved from simple solutions involving single devices such as a single internet-connected video camera to more advanced systems involving real-time analytics, artificial intelligence, and hardware such as smart sensors and actuators. IoT is also key for the so-called Industry 4.0, where "smart" objects (machines and products) are leading to a new paradigm shift in industrial production [1].

The application of IoT to the industrial sector is now generally called *Industrial Internet of Things* (IIoT) [2]. In IIoT, a massive number of (smart) industrial machines, actuators, and sensors connect to each other to form a network of smart IoT-based devices with some computing power, communications capabilities, and data storage and caching [3], which can be potentially shared to jointly perform computational tasks. IIoT solutions can be used to improve connectivity, efficiency, profits, scalability, and data speeds for industrial applications, thereby enhancing predictive maintenance, increasing safety, and boosting operational efficiencies.

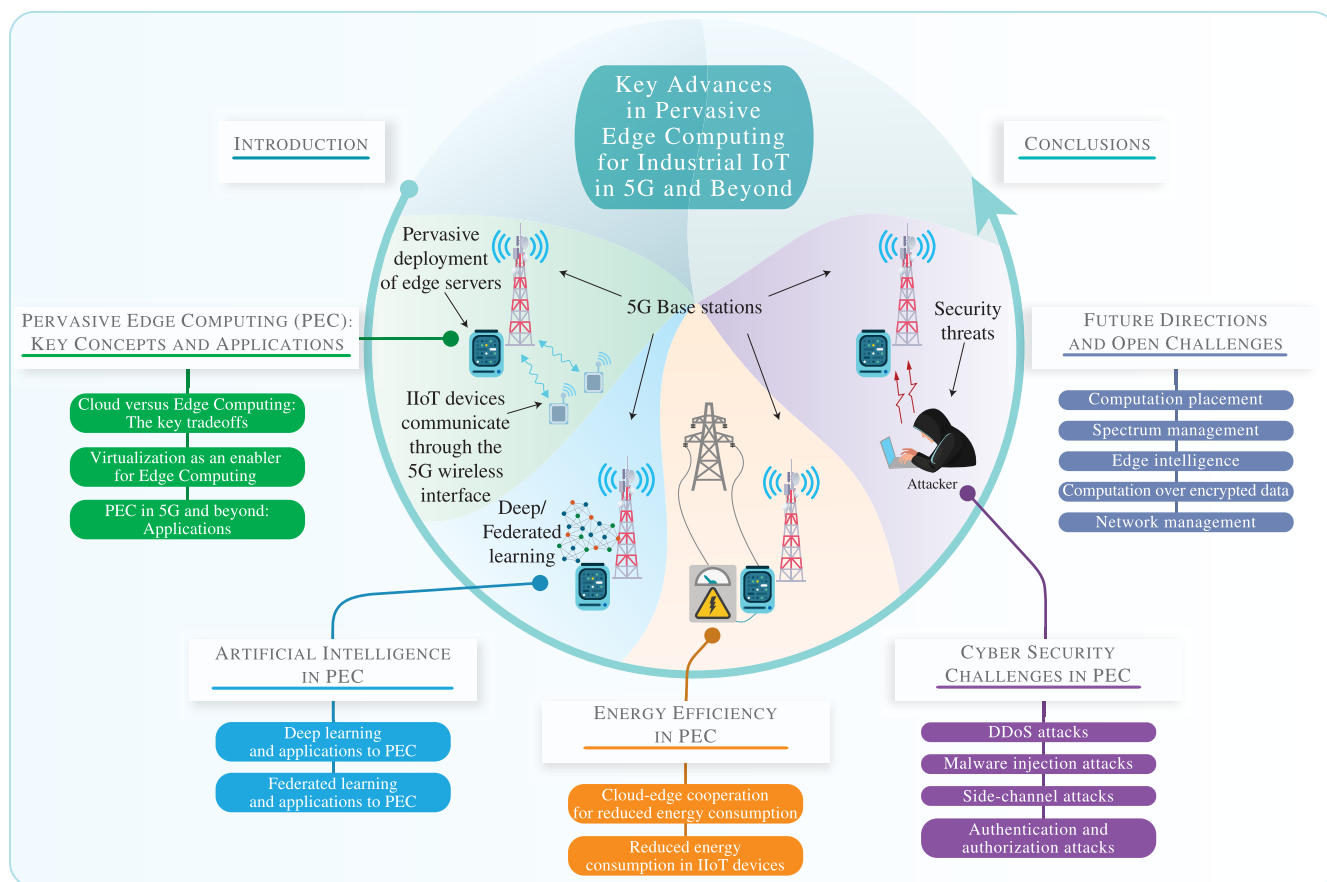The associate editor coordinating the review of this manuscript and approving it for publication was Xiaodong Xu.

**FIGURE 1.** Schematic presentation of this paper's structure and contributions.

Nevertheless, the full potential of IIoT can be realized only if it is enabled by a flexible communication system capable of supporting different requirements, from ultra-reliable low-latency communications (URLLC) to massive connectivity [4]. In this regard, recent advances in fifth-generation (5G) communication technologies have emerged at the center of IIoT applications [5] by offering greater bandwidth, faster data transmission, and improved spectral efficiency supported by localized private networks and micro-operators [6]. Nevertheless, new developments in both radio access technologies and core network solutions are needed, moving beyond the currently dominant design of cellular systems based on human-type communications (i.e., long data streams, with dominance of downlink) and cloud computing (i.e., centralized data processing units working as X-as-a-Service [7]) toward machine-type communications and edge computing [4].

This article aims to systematically review the state-of-the-art of edge computing enabled by 5G while indicating potential future developments beyond it in relation to IIoT. We will especially focus on the emerging network architecture that is based on *pervasive edge computing* (PEC), where virtually all devices that compose the radio access network (from sensors

to edge servers co-located with the base stations/gateways) perform computational tasks. As will be discussed in more detail subsequently, the main advantages of edge computing are related to the benefits of decreasing the traffic offered to the core network (which is used to access the computational capabilities of the cloud) and providing low latency for specific applications needed in industrial automation. It is worth saying that while there are several surveys dealing with edge computing and associated concepts (e.g., [8]–[10]), they are not focused on the most recent developments related to PEC in industrial applications. Our contribution (depicted in Fig. 1) covers this current gap by articulating the recent advances in radio access and network technologies, especially those related to the important areas of artificial intelligence (AI), federated learning (FL), energy efficiency, and cyber security for different industrial applications.

The rest of this survey is organized as follows. In Section II, we first clarify the meaning of PEC and elucidate its key concepts, and then, we list some of its key applications to 5G, describing the recent advancements. Section III focuses on machine learning (ML) tools that are used to enhance the intelligence of IIoT processes enabled by PEC. Section IV deals with the question of energy efficiency in PEC. Energy
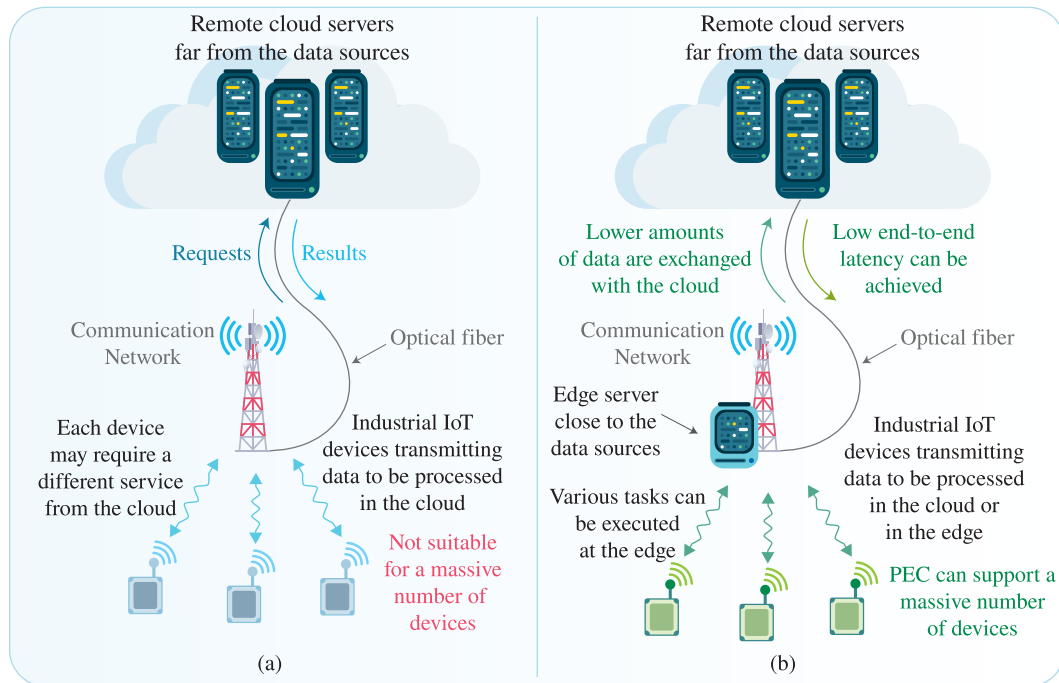
**FIGURE 2.** Cloud versus Edge Computing. (a) Typical cloud computing network where the end devices, which sense their operation environments, forward the collected data to the cloud through the communication network for further processing and inference. (b) Typical edge computing network where the edge devices in the network proximity of end devices process the collected data, and the cloud acts as a complementary processing and storage unit.

efficiency is an important issue today since the large number of distributed devices may lead to drastic increase in computational workloads and, in turn, energy usage. Further, cyber security is imperative for industries to adopt 5G and PEC, and hence, in Section V, we discuss the challenges related to cyber security and the research work addressing these challenges. In each section from Section II–V, the most relevant research works appropriate for the section are summarized in a table (Tables 2–5, respectively). Finally, in Section VI, we discuss the open challenges that still need to be addressed before we have a 5G-based scalable PEC solution for IIoT.

## II. PERVASIVE EDGE COMPUTING (PEC): KEY CONCEPTS AND APPLICATIONS TO BEYOND 5G

The term edge computing has already been widely used by the research and industrial community alike. However, a clear definition is still lacking. In this section, we provide the definition of edge computing used throughout this article explaining why we denominate it as pervasive edge computing.

The main idea behind the *edge computing* paradigm is to exploit the storage and computing capabilities of different devices at (or near) the network edge. A natural question that arises is what is an *edge*? Edge can be defined as any computing and networking resource, such as a smart phone or a 5G base station, that lies between the data source and the "cloud" (i.e., the core network). In other words, edge in this article refers to the edge of the core network including all the devices related to the radio access network. Since these (potential) edge computing elements are widely spread,

we denominate this architecture as *pervasive edge computing* (PEC), and it also incorporates the more restrictive (or fuzzy) concepts of mobile edge computing (MEC), edge servers, edge nodes, and fog network. Fig. 2 illustrates the process of edge computing (which is decentralized with computational tasks distributed among the nodes at the edge of the core network) and contrasts it with cloud computing (which is mainly centralized with computational tasks being performed in the servers at the core network and the Internet).

PEC is beneficial as it moves data processes away from centralized servers [11]. As a result, some IIoT applications do not need to send their data through the core network, avoiding congestion and potentially high delays. Moreover, PEC can also pre-process data by filtering during the acquisition phase, thereby improving the speed of data analysis and the decision-making processes [12]–[14]. Local processing provided by PEC also helps to protect sensitive data that are better processed on an edge device instead of sending to a cloud. Note that although PEC provides significant benefits to IIoT, cloud computing cannot be eliminated completely because having a centralized location for the data storage and analysis still has many benefits in different industrial application. PEC is important for offloading some tasks from the core network and also to fulfill strict latency requirements, but the remaining data may still have to be sent to the cloud for processing because of its better processing capabilities. In the following section, we will describe some of the tradeoffs of edge and cloud computing architectures.

## A. CLOUD VERSUS EDGE COMPUTING: THE KEY TRADEOFFS

As previously discussed, the edge computing paradigm has emerged not to replace conventional local and cloud computing solutions, but to complement their capabilities. The truth is that PEC alone cannot provide a universal solution that can tackle all issues of future networks. Instead, it comes with various tradeoffs that need to be investigated and well understood. Cooperation among local, edge, and cloud computing will be essential to meet the diverse requirements of IIoT. Therefore, it is important to provide the readers with a fundamental understanding of the benefits and drawbacks that can be potentially provided by the two architecture paradigms, edge and cloud computing. In this subsection, we shed light on the main tradeoffs of cloud computing and edge computing solutions.

### 1) LATENCY

When PEC comes into play, the first benefit that one can think of is the lower end-to-end latency that can be achieved. Indeed, as widely demonstrated [15]–[19], in contrast to cloud computing, the distances that data packets need to travel are, in most cases, shortened with edge servers installed closer to end-user applications; this can greatly contribute to reducing latency. However, there are scenarios in which such a benefit might not be attained. Latency does not depend only on the distance between the user and the processing server. It also depends on other factors such as the computational complexity of tasks, processing power of edge servers, and edge traffic. For instance, if the edge network is congested, or if the time spent to process the offloaded computational tasks is too high, it might be more advantageous to opt for some cloud computing solution. This tradeoff can be visualized in Fig. 3 that shows the latency versus the central processing unit (CPU) cycles required per bit offloaded by a single device in a wireless system assisted by either cloud or edge computing. This simple example is generated based on the system models proposed in [15], [16], in which the total edge computing latency can be represented by

$$T^{\text{Edge}} = \frac{b}{B^{\text{Edge}} \log_2(1 + \gamma^{\text{Edge}})} + \frac{bC}{f^{\text{Edge}}}, \qquad (1)$$

and the cloud computing latency by

$$T^{\text{Cloud}} = \frac{b}{B^{\text{Cloud}} \log_2(1 + \gamma^{\text{Cloud}})} + \tau^{\text{Cloud}}, \qquad (2)$$

where $b$ is the total number of bits of the offloaded task; $C$ is the number of CPU cycles required to compute one bit; $B^{\text{Edge}}$, $B^{\text{Cloud}}$ and $\gamma^{\text{Edge}}$, $\gamma^{\text{Cloud}}$ represent, respectively, the bandwidths and signal-to-noise ratios (SNR) of the uplink channels between device and edge server and device and cloud gateway; $f^{\text{Edge}}$ denotes the CPU's clock frequency in the edge server; and $\tau^{\text{Cloud}}$ is a fixed latency due to the long distance between the cloud gateway and the central cloud server. As shown in Fig. 3, because cloud servers dispose
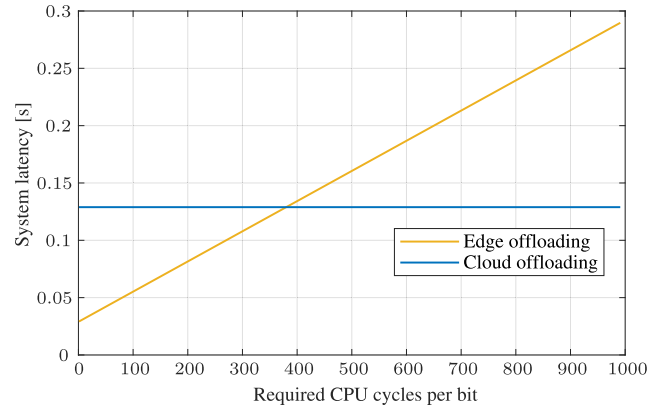


**FIGURE 3.** System latency versus required CPU cycles per bit in wireless systems assisted by edge and cloud computing ($b = 1$ Mbit, $B^{\text{Edge}} = 10$ MHz $B^{\text{Cloud}} = 10$ MHz, $\gamma^{\text{Edge}} = \gamma^{\text{Cloud}} = 10$ dB, $f^{\text{Edge}} = 6$ GHz, $\tau^{\text{Cloud}} = 100$ ms).

of abundant computational resources, the latency in the system assisted by cloud computing is not affected when the task complexity increases. On the other hand, the latency in the system assisted by edge computing escalates with the increase in the number of CPU cycles, i.e., when the computational complexity of the task becomes high.

A dynamic mobile scenario is another example where it can be challenging to provide low end-to-end latency. This is because mobile IIoT devices might not always be able to fully use edge servers (which have a fixed location) in their vicinity. As proposed in [20], a solution for this issue can be achieved by employing some multi-hop strategy so that the data can reach the nearest server. However, such an approach comes also at the cost of increased communication latency. In summary, to efficiently meet the latency requirements of different IIoT applications, factors such as task complexities, the processing power of servers, and the network topology must be carefully taken into consideration when designing a PEC-assisted network.

### 2) BACKHAUL BANDWIDTH

To satisfactorily attend a massive number of connections via a centralized cloud architecture, stringent bandwidth requirements in the backhaul lines are needed. Otherwise, severe congestion and packet losses could occur, resulting in unstable and unreliable cloud service provisioning. To alleviate such an issue, cloud providers will have to make heavy investments in communication infrastructures, which can be financially unviable. Fortunately, PEC offers cheaper efficient alternatives for reducing backhaul data traffic. By distributing the computational workload among different edge servers, lower amounts of data are required to be exchanged with the cloud. Edge data caching, data cleansing, and compression are other efficient approaches for tackling the traffic issue [21]–[23]. A combination of all these strategies can effectively relax the backhaul bandwidth requirements and decrease the costs of communication infrastructure.

The simulation examples illustrated in Fig. 4 show the potential benefits that the cooperation between cloud and
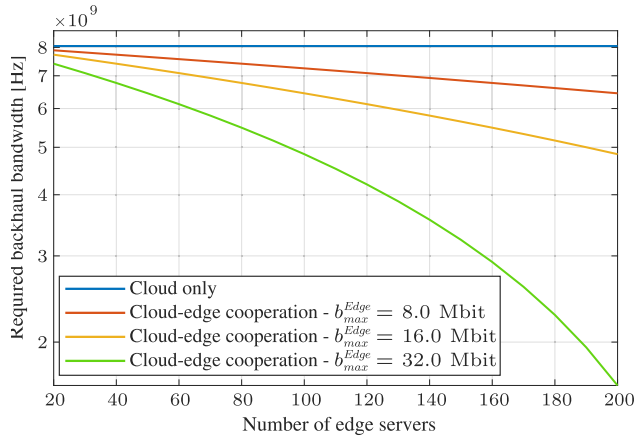
**FIGURE 4.** Required backhaul bandwidth versus number of edge servers for different values of $b_{max}^{Edge}$ ($D = 500$ devices, $\gamma^{BH} = 30$ dB, $T^{BH} = 100$ ms).



**FIGURE 5.** Outage probability curves for wireless systems assisted by cloud and edge computing considering different numbers of edge servers ($B^{BH} = 2$ MHz, $B^{Edge} = 50$ kHz, $R^{target} = 5$ Mbit/s).

edge computing can provide to reduce the required backhaul bandwidth. In these examples, we consider a scenario with $D$ devices and $N$ edge servers, in which the $i$th device offloads $b_i$ bits to be computed in a cooperative manner by cloud and edge servers. Specifically, bits are transmitted to the cloud server only when the sum of the tasks exceeds the combined computational capacity of the edge servers, i.e., when $\sum_{i=1}^{D} b_i > Nb_{max}^{Edge}$, where $b_{max}^{Edge}$ is the maximum number of bits that each edge server is able to process. With these assumptions, we compute the required backhaul bandwidth by averaging $1 \times 10^6$ realizations of the following formula

$$
B^{BH}
= \begin{cases}
\dfrac{\sum_{i=1}^{D} b_i - Nb_{max}^{Edge}}{T^{BH} \log_2(1 + \gamma^{BH})}, & \text{if } \sum_{i=1}^{D} b_i > Nb_{max}^{Edge} \\
0, & \text{otherwise,}
\end{cases}
$$

$$(3)$$

where the number of bits $b_i$ is drawn from a uniform distribution between 100 kbit and 32 Mbit, and $\gamma^{BH}$ and $T^{BH}$ are, respectively, the backhaul SNR and the desired backhaul latency. These results highlight the attractive capabilities of edge computing to alleviate backhaul requirements. As one can see, remarkable reductions in bandwidth are achieved when increasing the number of edge servers. These savings become even more prominent when the processing power of edge servers gets high, i.e., when $b_{max}^{Edge}$ is increased.

### 3) PRIVACY AND SECURITY

Privacy and security issues are critical concerns that arise with PEC-assisted systems [24]–[32]. The pervasive deployment of edge servers, geographically distributed and closer to the end users, can introduce numerous vulnerabilities to the network; moreover, these vulnerabilities can be hard to track. Edge servers may have limited computational capabilities and might lack physical protection, which creates a favorable scenario for hacker invasions and eavesdropping [24]. On the other hand, the centralized architecture of the
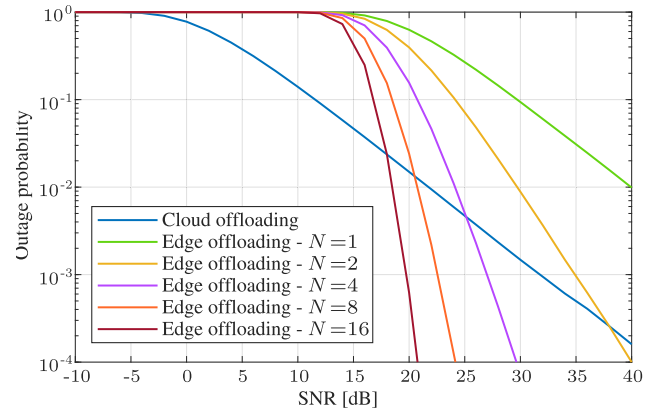
cloud computing paradigm together with the high computational power of its cloud servers enable the implementation of strong security measures, including sophisticated encryption techniques and very safe physical infrastructures. As a result, in general, it is more challenging to hack and to physically violate cloud servers [33]. An in-depth discussion highlighting important recent works on this topic is provided in Section V.

### 4) ROBUSTNESS TO FAILURES

One strong aspect of PEC is that it enables robustness to failures. Due to the branch architecture of PEC, it becomes very hard to shut down the entire network. For instance, if an electricity outage happens in a particular area of the grid, the edge computing services of other areas will continue to operate normally, unaffected. In contrast, if a given IIoT network relies solely on centralized cloud computing, when the electricity supply fails due to any natural disaster happening in the cloud infrastructure, or the backhaul communication link becomes unstable, the whole network will fail [34], [35]. To exemplify the robustness of PEC, in Fig. 5, we show the outage probability curves experienced in two wireless systems assisted by cloud and edge computing. In the first system, we consider that one IIoT device is assisted by only cloud computing, in which the device transmits its data to a gateway, such as a base station, that communicates with a central server through a wireless backhaul link with bandwidth $B^{BH}$. In particular, we assume that the cloud-assisted IIoT device experiences outage if the data rate achieved in the backhaul link is less than its target rate, $R^{target}$. On the other hand, the second system is assisted by only edge computing, where one IIoT device offloads its data to $N$ edge servers through $N$ wireless links, each with bandwidth $B^{Edge}$, such that $B^{Edge} < B^{BH}$. Differently from the first scenario, in this edge-assisted system, the IIoT device faces outage only if the data rates achieved in all $N$ links are less than $R^{target}$. For illustrative purposes, the SNR observed in the backhaul link for the cloud-assisted system is defined by $\gamma^{BH} = \frac{|h^{BH}|^2}{\sigma_n}$,

and for edge-assisted system by $\gamma_n^{\text{Edge}} = \frac{|h_n^{\text{Edge}}|^2}{\sigma_n}$, where $\sigma_n$ represents the noise variance, and $h^{\text{BH}}$ and $h_n^{\text{Edge}}$ denote, respectively, the channel coefficients for backhaul link and for the $n$th edge link, which are modeled by complex Gaussian random variables with zero mean and unit variance. Under these definitions, the outage probability for the cloud-assisted system is calculated by

$$P^{\text{Cloud}} = P[B^{\text{BH}} \log_2(1 + \gamma^{\text{BH}}) < R^{\text{target}}], \qquad (4)$$

and for the edge-assisted system by

$$P^{\text{Edge}} = \prod_{n=1}^{N} P[B^{\text{Edge}} \log_2(1 + \gamma_n^{\text{Edge}}) < R^{\text{target}}]. \qquad (5)$$

As shown in Fig. 5, even though the system assisted by cloud computing achieves the best performance when the SNR is low, the one assisted by edge computing becomes more robust in the moderate to high SNR regime as the number of edge servers increases, outperforming the cloud-assisted counterpart. These results provide a clear example of the tolerance to failures of a PEC system. Nevertheless, despite the resiliency that PEC can provide to IIoT, systemic software failures can still happen, and this can result in a generalized network collapse, as reported in [34], [36], [37].

### 5) MONETARY COST

Cloud computing providers usually charge the costumers for data transmissions, storage, and computation services [38]. Therefore, if the number of transmissions grows, if the amount of data that needs to be stored increases, or if computation tasks become more complex (which is likely to happen in IIoT), it can become excessively expensive to rely only on cloud services [39]. On the other hand, by computing and caching data and tasks at the edge, the PEC-assisted networks can decrease the traffic to cloud servers and effectively reduce the monetary costs with cloud services. However, note that even though PEC can reduce the expenditure with cloud services, shifting to a pervasive decentralized architecture can also lead to an increase in costs for the installation of new hardware and additional energy consumption [40].

These tradeoffs show that the design of a PEC-assisted network should be optimized based on the characteristics and requirements of each specific application. At the same time, standardized flexible solutions, such as those offered by 5G (and beyond) communication systems, are crucial for guaranteeing the heterogeneous quality of services of the future IIoT, making 5G a key enabler of PEC.

### B. VIRTUALIZATION AS AN ENABLER FOR EDGE COMPUTING

As we discussed earlier, edge and cloud, and the resources in between, complement each other. Depending on an application's requirements, a service provider deploys its services on the available network nodes considering their objectives, tradeoffs, and constraints. In contrast to the mostly homogeneous resources of the cloud, the edge infrastructure might be composed of multiple parties including end users in the case of PEC; this leads to an infrastructure with very diverse node properties. The challenges emerging due to this heterogeneity can largely be mitigated by virtualization techniques.

Virtualization techniques, be it using a virtual machine (VM) or a container, offer three key benefits for edge computing [41]: (i) hiding of resource heterogeneity and thus enabling of platform independence; (ii) ease of service deployment and management via resource abstraction; and (ii) isolation. First, virtualization hides the underlying hardware heterogeneity by ensuring an identical execution environment via the specification of a virtual environment [41]. Second, virtualization facilitates resource abstraction and hardware emulation, thereby simplifying the interaction between the services and the underlying hardware. A hypervisor handles the resource management for the services running in a virtual machine. Finally, virtualization achieves various levels of isolation, e.g., at hardware level or operating system level, among the different services hosted on the same node. Note that virtualization approaches might differ from each other in their storage overhead, memory cost, and initialization latency according to the virtualization is implemented; for example, a VM hosts its own OS and therefore is considered heavyweight compared to a container that uses the hardware and the kernel of its host machine [41]. Lightweight containers have gained higher support for edge computing due to their lower resource footprint in comparison to VMs [42].

The deployment may become suboptimal with evolving network dynamics, e.g., changes in the user or edge node locations, or in the network conditions, leading to inefficient operations. Hence, services need to be swiftly migrated to other network locations, e.g., closer to the request locations. Note that a service can consist of multiple tasks and the tasks can be deployed on different nodes based on their computation and communication requirements as well as the dependencies among the tasks. Hence, a service provider has to profile the requirements of the sub-tasks of a service, decide on which network nodes to deploy each task, and migrate the tasks to the new locations seamlessly or with minimal impacts on the ongoing sessions [43]. While computation offloading offers many benefits to the resource-limited devices by augmenting them with resourceful network nodes, it poses significant challenges such as smooth service migration especially for mobile end users. Stateful service migration includes transfer of both the execution environment, e.g., a VM or container, and the application-related data such as runtime memory states. Hence, when the data to be transmitted from one host to another is large, it will result in a long migration delay. To ensure service continuity, live migration is desirable as opposed to cold migration which suspends a service during the time the service is migrated to another host [44].

There is increasing literature on improving service migration performance, e.g., [45], [46], and [42]. Ma *et al.* [45], [46] proposed to leverage the layered storage system of Docker containers: the base image layers of a container can

**TABLE 1.** Description of acronyms used in "PEC in 5G and Beyond: Applications" section .

| Acronyms | Brief Explanation |
|---|---|
| COP | Concatenated orthogonal preamble; a preamble structure used in the frame transmission of wireless communication networks, where the preamble is divided into multiple orthogonal sub-preambles. |
| FOVs | Field of view; the extent of the observable world that is seen at any given moment. |
| GF-NOMA | Grant-free non-orthogonal multiple access; a joint combination of both schemes, grant-free access and NOMA. |
| JPORA | Joint partial offloading and resource allocation; an optimization technique used to reduce the task execution latency of a specific computational process. |
| KPIs | Key performance indicator; a measurable value that demonstrates how effectively a company is achieving key business objectives. |
| MHVA | Multi-hop VANETs-assisted offloading strategy; a specific case of VANET in which multi-hop communication between different networks is considered. |
| MIMO | Multiple-input multiple-output; physical layer technique used in wireless communications for multiplying the capacity of a radio link using multiple transmission and receiving antennas. |
| mMIMO | Massive MIMO; an extension of MIMO, which improves the spectrum efficiency because it uses an extensive number of antennas ($L$) in the access point, when compared to the number of users ($K$) (in other words $L >> K$.). |
| mGFRA | Massive MIMO based grant-free random access; a complement of mMIMO with grant-free random access, which aims to improve the uplink channel access in terms of latency by avoiding the traditional random access procedure used in cellular networks. |
| MINLP | Mixed integer nonlinear programming; a type of optimization problem with continuous and discrete variables and nonlinear functions in the objective function. |
| NOMA | Non-orthogonal multiple access; wireless communication scheme that aims to enhance the spectrum efficiency. Using the same resource in terms of time, frequency, and space, it serves multiple users simultaneously. |
| OFDMA | Orthogonal frequency-division multiple access; a popular multiple access technique, characterized by the orthogonality of the resources. It is used in the physical layer design of 4G and 5G. |
| PST-ResNet | Deep spatio-temporal residual networks with a permutation operator; a particular deep-learning framework used in some cooperative autonomous driving solutions. |
| QoS | Quality of service; a popular network characteristic that aims to support specialized requirements in terms of packet loss, latency, and jitter. |
| SC-DSCS | Subpopulation collaboration based dynamic self-adaption cuckoo search; a search algorithm, in which the population of cuckoos is divided into two subgroups. |
| SDN | Software-defined networks; a novel technology in which the control plane of a network is configured by software. |
| SDR | Software-defined radio; a technology in which the radio technology is configured using software. |
| SIC | Successive interference cancellation; a decoding technique used to receive two or more packets simultaneously, mostly in a full interference scenario. |
| SJD | Successive joint decoding; a transmission reception strategy in which the source sequence uses the successfully decoded layers and the side information sequence. |
| SOP | Single orthogonal preamble; traditional preamble structure used in cellular networks, such as 4G. It is used to avoid collisions during the random access procedure. |
| URLLC | Ultra-reliable low-latency communication; used to reduce physical layer latency in 5G networks and beyond. |
| VANETs | Vehicular ad hoc network; a network terminology used in vehicular applications, in which the networks can be formed and information can be relayed among cars. |

be downloaded before migration from a cloud server while the container layer and runtime data are transferred from the source host to the destination host after initiating the migration. The authors showed a significant decrease in the migration time and transferred data size for an example scenario and discussed that pipelining can also introduce further improvements in the migration delay. Bellavista *et al.* [42] proposed a flexible service migration framework that can operate in various modes, e.g., application-agnostic vs. application-aware. They showed that understanding the service characteristics and leveraging certain properties helps to migration latency. Please refer to [47] and [44] for an

elaborate discussion on service migration, [41] for virtual machine management, and [43] for computation offloading approaches.

### C. PEC IN 5G AND BEYOND: APPLICATIONS

The discussions within 3rd Generation Partnership Project (3GPP) indicate the path for current and future developments of 5G and beyond technologies. We identified four types of PEC applications, which are enabled by these technologies. For greater clarity, the abbreviations used in this section are listed, expanded, and explained in Table 1.

**TABLE 2.** Summary of the relevant research work related to Section II-C "PEC in 5G and Beyond: Applications".

| Research work(s) | Topic(s) | Brief overview |
|---|---|---|
| [17], [48], [49] (2020) | Mission Critical Applications | These studies examine (i) spectral efficiency on massive MIMO based on grant-free random access; (ii) collision treatment on grant-free non-orthogonal multiple access; and (iii) minimization of latency used on the execution of execution tasks when OFDMA is used in the communication network. |
| [21], [50], [51] (2019) | Augmented Reality / Virtual Reality | In these studies, (i) joint radio communication, caching, and computing decision problem are proposed to optimize resource allocation at edge access points and mobile VR devices; (ii) underlying dynamic rendering-module placement problem on mobile VR group gaming services using a model predictive control is studied; (iii) closed-form expression is proposed for the optimal joint policy that identifies key tradeoffs in terms of computing, communication, and caching capabilities, when FOVs are homogeneous. |
| [6], [52]–[56] (2019–2020) | Network Optimization | Key features of 5G, such as network flexible radio access network slicing, automated radio access technology selection, and mobile edge caching and content delivery are addressed. In the context of the core network, network slicing is optimized by machine learning methodologies to achieve the desired network orchestration. In other aspects, PEC plays a key role in the definition of future technologies, such as cell-free massive MIMO. |
| [20], [57]–[59] (2019–2020) | Vehicular Computing Applications | In these studies, (i) a computation offloading scheme for vehicles is proposed using multi-hop vehicular ad hoc networks; (ii) a cooperative autonomous driving based on PEC is studied; (iii) A prototype system using a 5G network and a PEC architecture is developed to support 3D dynamic map services; (iv) PEC is mentioned as a key element in 5G solutions that aim to support AR, VR, and URLLC use cases. |

### 1) MISSION-CRITICAL APPLICATIONS

An active PEC-related research field is concerned with mission critical settings that have strict requirements of latency, availability, and reliability. For example, URLLC in 5G aims to support mission-critical applications. Mission-critical applications are almost always related to industrial applications that have greater emphasis on feedback control loop and automation.

Different advanced multiple access techniques are being proposed to improve the system performance in terms of latency and reliability. Grant-free access and non-orthogonal multiple access (NOMA) are examples of such promising techniques. For instance, in [48], a closed-form expression for the spectral efficiency of two preamble structures in a multiple-input multiple-output (MIMO) based grant-free random access (mGFRA) scenario was obtained. The first preamble structure was named as concatenated orthogonal preamble (COP) and the second as single orthogonal preamble (SOP). The authors concluded that there is a threshold between both preambles in terms of the number of antennas in the massive MIMO (mMIMO) scheme. In [49], the authors proposed a framework to treat collisions in a grant-free NOMA (GF-NOMA) scheme. The authors used Poisson point processes and order statistics to derive simplified expressions of the outage probability and throughput of the

system for both successive joint decoding (SJD) and successive interference cancellation (SIC).

Another key research area is the minimization of task execution latency. A mathematical model of the minimization of the sum of task execution latencies of devices, which operate under interference, was presented in [17]. Here, the authors provided an integrated framework for partial offloading and interference management using the orthogonal frequency-division multiple access (OFDMA) scheme. They formulated the total latency of minimization as a mixed integer nonlinear programming (MINLP) problem, in which desired energy consumption, partial offloading, and resource allocation constraints were considered. A novel iterative scheme named joint partial offloading and resource allocation (JPORA) on data segmentation was proposed to optimize the Quality of service(QoS)-aware communication. JPORA obtained the lowest latency as compared to other baseline schemes, while simultaneously reducing the energy consumption in the devices.

### 2) AUGMENTED AND VIRTUAL REALITY APPLICATIONS

5G radio access and computational resources must be brought closer to augmented reality (AR) and virtual reality (VR) applications, which will also employ PEC. Since these immersive media applications require very low latency,

typically lower than 20 ms, the delay budget of the network to deliver the requested content is very tight. Edge servers can both reduce the computation time and energy consumption on the end devices by performing computation-heavy tasks such as rendering or object detection in AR. The authors of [21] formulated a joint radio communication, caching, and computing decision problem to optimize resource allocation at edge access points and mobile VR devices. In [50], the underlying dynamic rendering-module placement problem on mobile VR group gaming services was studied using model predictive control. Here, the authors proposed a methodology based on graph theory to explore the connection between the placement problem and the minimal s-t cut problem.

Since a ultra-high transmission rate is required for immersive media applications, the authors of [51] formulated a joint caching and computing optimization on a PEC-based mobile VR delivery framework to support diverse fields of view (FOVs) in advance. They proposed a closed-form expression for the optimal joint policy that identifies key tradeoffs in terms of computing, communication, and caching capabilities, when FOV are homogeneous. In the case of heterogeneous FOVs, the framework considered a local optima transformation into a linearly constrained quadratic problem. It is important to remark that AR and VR applications are very important to the field of intelligent machines and remote control, which can, for example, can increase the safety of operations in mining and maritime vessels.

### 3) NETWORK OPTIMIZATION APPLICATIONS

In 5G cellular networks, edge computing is a key enabler to support specialized key performance indicators (KPIs) such as low latency, high connection density, and bandwidth efficiency. In addition, the evolution of virtual network functions (VNF) running on general purpose edge infrastructure creates novel technical possibilities, such as virtualization of a portion of the access network with functionalities deployed close to the end users. In the context of network slicing, which facilitates the creation of a logical end-to-end isolated network to support specific applications, the authors in [52] established a combinatorial optimization model that natively supports multiple network slices with different QoS requirements. The algorithm addressed a combinatorial problem that was a multi-period variant of the generalized assignment problem. They also showed that the network performance benefits from a multi-sliced approach that is more suitable for capturing the distinct spatiotemporal pattern of each slice than conventional single-slice models.

The distributed network topology that characterizes PEC and new improvements in software defined networks (SDN) and software defined radio (SDR) will enable cutting-edge technologies, with the potential to improve the performance of wireless communication systems beyond the current KPI requirements. One promising research topic is orchestration on network slicing within private 5G networks with local operators [6]. This featured approach, which is implemented

via advanced ML techniques, aims to optimize the performance of complex systems such as the combination of different sub-slices on the end-to-end network slice. Shen *et al.* has highlighted some benefits and potentials of AI-based techniques on next-generation wireless networks [53]. Three challenging scenarios were addressed using AI, including flexible radio access network slicing, automated radio access technology selection, and mobile edge caching and content delivery.

A promising beyond 5G technology is cell-free massive MIMO, which aims to deploy distributed access points in contrast to the traditional cellular deployment of current broadband wireless networks. In [54], a PEC implementation was considered with a diversity of computational/processing requirements of the users. The authors considered access points with PEC servers and a central server with cloud computing capability. It is also important to consider that local industrial network operators for 5G are key enablers for realizing the PEC to its fullest extent as the most suitable way to flexibly manage the network resources in virtual slices reserved to different applications [55]. Finally, PEC is expected to support every layer of mobile networks to renovate the common computing architecture, based on the analysis of the wireless big data, as stated in [56].

### 4) VEHICULAR COMPUTING APPLICATIONS

The safe and efficient deployment of autonomous vehicles is a key application area of PEC and 5G networks, and it has been attracting strong research interest [57]. For example, cooperative autonomous driving based on PEC was proposed in [58]; here, the authors developed a prototype system that was based on a 5G next-generation radio access network, a PEC server providing high definition 3D dynamic map service, and a cooperative driving vehicle platoon. They also performed several field tests, where the results indicated that the combination of 5G-vehicle-to-everything (V2X), PEC, and cooperative autonomous driving can provide important improvements to the system. However, these practical deployments also presented some challenges. To address them, the authors proposed two AI-based approaches. The first one was a deep-learning-based tool called deep spatio-temporal residual networks with a permutation operator (PST-ResNet), and the second was a swarm intelligence-based optimization tool called subpopulation collaboration based dynamic self-adaption cuckoo search (SC-DSCS).

In [20], the authors proposed a computation offloading scheme for vehicles using the multi-hop vehicular ad hoc network (VANET), called multi-hop VANETs-assisted (MHVA) offloading strategy. A reliability model for multi-hop routing was developed using link correlation theory applied to VANET. The offloading strategy based on a binary search algorithm was optimized in order to identify the lowest relaying and reduced computing cost. The simulation results of multi-hop VANETs-assisted (MHVA) offloading strategy showed better offloading performance in terms of delay and cost, when compared to typical strategies. Although vehicular

technologies are being studied in scenarios related to transportation, similar solutions could be used in private industrial networks considering mobile (autonomous) robots or large-scale plants related to, for example, mining.

PEC will play a key role in beyond 5G technologies. As in 5G, it is expected that PEC will operate as an intermediate layer providing low latency and local data processing for critical and resource constrained applications. AR, VR, network optimization, V2X communication, energy efficiency, and offloading for URLLC can be cited as particularly relevant use cases [59].

## III. ARTIFICIAL INTELLIGENCE IN PEC

IoT devices are normally heterogeneous devices and difficult to coordinate. As a result, a major challenge in PEC is the efficient resource co-ordination and communication between different types of heterogeneous edge devices, from computers equipped with graphical processing units (GPUs) to smart phones with mobile processors to devices with just small single-board computers like Raspberry Pi [67]. Secondly, edge devices have to coordinate with the cloud under dynamic network conditions and different settings to ensure satisfactory application-level performance. Several AI techniques have been proposed to meet these challenges. Among them, deep learning (DL) is very promising due to its potential to create hybrid approaches that combine cloud and edge computing.

In this section, we first briefly present the recent status of DL research applied to PEC before moving to our major focus, the relatively new concept of federated learning (FL). FL is an extremely important technique that has a nearly symbiotic role with PEC for addressing privacy issues. Privacy is a major concern with smart devices because of the need to interact and share data with the cloud and third-party platforms. FL takes advantage of PEC to maintain user privacy while performing the required computations efficiently.

### A. DEEP LEARNING AND APPLICATIONS TO PEC

The application of DL at the network edge offers many improvements to PEC as DL and PEC together can support numerous applications, including computer vision, smart surveillance, natural language processing, network functions, and VR and AR. In a recent paper, Chen *et al.* surveyed research literature at the confluence of DL and PEC [67]. They highlighted that the researches so far have been based on three major architectures: (1) centralized edge server-based architectures, where data from end devices are sent to one or more edge servers for computation; (2) semi-distributed architectures in which the computation is shared among end devices, edge servers, and the cloud via a joint computation mechanism; and (3) distributed architectures based on on-device computations, where deep neural networks are executed on the end device itself.

In the future, edge intelligence is expected to move DL implementations from the cloud to the edge, forming the so-called edge DL [68], [69]. Research into the industrial

applications of this convergence of DL and the edge is still nascent but is expected to increase significantly, especially with the implementation of 5G. In [70], the authors proposed an edge computing-based DL model that migrates the DL process from the cloud to the edge in an IIoT network using edge computing concepts. The DL model is optimized so that computational power requirements are reduced. By deploying a testbed implementation with their proposed convolutional neural network model and real-world IIoT dataset, the authors showed that network traffic overheads are reduced without compromising the classification accuracy.

In another study, a DL-based method was used to detect hazardous conditions in supermarkets, such as spilled liquids or fallen items on floors [71]. They showed that their lightweight DL model can be deployed on edge devices for quick computations. Similar studies based on intelligent visual recognition using DL implemented at the edge have been applied to industrial electrical equipment [72], health monitoring [73], and flat surface texture inspection [74].

### B. FEDERATED LEARNING AND APPLICATIONS TO PEC

FL is a relatively recent ML paradigm that was motivated by the need to protect user privacy. FL aims to train a centralized model using training data that are distributed over a large number of client devices that themselves are a part of the training [66]. In other words, FL differs from classical ML learning approaches such as DL in the fact that the model training does not use a single processing device [75], [76]. In FL, end devices use their local data to cooperatively train an ML model (using ML techniques such as DL, support vector machines (SVM), artificial neural networks (ANN), etc.) required by an FL server. They then send the model updates to the FL server for aggregation. These steps are repeated in multiple rounds until a desirable accuracy is achieved [60]. Fig. 6 illustrates the difference between DL and FL. In general, FL involves the training of statistical models directly on remote devices, and it is motivated by the need to maintain information privacy [77].

Such a decentralized approach facilitates collaborative complex ML techniques while guaranteeing that the data will remain in the personal devices, thereby preserving privacy. Note that there is usually an underlying assumption that the data owners are honest: they use their *real* private data to do the training and submit the true local models to the FL server. The main advantages of the FL approach are as follows:

- *Bandwidth efficiency*: Less data is required to be transmitted to the cloud.
- *Privacy*: Raw (local) data need not be sent to the cloud any more.
- *Low latency*: Real-time decisions can be made locally instead of being made in the cloud.

An important question in FL research is whether its performance is comparable to that of traditional DL-based approaches for resource coordination. In [63], the authors examined this question by employing FL for the joint

**TABLE 3.** Summary of relevant research works on artificial intelligence in PEC.

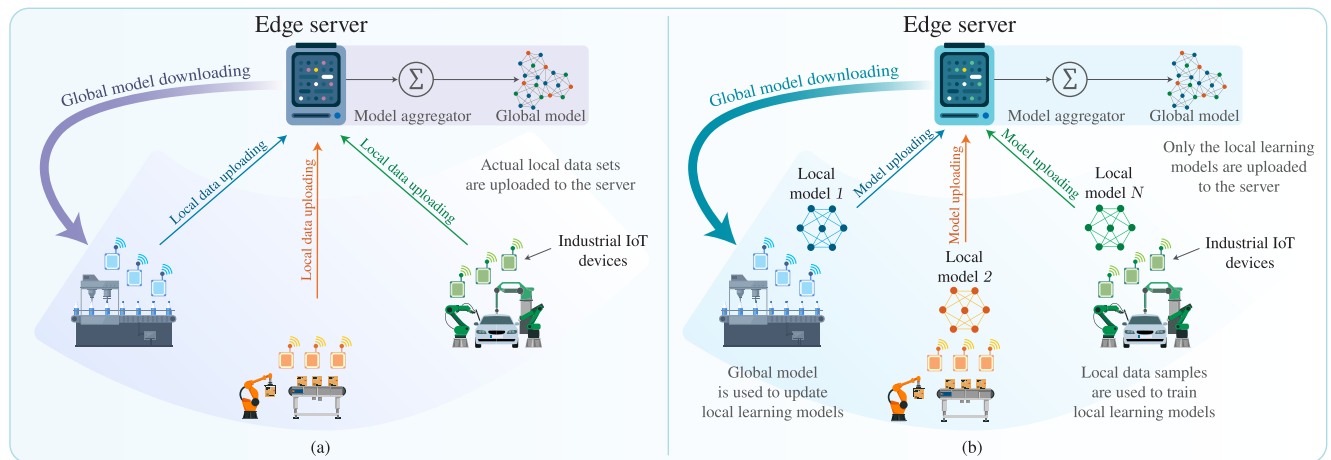| Research work(s) | Work area(s) | Brief overview |
|---|---|---|
| Lim et al. [60] (2020) | Survey focusing on federated learning (FL) in PEC and fundamental concepts | Highlights the challenges of implementing FL in PEC; reviews existing solutions; presents the applications of FL for mobile edge network optimization; and discusses the future directions in FL. |
| Fadlullah et al. [61] (2020) | Proposal of a heterogeneous platform for FL including 6G networks | A distributed heterogeneous computing platform is proposed, and privacy of information in UAV's is preserved by using a a two-stage FL algorithm. To avoid redundant learning transfer in the FL, the authors use an asynchronous weight updating method. The communication between ground and aerial vehicles is performed using a 6G tiny cell. |
| Lu et al. [62] (2020) | Proposal of an asynchronous FL in mobile edge computing | To address the security of FL in vehicular networks, a model is proposed to protect each update of local models by verification and weighted aggregation. Real-world datasets are used for the evaluation, and the model demonstrates high efficiency and high accuracy. |
| Ren et al. [63], [64], [65] (2019, 2020) | Proposal to train ML algorithms in IoT devices and vehicular edge computing | The paper shows how FL can decrease the transmission costs of devices and central server, as well as the utility of the latter, as compared with baseline approaches. |
| Konecny et al. [66] (2016) | Proposal of a strategy to improve communication networks in FL | The authors propose reductions in the uplink communication costs using two approaches: structures updates and sketched updates. The paper uses both convolutional and recurrent networks with FL occurring at the edge nodes (mobile devices) and aggregation occurring at a central server. |



**FIGURE 6.** (a) Traditional deep learning: end devices send their local data to a global server that collects all data and performs the complete model training. (b) Federated learning: end devices use local data to cooperatively train a local model and then send it to the global server for aggregation. These steps are repeated in multiple rounds until a desirable accuracy is achieved.

allocation of communication and computing resources by guiding the training of deep reinforcement learning agents. The IoT devices in their model harvested energy units from edge nodes and stored them in their energy queue. The authors demonstrated that the fluctuation range of FL-based DL with respect to utility variation is bigger than that from centralized training. Their results confirm that the performance of FL-based DL training for computation offloading approaches the results from centralized DL training.

The efficient utilization of limited computation and communication resources to increase the optimal learning

performance of different applications was examined in [65]. The authors considered a typical edge computing architecture where the edge nodes are interconnected with the remote cloud via network elements such as gateways and routers. The raw data was collected and stored at multiple edge nodes and FL learning was performed. In the FL approach in this work, the frequency of global aggregation was configurable; that is, it was possible to aggregate at an interval of one or multiple local updates. Each local update consumes computation resources and each global aggregation consumes communication resources of the network. The amount of

consumed resources may vary over time, and there is a complex relationship among the frequency of global aggregation, the model training accuracy, and resource consumption. This is a tradeoff between the model resource optimization and precision. In this work, the convergence bound of gradient-descent based FL was analyzed from a theoretical perspective, and a control algorithm that learns data distribution was proposed. This algorithm was tested using real datasets both on a hardware prototype and simulated environment.

Based on the distributed topology that characterizes the PEC paradigm, different architectures are being considered to deploy FL techniques in order to improve the trade-off between energy consumption, computation capacity, and training capabilities of distributed learning nodes connected by a customized communication network. For instance, in [61], a two-stage FL algorithm was proposed for a collaborative scenario between user equipment (UE), unmanned aerial vehicles (UAV)/ base stations (BS), and a heterogeneous computing platform (HCP) to predict content caching. Here, an asynchronous weight updating method was adopted to reduce the effects of redundant learning in FL.

FL in PEC also increases the reliability and security for some critical applications such as vehicular edge computing (VEC). VEC faces a major challenge that the accuracy of image quality can decrease during the model aggregation. To address this, the authors in [64] proposed a selective model aggregation exploiting a geometric model that illustrates the relationship between the object of interest and the camera in each vehicular agent. FL was used to train image classification and to tackle the asymmetry caused by it, and a model selection procedure was formulated as a two-dimensional contract theory problem. Then, the contract problem was transformed into a problem that can be tracked by relaxing and simplifying the complicated constraints. In [62], FL was applied to urban informatics tasks where the vehicular network consisted of a macro base station, a number of roadside units, and moving vehicles. The authors focused on three aspects: enhancing the privacy of the updated models in FL, development of a new asynchronous FL architecture by leveraging distributed peer-to-peer update schemes, and boosting the proposed FL process. The proposed FL method not only gave higher accuracy than the compared methods such as convolution networks, GraphStar, and text graph convolution networks, but it could also execute parallel local training; moreover, increasing the number of data providers did not affect the accuracy.

## IV. ENERGY EFFICIENCY IN PEC
Despite the benefits of PEC, this new computing paradigm also raises concerns such as energy efficiency. To accommodate the massive number of IIoT connections efficiently, a large number of distributed servers must be installed. As a result, energy consumption can increase drastically if the resources and computational workload are not well distributed within the PEC-assisted IIoT network. Moreover, as illustrated in Fig. 7, in contrast to the conventional
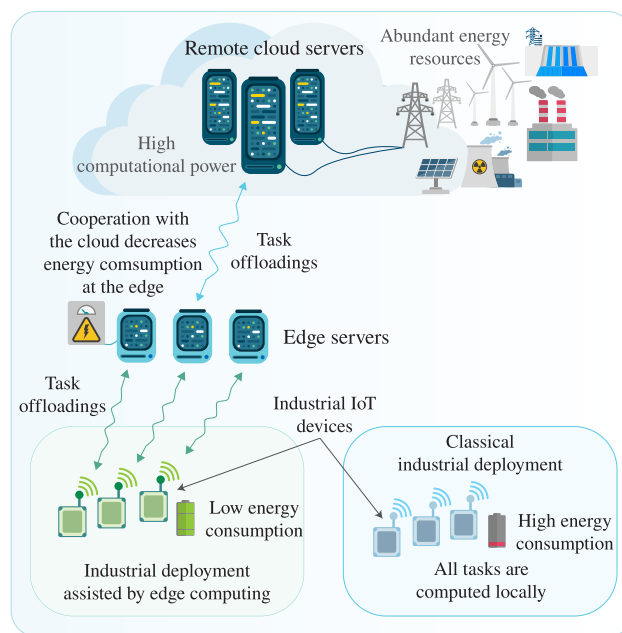


**FIGURE 7.** Energy consumption in industrial IoT (IIoT) deployments. Edge computing can provide remarkable energy savings to IIoT devices, and the cooperation with cloud computing can reduce energy consumption at the edge servers.

remote cloud servers that possess abundant energy resources, IIoT devices and edge serves may have access to only limited power supplies and limited computational capabilities. Therefore, energy efficiency has become a critical concern calling for an energy-centric design of PEC solutions. We will now present and discuss the most recent literature that is relevant to energy efficiency.

### A. CLOUD-EDGE COOPERATION FOR REDUCED ENERGY CONSUMPTION
The work in [78] addressed an industrial scenario where multiple IIoT devices are assisted by both edge and remote central servers. Aiming to minimize the energy consumption at the edge, the authors developed two optimal dynamic algorithms for offloading the computation-intensive tasks from the edge nodes to the remote server. Specifically, the first algorithm used an improved gradient method for achieving faster convergence, while the second one employed the concept of dynamic voltage scaling for further maximizing the energy gains at the edge servers. Simulation studies in [78] showed that the proposed algorithms outperform the conventional approaches both in terms of energy consumption and convergence time.

In [15], the authors studied a heterogeneous network consisting of a macro base station equipped with a central server, multiple multi-antenna small base stations equipped with edge servers, and multiple single-antenna energy-constrained devices with low computational capabilities. The authors provided an optimization approach to minimize the energy consumption of the network by jointly optimizing the devices' transmit powers, server selection, and edge servers' receive

beamformers. To solve the non-convex formulated problem, the authors proposed an iterative algorithm based on decomposition and successive pseudo-convex approach.

In some IIoT settings such as in hard-to-reach mining sites, it might be very challenging to provide single-hop connectivity. Hence, enabling PEC also for multi-hop networks is necessary for IIoT. Reference [79] focused on such a setting where IIoT devices can operate as relays in a distributed cooperative edge–cloud network. This work applied game-theoretic concepts to develop two distributed offloading algorithms for minimizing the tasks' computation time and energy consumption. The authors showed that their proposal can ensure the required QoS of each IIoT device and achieve the Nash equilibrium. Their numerical results demonstrated the superiority of the developed algorithms over benchmark schemes in terms of stability, scalability, time processing, and energy gains. With the idea of distributing the computational tasks among the local edge server, neighbor edge servers, and cloud, the authors of [80] formulated an optimization problem for minimizing the network's energy consumption under per-task time processing constraints. To solve the problem, the authors devised an optimal task allocation algorithm based on Lyapunov drift-plus-penalty theory for queuing systems. The obtained solution showed significant energy efficiency gains and improvements in end-to-end latencies.

### B. REDUCED ENERGY CONSUMPTION IN IIoT DEVICES

While the researches discussed in the previous subsection focused on enhancing the energy efficiency of the entire network and is mainly concentrated on the servers, in this subsection, we survey important contributions that optimize the energy consumption at the IIoT devices. For instance, the work in [81] employed a task caching strategy at the edge servers to avoid unnecessary offloading transmissions and computations of repeated tasks. Based on the proposed caching-enabled system design, and aiming to minimize the devices' energy consumption under delay constraints, the authors formulated a problem to jointly optimize caching, computation, and communication resources. The problem turned out to be a mixed-integer non-convex optimization formulation, which is difficult to solve optimally. To tackle the complex formulation, block coordinate descent and convex optimization techniques were explored and a sub-optimal iterative algorithm was developed. The simulation results showed that substantial energy efficiency gains and a significant reduction in processing time can be achieved with the proposed solution.

The energy efficiency of IIoT devices with multitask capabilities was studied in [82]. In the proposed system model, each device was able to simultaneously offload multiple tasks with different requirements to the edge servers using a NOMA technique. Specifically, a two-step optimization approach was proposed for minimizing the total energy consumption at the IIoT devices, i.e., energy spent with NOMA transmissions and local task computations. In the first step, a problem was formulated for jointly optimizing the number of bits to be offloaded, the computation rate, and the transmission time for each task. Even though the original problem was non-convex, after several simplifications, the authors developed a layered algorithm for computing the optimal solution. In the second step, for further minimizing the devices' energy consumption, an index-swapping algorithm was developed for determining the optimal assignment of the tasks to the most suitable edge servers. Numerical results revealed the effectiveness of the proposed algorithms for improving the energy efficiency of IIoT devices.

PEC-assisted vehicular networks with device-to-device offloading capabilities were investigated in [83]. In this work, by modeling both the network mobile traffic and the computational workload, the authors investigated the trade-offs between energy consumption and system delay. In the considered network design, smart vehicles were capable of optimizing their energy consumption and computation time by properly offloading their most demanding tasks to edge servers or to neighboring vehicles that were willing to share their computational resources. To find such an optimal offloading strategy, the concepts of Markov decision processes were exploited to formulate an energy/time cost minimization problem. Then, two reinforcement learning-based algorithms were proposed to compute the optimal solution of the formulated problem. Simulation results showed that the energy and delay performances achieved with the developed solutions were remarkably superior to those observed with the considered benchmark schemes.

The work in [97] also studied the energy efficiency of PEC in a vehicular network. However, differently from the previous reference, the authors focused on reducing the energy consumption of in-vehicle battery-powered devices and not on the vehicle itself. Specifically, a task offloading optimization problem was formulated to minimize the energy consumption of in-vehicle devices under energy and latency constraints. Due to the fractional form of the objective function and the complicated constraints, the optimization problem turned out to be NP-hard. To tackle such a challenge, the complex original problem was transformed into an equivalent consensus problem with separable objectives. The transformed problem was further decomposed into multiple tractable sub-problems, which allowed the authors to achieve a low complexity solution based on the alternating direction method of multipliers. The obtained solution enabled the in-vehicle devices to solve the sub-problems simultaneously in a distributed fashion. To demonstrate the effectiveness of the developed solution, the authors provided simulation results based on a real-world topology. The results showed that significant reductions in energy consumption were achieved with the proposed approach.

## V. CYBER SECURITY CHALLENGES IN PEC

Despite the rapid growth in research and many technological advancements, PEC still faces numerous problems with respect to security risks and privacy challenges, as discussed in [24]–[32]. These challenges are particularly critical in

**TABLE 4.** Summary of relevant research works on energy efficiency in PEC.

| Research work(s) | Work area(s) | Brief overview |
|---|---|---|
| X. Hu *et al.* [15] (2020) | Heterogeneous cellular networks | Investigated the integration of cloud and edge computing for reducing the energy consumption of a two-tier heterogeneous cellular network. |
| S. Chen *et al.* [78] (2019) | Cloud-edge task computation offloading in Industrial IoT (IIoT) | Proposed an accelerated gradient algorithm with Lagrangian dual theory to minimize energy and delay overheads of an IIoT network. |
| Z. Hong *et al.* [79] (2019) | Game theory in multi-hop IIoT–edge–cloud computing systems | Exploited game theory to develop two distributed offloading algorithms for minimizing the tasks' computation time and energy consumption in a multi-hop cooperative IIoT network assisted by edge and cloud computing. |
| P. Liu *et al.* [81] (2019) | Task caching in edge computing | Proposed a task caching strategy at the edge servers and minimized the devices' energy consumption by jointly optimizing caching, computation, and communication resources under delay constraints. |
| Y. Wu *et al.* [82] (2020) | NOMA for computation offloading in multitask-enabled IIoT devices | Proposed a two-step optimization approach for minimizing the total energy consumption of multitask-enabled IIoT network devices, including energy spent with NOMA transmissions and local task computations. |
| Y. Wang *et al.* [83] (2019) | Device-to-device offloading in PEC-assisted vehicular networks | Exploited concepts of Markov decision processes to minimize energy and delay costs in a vehicular IIoT network, considering mobile traffic and computational workload. |

**TABLE 5.** Summary of relevant research works on cyber security in PEC.

| Research work (s) | Work area (s) | Brief overview |
|---|---|---|
| Yahuza *et al.* [25] (2020) | A systematic review of edge computing paradigm emphasizing privacy and security requirements in the light of state-of-the-art technologies and solutions | Comprehensively elaborated security requirements in edge computing including the technological trends and taxonomy of security attacks; described state-of-the-art solutions to mitigate cyberattacks (considering emerging technologies) and presented future research opportunities. |
| Jiale *et al.* [24] (2019) | A survey focusing on the challenges and countermeasures of data security in edge computing | Provided various concepts and architectures of edge computing; discussed data security requirements and solutions in edge computing particularly based on cryptography approaches. |
| Yinhao *et al.* [25] . (2018) | Presentation of the basic structure of the edge computing paradigm and most influential privacy and security attacks | Described layers of edge computing; root causes of the four most influential security attacks; and defensive mechanisms to limit security breaches in PEC. |
| [84]–[88] (2011, 2017–2019) | Background of DDoS attacks in IIoT and solution designs and frameworks for the mitigation of DDoS | Presented an overview of the DDoS attacks along with various types of DDoS attacks and proposed solutions to limit the impact of DDoS |
| [89], [90] (2013, 2015) | Proposal of an infrastructure to detect anomalies (in cloud computing), particularly malware injection and DDoS. | Proposed anomalies detection techniques and discussed real-world case studies of the various attacks and feasible solutions. |
| [91]–[96] (2015–2018) | Discussion on the impacts of side-channel attacks in cloud and edge computing, and the essential countermeasures | Surveyed side-channels attacks, its types, and countermeasures, and proposed solutions to mitigate various adversarial attacks. |

industry applications. In this section, we overview the studies focusing on the security and privacy aspects of PEC from a holistic perspective, keeping in mind that the discussions apply also to industrial applications.

To systematically investigate the security aspects and to identify potential security risks of PEC, Xiao *et al.* [24] provided a four-layer architecture as follows: edge server security, network security, devices security, and infrastructural security. On the edge server side, the security concern could be, for example, that adversaries attempt to access the

edge servers and manipulate the services or that they control the edge servers and exploit their privileges even as legitimate administrators. Subsequently, attackers can execute attacks such as denial of service, man in the middle, etc. Similarly, edge devices could be infected by adversaries with a malware injection, with these malicious devices consequently posing security challenges to edge servers, edge networks, and core infrastructure.

In this vein, Zhang *et al.* [25] elaborated on the security challenges faced by edge computing and their defense

mechanisms. The authors inferred the root cause of security threats and challenges and explained the question of *why are these attacks more common in edge computing than traditional cloud computing?* The key reasons discussed in the paper are as follows: (i) weak or low power computation, (ii) resource constraints, (iii) protocol heterogeneity, and (iv) distributed access control.

In a similar vein, Yahuza *et al.* [26] thoroughly reviewed the cyber security aspects in PEC and pointed out the various possible cyberattacks in the edge computing paradigm, including message alteration, camouflaging, networking attacks, physical attacks, and reputation tarnishing. Also, the authors discussed two types of methods to evaluate these attacks—"with tools" including algorithmic proofs, simulations, and prototype implementations, etc., and "without tools" including mathematical analysis and informal security proofs. Further, many other works have highlighted the aforementioned discussion from different perspectives, for instance, edge security in terms of data analytics [27], secure IoT services [28], and others [29]–[32]. It is worth mentioning that such challenges could be due to misconfiguration, design flaws, implementation bugs, data correlations, and missing fine-grained access controls.

The four most important state-of-art security challenges in PEC are distributed denial of service (DDoS) attacks, malware injection attacks, side-channel attacks, and authentication & authorization attacks. Next, we explain each of these attacks.

### A. DDoS ATTACKS
DDoS attacks are a type of cyberattack in which an attacker attempts to distort normal services by flooding the internet traffic and making the service temporarily unavailable to the end users [84]. This type of attack is broadly classified into two types—flooding attack and logical attack. In a flooding attack, an attacker frequently sends malicious packets to edge devices or servers (victims) from (an electronic) source and makes the victims unable to handle these packets. As a result, the victims cannot respond to any legitimate requests on time [85], [86]. In a logical attack, the attacker sends malicious packets and misleads the application/protocol of the target machine by reflecting that all resources are fully occupied.

In comparison to cloud computing, PEC is more prone to such attacks as it provides services to edge devices that cannot maintain a strong defense system because of heterogeneous malware and computational limitations. Moreover, the DDoS attacker often intends to attack edge devices and then use them as a weapon against (edge) servers. In this aspect, a prominent example is Mirai botnet, which infected 65000 IoT devices and then exploited these devices; this DDoS attacker launched attacks against well-known services such as OVH, Dyn, and Krebs [87]. Bhardwaj *et al.* [88] demonstrated a proactive strategy to limit the impact of DDoS attacks by leveraging a proposed ShadowNet approach. The proposed approach consists of three components—edge

function, locally derived information, web service—and is unique in terms of the fast detection of the attack and defense responses. The authors presented that the proposed approach detects IoT DDoS attacks up to 10 times faster and enables reductions in the impacts of the damage by 82% of the internet traffic. Similarly, Zhou *et al.* [86] analyzed the DDoS attack mitigation in IIoT under the fog computing concept. The authors addressed the requirements of response time and the constraints related to the computational capabilities of devices in the IIoT network. A three-level architecture was proposed to mitigate the DDoS attack, which was then implemented in the Mero control system to yield effective results.

### B. MALWARE INJECTION ATTACKS
In malware injection attacks, the attacker aims to access the victim's service requests and to transfer malware into the network or computing systems [89]. Such attacks pose significant threats to data integrity and system security. In particular, (low level) edge servers and edge devices are more prone to such types of attacks. The injection of malware or malicious code at the edge server end is termed as server-side injection (SSI). SSI is classified into four types—SQL injection, extensible markup language (XML), server-side request forgery (SSRF), cross-site request forgery (CSRF), and cross-site scripting (XSS). The injection at the user side, in which an attacker injects malicious code into IoT devices, is termed as device-side injection (DSI). Examples of such attacks are remote code execution (RCE) and reaper [24], [90].

### C. SIDE-CHANNEL ATTACKS
In the side-channel attack, the attacker exploits publicly available data (not sensitive data) and correlates it with the user's private data "secretly" to infer confidential data. In this attack, an attacker continuously seizes the information from PEC infrastructure and uses it as an input to the ML/DL models or anonymous algorithms that produce the desired output (sensitive information). This type of adversarial attack can happen at any node of the network, and attackers exploit multiple techniques for side-channel attacks, for example, cache attack, timing attack, and electromagnetic attack [91]–[93]. The susceptibility of DL or ML-based systems and devices with edge intelligence to adversarial attacks is well known and has been studied intensely [94], [95]. In [96], the authors described a framework for edge learning as a service (EdgeLaaS) for healthcare infrastructures and emphasized the need for securing such data-sensitive systems against adversarial attacks. In [94], the authors exhaustively reviewed the fundamental methods to generate *adversarial examples*—intentionally designed inputs to ML models constituting an attack to force the model to make errors—and proposed a taxonomy for these methods. Additionally, the authors provided insights into adversarial attacks' applications to reinforcement learning, generative models, malware detection, etc. State-of-the-art approaches (e.g., input reconstruction, network verification, network distillation, etc.) to
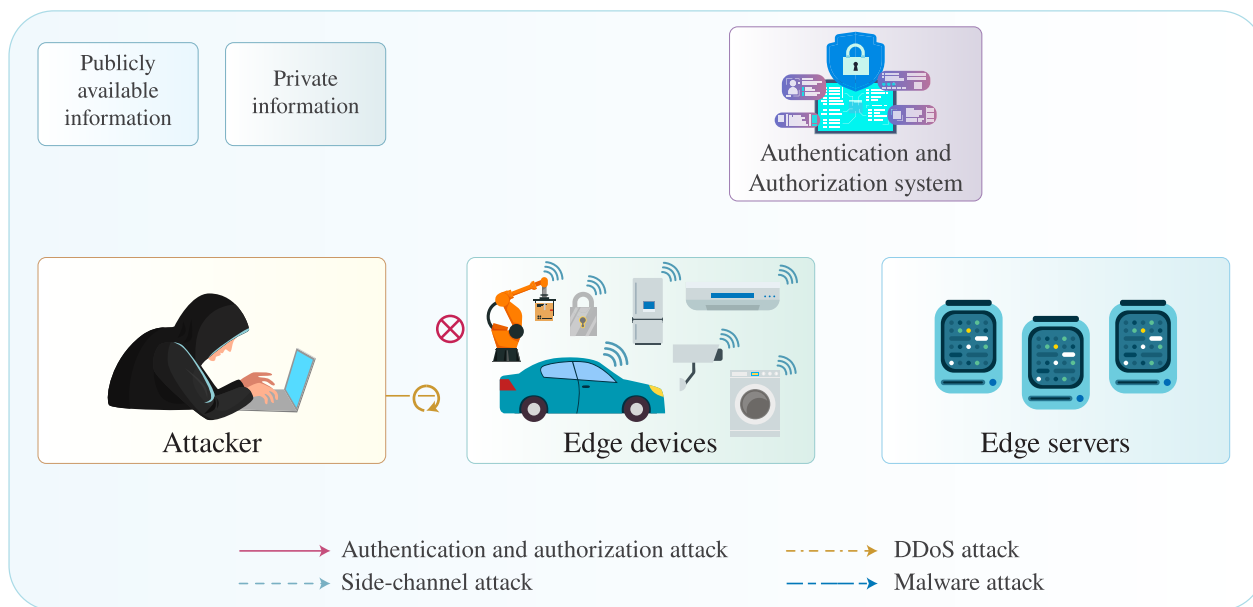
**FIGURE 8.** An overview of cyber security and cyberattacks in pervasive edge computing.

cope with such adversarial attacks are also discussed. In [95], Li *et al.* proposed a defensive framework (decentralized swift vigilance) in industrial systems to detect adversarial attacks swiftly and eliminate the risk failure without relying on complex reinforcement models. Concurrently, their work incorporated MEC and generative adversarial networks to achieve ultra low latency and privacy protections in industrial scenarios.

### D. AUTHENTICATION AND AUTHORIZATION ATTACKS
Authentication refers to the confirmation or verification of the identity of the entity who requests for certain services. And authorization refers to the process of ensuring the rights and access of an entity that should be within certain constraints and boundaries. In authentication and authorization attacks, the adversary aims to achieve access to the desired resources with fake credentials. Generally, PEC authentication is often accomplished among edge devices and servers, while authorization is the permission that edge servers grant to certain edge devices or its services/applications. The work in [24] classifies such attacks into four categories—dictionary attacks, authentication protocols attacks, authorization protocols attacks, and over-privileged attacks. Table 5 summarizes research works that have been cited in Section V, and Fig. 8 presents the aforementioned state of the art of cyber security attacks in PEC.

In addition to cyber security, there are still several open challenges that need to be solved to have a 5G-based PEC solution for IIoT, and these will be discussed in the next section.

## VI. FUTURE DIRECTIONS AND OPEN CHALLENGES
In addition to the IIoT-related challenges [98], [99], there are various other roadblocks to overcome for PEC to

become a mainstream and widely adopted solution. Below, we list these open research directions that merit further investigation.

### A. COMPUTATION PLACEMENT
While PEC offers numerous benefits discussed throughout this article, how to place the computation tasks on the available PEC nodes has to be determined considering the various tradeoffs discussed in Section II-A. Moreover, computation tasks such as ML processes might consist of multiple components with certain input-output dependencies. Hence, the computation tasks should be placed accordingly in the continuum of the pervasive edge and the cloud considering the tradeoffs, available resources, and the dependency among the computation tasks. Prior researches on edge service placement mostly consider smaller scale settings, and security aspects are largely overlooked, making such solutions ill-suited for IIoT settings where scalability and security are crucial. Moreover, the IIoT devices might consist of low-end sensors with limited uplink bandwidth and energy resources. Hence, the computation placement should take these peculiarities into account. For large-scale IoT solutions where the data produced by smart objects might be requested by many consumers, the efficiency of the communication between data producers and the consumers needs special attention, in particular, with respect to the existence of low-end data producers. In this respect, broker-based communication, also known as Pub/Sub architecture, facilitates transparency, mobility of data providers and consumers, scalability, and energy savings at the low-end data producers by decoupling the data producers and consumers. An emerging question is the placement of the brokers and computation tasks jointly.

## B. SPECTRUM MANAGEMENT

Data collection from the devices toward the PEC nodes results in increased uplink traffic in an industrial network. Moreover, given that some computations might be offloaded to the fog or cloud, there will be network traffic from the access network toward the core network and the Internet. As a result, the PEC policies will affect the amount of spectrum resources needed for the radio access network in the uplink and downlink as well as in the backhaul. Given that 5G envisions self-backhauling for small cells [100], there is a need for spectrum management schemes that dynamically allocate resources based on the needed capacity in each segment of an industrial network.

## C. EDGE INTELLIGENCE

Recent advances in communication systems have opened a path for FL in edge computing. However, several major challenges remain to be addressed for the edge intelligence to be offered ubiquitously. In particular, the deployment of FL at scale is not straightforward. Some of the key concerns are as follows. First, the learning framework needs to be robust against disruptions. That is, it should be able to handle cases where some nodes performing the computation lose connectivity temporarily or go completely offline due to, e.g., energy failures. Second, even if the nodes are reachable and perform their tasks, the results from each participating node need to be communicated efficiently, e.g., within certain delay bounds, to the rest of the nodes that rely on these nodes. Under dense deployments and low spectrum reuse factor, interference management plays a key role in ensuring the performance guarantees of FL. Third, due to the heterogeneity of the data produced by a variety of devices, the data that has to be processed and fed into learning schemes is mostly non-standardized. Finally, since some critical decisions about the operation of an industrial system for automation might rely on the outputs of the learning schemes, it is paramount that communication security is ensured and that learning schemes are robust against malicious or noisy inputs.

Significant research efforts are being made to deal with these problems across the world, and communication security, asynchronous FL techniques, improved ML algorithms, statistical analysis, and algorithms for communication reduction are some of the candidate solutions that will eventually accelerate the use of FL in edge computing.

## D. COMPUTATION OVER ENCRYPTED DATA

Even when the data to be processed is sensitive, there might be cases where a remote cloud is a favorable computation location due to the ample computation resources available for this sensitive data. To unlock the performance benefits offered by the cloud while preserving the data confidentiality, there is a need for performing computations over the encrypted input. There is a growing interest in privacy-preserving approaches such as solutions in [101] or [102]. However, to the best of our knowledge, it is still a widely unexplored research area.

## E. NETWORK MANAGEMENT

Due to the stringent performance and security requirements of IIoT networks, it is widely argued that industry plant owners would not prefer to outsource the control of their network to a third party, i.e., a 5G operator [98]. Private cellular networks aim to address this concern. However, there are many open questions such as the operation and deployment models, dynamic provisioning of the spectrum resources for the radio access network and the backhaul, as well as the deployment of computation units to meet the performance requirements of a particular vertical industry [103].

## VII. CONCLUSION

In this article, we have surveyed emerging technologies related to industrial internet-of-things (IIoT) enabled by 5G and beyond communication networks. We have discussed the main advantages of this paradigm—core network offloading (and benefits therefrom) and low latency for delay-sensitive applications (e.g., automatic control)—and reviewed the state-of-the-art in the PEC paradigm and its applications to the IIoT domain. We have also surveyed and described researches on distributed artificial intelligence methods, energy efficiency, and cyber security, which are three important research areas related to PEC. We identified the main open challenges that must be solved to have a scalable PEC-based IIoT network that operates efficiently under different conditions.

PEC deployments clearly have several interesting future directions, and academic and industrial researches into PEC are ongoing at a fast pace. This is motivated by the fact that PEC provides an extremely suitable and important deployment model for industrial communication networks, especially considering the recent trends of private industrial 5G networks incorporating local operations and flexible management. Nevertheless, PEC deployments also face many challenges that still need to be solved in order to have an effective solution that could be deployed in larger scales across different industrial domains, especially in terms of computational performance, energy efficiency, and cyber security.

## REFERENCES

[1] H. Lasi, P. Fettke, H. G. Kemper, T. Feld, and M. Hoffmann, "Industry 4.0," *Bus. Inf. Syst. Eng.*, vol. 6, no. 4, pp. 239–242, 2014.

[2] P. Corcoran and S. K. Datta, "Mobile-edge computing and the Internet of Things for consumers: Extending cloud computing and services to the edge of the network," *IEEE Consum. Electron. Mag.*, vol. 5, no. 4, pp. 73–74, Oct. 2016.

[3] P. Daugherty, P. Banerjee, W. Negm, and A. E. Alter. (2014). *Driving Unconventional Growth through the Industrial Internet of Things*. Accenture. p. 20. Accessed: Aug. 16, 2016. [Online]. Available: https://www.accenture.com/au-en/_acnmedia/Accenture/next-gen/reassembling-industry/pdf/Accenture-Driving-Unconventional-Growth-through-IIoT.pdf

[4] P. Popovski, K. F. Trillingsgaard, O. Simeone, and G. Durisi, "5G wireless network slicing for eMBB, URLLC, and mMTC: A communication-theoretic view," *IEEE Access*, vol. 6, pp. 55765–55779, 2018.

[5] S. Mumtaz, A. Bo, A. Al-Dulaimi, and K.-F. Tsang, "Guest editorial 5G and beyond mobile technologies and applications for industrial IoT (IIoT)," *IEEE Trans. Ind. Informat.*, vol. 14, no. 6, pp. 2588–2591, Jun. 2018.

[6] Y. Siriwardhana, P. Porambage, M. Liyanage, J. S. Walia, M. Matinmikko-Blue, and M. Ylianttila, "Micro-operator driven local 5G network architecture for industrial Internet," in *Proc. IEEE Wireless Commun. Netw. Conf. (WCNC)*, Apr. 2019, pp. 1–8.

[7] N. Fernando, S. W. Loke, and W. Rahayu, "Mobile cloud computing: A survey," *Future Gener. Comput. Syst.*, vol. 29, no. 1, pp. 84–106, Jan. 2013, doi: 10.1016/j.future.2012.05.023.

[8] K. Kaur, S. Garg, G. S. Aujla, N. Kumar, J. J. P. C. Rodrigues, and M. Guizani, "Edge computing in the industrial Internet of Things environment: Software-defined-networks-based edge-cloud interplay," *IEEE Commun. Mag.*, vol. 56, no. 2, pp. 44–51, Feb. 2018.

[9] Q.-V. Pham, F. Fang, V. N. Ha, M. J. Piran, M. Le, L. B. Le, W.-J. Hwang, and Z. Ding, "A survey of multi-access edge computing in 5G and beyond: Fundamentals, technology integration, and state-of-the-art," 2019, *arXiv:1906.08452*. [Online]. Available: http://arxiv.org/abs/1906.08452

[10] N. Hassan, K.-L.-A. Yau, and C. Wu, "Edge computing in 5G: A review," *IEEE Access*, vol. 7, pp. 127276–127289, 2019.

[11] Y. Miao, G. Wu, M. Li, A. Ghoneim, M. Al-Rakhami, and M. S. Hossain, "Intelligent task prediction and computation offloading based on mobile-edge cloud computing," *Future Gener. Comput. Syst.*, vol. 102, pp. 925–931, Jan. 2020, doi: 10.1016/j.future.2019.09.035.

[12] W. Shi, J. Cao, Q. Zhang, Y. Li, and L. Xu, "Edge computing: Vision and challenges," *IEEE Internet Things J.*, vol. 3, no. 5, pp. 637–646, Oct. 2016.

[13] P. Popovski, O. Simeone, F. Boccardi, D. Gunduz, and O. Sahin, "Semantic-effectiveness filtering and control for post-5G wireless connectivity," 2019, *arXiv:1907.02441*. [Online]. Available: http://arxiv.org/abs/1907.02441

[14] D. Gutierrez-Rojas, M. Ullah, I. T. Christou, G. Almeida, P. H. J. Nardelli, D. Carrillo, J. M. Sant'Ana, H. Alves, M. Dzaferagic, A. Chiumento, and C. Kalalas, "Three-layer approach to detect anomalies in industrial environments based on machine learning," 2020, *arXiv:2004.09097*. [Online]. Available: http://arxiv.org/abs/2004.09097

[15] X. Hu, L. Wang, K.-K. Wong, M. Tao, Y. Zhang, and Z. Zheng, "Edge and central cloud computing: A perfect pairing for high energy efficiency and low-latency," *IEEE Trans. Wireless Commun.*, vol. 19, no. 2, pp. 1070–1083, Feb. 2020.

[16] J. Ren, G. Yu, Y. He, and G. Y. Li, "Collaborative cloud and edge computing for latency minimization," *IEEE Trans. Veh. Technol.*, vol. 68, no. 5, pp. 5031–5044, May 2019.

[17] U. Saleem, Y. Liu, S. Jangsher, X. Tao, and Y. Li, "Latency minimization for D2D-enabled partial computation offloading in mobile edge computing," *IEEE Trans. Veh. Technol.*, vol. 69, no. 4, pp. 4472–4486, Apr. 2020.

[18] Y. Ai, M. Peng, and K. Zhang, "Edge computing technologies for Internet of Things: A primer," *Digit. Commun. Netw.*, vol. 4, no. 2, pp. 77–86, Apr. 2018.

[19] K. Ha, Z. Chen, W. Hu, W. Richter, P. Pillai, and M. Satyanarayanan, "Towards wearable cognitive assistance," in *Proc. 12th Annu. Int. Conf. Mobile Syst., Appl., Services (MobiSys)*, 2014, pp. 68–81.

[20] Z. Deng, Z. Cai, and M. Liang, "A multi-hop VANETs-assisted offloading strategy in vehicular mobile edge computing," *IEEE Access*, vol. 8, pp. 53062–53071, 2020.

[21] T. Dang and M. Peng, "Joint radio communication, caching, and computing design for mobile virtual reality delivery in fog radio access networks," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 7, pp. 1594–1607, Jul. 2019.

[22] J. Kwak, Y. Kim, L. B. Le, and S. Chong, "Hybrid content caching in 5G wireless networks: Cloud versus edge caching," *IEEE Trans. Wireless Commun.*, vol. 17, no. 5, pp. 3030–3045, May 2018.

[23] J. Ren, Y. Ruan, and G. Yu, "Data transmission in mobile edge networks: Whether and where to compress?" *IEEE Commun. Lett.*, vol. 23, no. 3, pp. 490–493, Mar. 2019.

[24] Y. Xiao, Y. Jia, C. Liu, X. Cheng, J. Yu, and W. Lv, "Edge computing security: State of the art and challenges," *Proc. IEEE*, vol. 107, no. 8, pp. 1608–1631, Aug. 2019.

[25] J. Zhang, B. Chen, Y. Zhao, X. Cheng, and F. Hu, "Data security and privacy-preserving in edge computing paradigm: Survey and open issues," *IEEE Access*, vol. 6, pp. 18209–18237, 2018.

[26] M. Yahuza, M. Y. I. B. Idris, A. W. B. A. Wahab, A. T. S. Ho, S. Khan, S. N. B. Musa, and A. Z. B. Taha, "Systematic review on security and privacy requirements in edge computing: State of the art and future research opportunities," *IEEE Access*, vol. 8, pp. 76541–76567, 2020.

[27] D. Liu, Z. Yan, W. Ding, and M. Atiquzzaman, "A survey on secure data analytics in edge computing," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4946–4967, Jun. 2019.

[28] T. Wang, G. Zhang, A. Liu, M. Z. A. Bhuiyan, and Q. Jin, "A secure IoT service architecture with an efficient balance dynamics based on cloud and edge computing," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4831–4843, Jun. 2019.

[29] R. Ullah, M. A. U. Rehman, and B.-S. Kim, "Design and implementation of an open source framework and prototype for named data networking-based edge cloud computing system," *IEEE Access*, vol. 7, pp. 57741–57759, 2019.

[30] R. W. Lucky, "Automatic equalization for digital communication," *Bell Syst. Tech. J.*, vol. 44, no. 4, pp. 547–588, Apr. 1965.

[31] P. Mach and Z. Becvar, "Mobile edge computing: A survey on architecture and computation offloading," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 3, pp. 1628–1656, 3rd Quart., 2017.

[32] N. Abbas, Y. Zhang, A. Taherkordi, and T. Skeie, "Mobile edge computing: A survey," *IEEE Internet Things J.*, vol. 5, no. 1, pp. 450–465, Feb. 2018.

[33] C. Esposito, A. Castiglione, F. Pop, and K.-K.-R. Choo, "Challenges of connecting edge and cloud computing: A security and forensic perspective," *IEEE Cloud Comput.*, vol. 4, no. 2, pp. 13–17, Mar. 2017.

[34] J. Zhang, H.-W. Lee, and E. Modiano, "On the robustness of distributed computing networks," in *Proc. 15th Int. Conf. Design Reliable Commun. Netw. (DRCN)*, Mar. 2019, pp. 122–129.

[35] Y. Sharma, B. Javadi, W. Si, and D. Sun, "Reliability and energy efficiency in cloud computing systems: Survey and taxonomy," *J. Netw. Comput. Appl.*, vol. 74, pp. 66–85, Oct. 2016.

[36] P. Gill, N. Jain, and N. Nagappan, "Understanding network failures in data centers: Measurement, analysis, and implications," *ACM SIGCOMM Comput. Commun. Rev.*, vol. 41, no. 4, pp. 350–361, Aug. 2011.

[37] H. Yang, W. Bai, A. Yu, Q. Yao, J. Zhang, Y. Lin, and Y. Lee, "Bandwidth compression protection against collapse in fog-based wireless and optical networks," *IEEE Access*, vol. 6, pp. 54760–54768, 2018.

[38] R. Makhlouf, "Cloudy transaction costs: A dive into cloud computing economics," *J. Cloud Comput.*, vol. 9, no. 1, p. 1, Dec. 2020.

[39] H. Yuan and M. Zhou, "Profit-maximized collaborative computation offloading and resource allocation in distributed cloud and edge computing systems," *IEEE Trans. Autom. Sci. Eng.*, early access, Jul. 14, 2020, doi: 10.1109/TASE.2020.3000946.

[40] S. Wang, X. Zhang, Z. Yan, and W. Wenbo, "Cooperative edge computing with sleep control under nonuniform traffic in mobile edge networks," *IEEE Internet Things J.*, vol. 6, no. 3, pp. 4295–4306, Jun. 2019.

[41] Z. Tao, Q. Xia, Z. Hao, C. Li, L. Ma, S. Yi, and Q. Li, "A survey of virtual machine management in edge computing," *Proc. IEEE*, vol. 107, no. 8, pp. 1482–1499, Aug. 2019.

[42] P. Bellavista, A. Corradi, L. Foschini, and D. Scotece, "Differentiated Service/Data migration for edge services leveraging container characteristics," *IEEE Access*, vol. 7, pp. 139746–139758, 2019.

[43] L. Lin, X. Liao, H. Jin, and P. Li, "Computation offloading toward edge computing," *Proc. IEEE*, vol. 107, no. 8, pp. 1584–1607, Aug. 2019.

[44] Z. Rejiba, X. Masip-Bruin, and E. Marín-Tordera, "A survey on mobility-induced service migration in the fog, edge, and related computing paradigms," *ACM Comput. Surveys*, vol. 52, no. 5, pp. 1–33, Oct. 2019. [Online]. Available: https://doi-org.ezproxy2.utwente.nl/10.1145/3326540

[45] L. Ma, S. Yi, and Q. Li, "Efficient service handoff across edge servers via docker container migration," in *Proc. 2nd ACM/IEEE Symp. Edge Comput. (SEC)*. New York, NY, USA: Association for Computing Machinery, 2017, pp. 1–13. [Online]. Available: https://doi-org.ezproxy2.utwente.nl/10.1145/3132211.3134460

[46] L. Ma, S. Yi, N. Carter, and Q. Li, "Efficient live migration of edge services leveraging container layered storage," *IEEE Trans. Mobile Comput.*, vol. 18, no. 9, pp. 2020–2033, Sep. 2019.

[47] S. Wang, J. Xu, N. Zhang, and Y. Liu, "A survey on service migration in mobile edge computing," *IEEE Access*, vol. 6, pp. 23511–23528, 2018.

[48] J. Ding and J. Choi, "Comparison of preamble structures for grant-free random access in massive MIMO systems," *IEEE Wireless Commun. Lett.*, vol. 9, no. 2, pp. 166–170, Feb. 2020.

[49] R. Abbas, M. Shirvanimoghaddam, Y. Li, and B. Vucetic, "A novel analytical framework for massive grant-free NOMA," *IEEE Trans. Commun.*, vol. 67, no. 3, pp. 2436–2449, Mar. 2019.

[50] Y. Zhang, L. Jiao, J. Yan, and X. Lin, "Dynamic service placement for virtual reality group gaming on mobile edge cloudlets," *IEEE J. Sel. Areas Commun.*, vol. 37, no. 8, pp. 1881–1897, Aug. 2019.

[51] Y. Sun, Z. Chen, M. Tao, and H. Liu, "Communications, caching, and computing for mobile virtual reality: Modeling and tradeoff," *IEEE Trans. Commun.*, vol. 67, no. 11, pp. 7573–7586, Nov. 2019.

[52] C. Quadri, M. Premoli, A. Ceselli, S. Gaito, and G. P. Rossi, "Optimal assignment plan in sliced backhaul networks," *IEEE Access*, vol. 8, pp. 68983–69002, 2020.

[53] X. Shen, J. Gao, W. Wu, K. Lyu, M. Li, W. Zhuang, X. Li, and J. Rao, "AI-assisted network-slicing based next-generation wireless networks," *IEEE Open J. Veh. Technol.*, vol. 1, pp. 45–66, 2020.

[54] S. Mukherjee and J. Lee, "Edge computing-enabled cell-free massive MIMO systems," *IEEE Trans. Wireless Commun.*, vol. 19, no. 4, pp. 2884–2899, Apr. 2020.

[55] P. Ahokangas, M. Matinmikko-Blue, S. Yrjola, V. Seppanen, H. Hammainen, R. Jurva, and M. Latva-aho, "Business models for local 5G micro operators," *IEEE Trans. Cognit. Commun. Netw.*, vol. 5, no. 3, pp. 730–740, Sep. 2019.

[56] J. Zhu, M. Zhao, S. Zhang, and W. Zhou, "Exploring the road to 6G: ABC—Foundation for intelligent mobile networks," *China Commun.*, vol. 17, no. 6, pp. 51–67, 2020.

[57] S. Liu, L. Liu, J. Tang, B. Yu, Y. Wang, and W. Shi, "Edge computing for autonomous driving: Opportunities and challenges," *Proc. IEEE*, vol. 107, no. 8, pp. 1697–1716, Aug. 2019.

[58] H. Ma, S. Li, E. Zhang, Z. Lv, J. Hu, and X. Wei, "Cooperative autonomous driving oriented MEC-aided 5G-V2X: Prototype system design, field tests and AI-based optimization tools," *IEEE Access*, vol. 8, pp. 54288–54302, 2020.

[59] N. H. Mahmood, H. Alves, O. A. Lopez, M. Shehab, D. P. M. Osorio, and M. Latva-Aho, "Six key features of machine type communication in 6G," in *Proc. 2nd 6G Wireless Summit (6G SUMMIT)*, Mar. 2020, pp. 1–5.

[60] 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," *IEEE Commun. Surveys Tuts.*, vol. 22, no. 3, pp. 2031–2063, 3rd Quart., 2020.

[61] Z. M. Fadlullah and N. Kato, "HCP: Heterogeneous computing platform for federated learning based collaborative content caching towards 6G networks," *IEEE Trans. Emerg. Topics Comput.*, early access, Apr. 8, 2020, doi: 10.1109/TETC.2020.2986238.

[62] Y. Lu, X. Huang, Y. Dai, S. Maharjan, and Y. Zhang, "Differentially private asynchronous federated learning for mobile edge computing in urban informatics," *IEEE Trans. Ind. Informat.*, vol. 16, no. 3, pp. 2134–2143, Mar. 2020.

[63] J. Ren, H. Wang, T. Hou, S. Zheng, and C. Tang, "Federated learning-based computation offloading optimization in edge computing-supported Internet of Things," *IEEE Access*, vol. 7, pp. 69194–69201, 2019.

[64] D. Ye, R. Yu, M. Pan, and Z. Han, "Federated learning in vehicular edge computing: A selective model aggregation approach," *IEEE Access*, vol. 8, pp. 23920–23935, 2020.

[65] 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 J. Sel. Areas Commun.*, vol. 37, no. 6, pp. 1205–1221, Jun. 2019.

[66] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," 2016, *arXiv:1610.05492*. [Online]. Available: http://arxiv.org/abs/1610.05492

[67] J. Chen and X. Ran, "Deep learning with edge computing: A review," *Proc. IEEE*, vol. 107, no. 8, pp. 1655–1674, Aug. 2019.

[68] X. Wang, Y. Han, V. C. M. Leung, D. Niyato, X. Yan, and X. Chen, "Convergence of edge computing and deep learning: A comprehensive survey," *IEEE Commun. Surveys Tuts.*, vol. 22, no. 2, pp. 869–904, 2nd Quart., 2020.

[69] H. Li, K. Ota, and M. Dong, "Learning IoT in edge: Deep learning for the Internet of Things with edge computing," *IEEE Netw.*, vol. 32, no. 1, pp. 96–101, Jan. 2018.

[70] F. Liang, W. Yu, X. Liu, D. Griffith, and N. Golmie, "Toward edge-based deep learning in industrial Internet of Things," *IEEE Internet Things J.*, vol. 7, no. 5, pp. 4329–4341, May 2020.

[71] M. G. Sarwar Murshed, E. Verenich, J. J. Carroll, N. Khan, and F. Hussain, "Hazard detection in supermarkets using deep learning on the edge," 2020, *arXiv:2003.04116*. [Online]. Available: http://arxiv.org/abs/2003.04116

[72] C.-F. Lai, W.-C. Chien, L. T. Yang, and W. Qiang, "LSTM and edge computing for big data feature recognition of industrial electrical equipment," *IEEE Trans. Ind. Informat.*, vol. 15, no. 4, pp. 2469–2477, Apr. 2019.

[73] S. Tuli, N. Basumatary, S. S. Gill, M. Kahani, R. C. Arya, G. S. Wander, and R. Buyya, "HealthFog: An ensemble deep learning based smart healthcare system for automatic diagnosis of heart diseases in integrated IoT and fog computing environments," *Future Gener. Comput. Syst.*, vol. 104, pp. 187–200, Mar. 2020.

[74] Y. Wang, M. Liu, P. Zheng, H. Yang, and J. Zou, "A smart surface inspection system using faster R-CNN in cloud-edge computing environment," *Adv. Eng. Informat.*, vol. 43, Jan. 2020, Art. no. 101037.

[75] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Process. Mag.*, vol. 37, no. 3, pp. 50–60, May 2020.

[76] N. H. Tran, W. Bao, A. Zomaya, M. N. H. Nguyen, and C. S. Hong, "Federated learning over wireless networks: Optimization model design and analysis," in *Proc. IEEE Conf. Comput. Commun. (INFOCOM)*, Apr. 2019, pp. 1387–1395.

[77] Q. Yang, Y. Liu, T. Chen, and Y. Tong, "Federated machine learning: Concept and applications," *ACM Trans. Intell. Syst. Technol. (TIST)*, vol. 10, no. 2, pp. 1–19, 2019.

[78] S. Chen, Y. Zheng, W. Lu, V. Varadarajan, and K. Wang, "Energy-optimal dynamic computation offloading for industrial IoT in fog computing," *IEEE Trans. Green Commun. Netw.*, vol. 4, no. 2, pp. 566–576, Jun. 2020.

[79] Z. Hong, W. Chen, H. Huang, S. Guo, and Z. Zheng, "Multi-hop cooperative computation offloading for industrial IoT–edge–cloud computing environments," *IEEE Trans. Parallel Distrib. Syst.*, vol. 30, no. 12, pp. 2759–2774, Dec. 2019.

[80] M. Guo, L. Li, and Q. Guan, "Energy-efficient and delay-guaranteed workload allocation in IoT-edge-cloud computing systems," *IEEE Access*, vol. 7, pp. 78685–78697, 2019.

[81] P. Liu, G. Xu, K. Yang, K. Wang, and X. Meng, "Jointly optimized energy-minimal resource allocation in cache-enhanced mobile edge computing systems," *IEEE Access*, vol. 7, pp. 3336–3347, 2019.

[82] Y. Wu, B. Shi, L. P. Qian, F. Hou, J. Cai, and X. S. Shen, "Energy-efficient multi-task multi-access computation offloading via NOMA transmission for IoTs," *IEEE Trans. Ind. Informat.*, vol. 16, no. 7, pp. 4811–4822, Jul. 2020.

[83] Y. Wang, K. Wang, H. Huang, T. Miyazaki, and S. Guo, "Traffic and computation co-offloading with reinforcement learning in fog computing for industrial applications," *IEEE Trans. Ind. Informat.*, vol. 15, no. 2, pp. 976–986, Feb. 2019.

[84] A. Srivastava, B. Gupta, A. Tyagi, A. Sharma, and A. Mishra, "A recent survey on ddos attacks and defense mechanisms," in *Proc. Int. Conf. Parallel Distrib. Comput. Technol. Appl.* Berlin, Germany: Springer, 2011, pp. 570–580.

[85] K. Bhardwaj, J. C. Miranda, and A. Gavrilovska, "Towards iot-ddos prevention using edge computing," in *Proc. USENIX Workshop Hot Topics Edge Comput. (HotEdge)*. Boston, MA, USA: USENIX Association, Jul. 2018, pp. 1–7. [Online]. Available: https://www.usenix.org/conference/hotedge18/presentation/bhardwaj

[86] L. Zhou, H. Guo, and G. Deng, "A fog computing based approach to DDoS mitigation in IIoT systems," *Comput. Secur.*, vol. 85, pp. 51–62, Aug. 2019.

[87] M. Antonakakis, T. April, M. Bailey, M. Bernhard, E. Bursztein, J. Cochran, Z. Durumeric, J. A. Halderman, L. Invernizzi, M. Kallitsis, and D. Kumar, "Understanding the mirai botnet," in *Proc. 26th USENIX Secur. Symp. (USENIX Security)*, 2017, pp. 1093–1110.

[88] K. Bhardwaj, J. C. Miranda, and A. Gavrilovska, "Towards IoT-ddos prevention using edge computing," in *Proc. USENIX Workshop Hot Topics Edge Comput. (HotEdge)*, 2018, pp. 1–7.

[89] M. R. Watson, N.-U.-H. Shirazi, A. K. Marnerides, A. Mauthe, and D. Hutchison, "Malware detection in cloud computing infrastructures," *IEEE Trans. Dependable Secure Comput.*, vol. 13, no. 2, pp. 192–205, Mar. 2016.

[90] C. Barron, H. Yu, and J. Zhan, "Cloud computing security case studies and research," in *Proc. World Congr. Eng.*, 2013, vol. 2, no. 2, pp. 1–6.

[91] Y. Xu, W. Cui, and M. Peinado, "Controlled-channel attacks: Deterministic side channels for untrusted operating systems," in *Proc. IEEE Symp. Secur. Privacy*, May 2015, pp. 640–656.

[92] T. Zhang, Y. Zhang, and R. B. Lee, "Cloudradar: A real-time side-channel attack detection system in clouds," in *Proc. Int. Symp. Res. Attacks, Intrusions, Defenses*. Berlin, Germany: Springer, 2016, pp. 118–140.

[93] Y. Lyu and P. Mishra, "A survey of side-channel attacks on caches and countermeasures," *J. Hardw. Syst. Secur.*, vol. 2, no. 1, pp. 33–50, Mar. 2018.

[94] X. Yuan, P. He, Q. Zhu, and X. Li, "Adversarial examples: Attacks and defenses for deep learning," *IEEE Trans. Neural Netw. Learn. Syst.*, vol. 30, no. 9, pp. 2805–2824, Sep. 2019.

[95] G. Li, K. Ota, M. Dong, J. Wu, and J. Li, "DeSVig: Decentralized swift vigilance against adversarial attacks in industrial artificial intelligence systems," *IEEE Trans. Ind. Informat.*, vol. 16, no. 5, pp. 3267–3277, May 2020.

[96] G. Li, G. Xu, A. K. Sangaiah, J. Wu, and J. Li, "EdgeLaaS: Edge learning as a service for knowledge-centric connected healthcare," *IEEE Netw.*, vol. 33, no. 6, pp. 37–43, Nov. 2019.

[97] Z. Zhou, J. Feng, Z. Chang, and X. Shen, "Energy-efficient edge computing service provisioning for vehicular networks: A consensus ADMM approach," *IEEE Trans. Veh. Technol.*, vol. 68, no. 5, pp. 5087–5099, May 2019.

[98] S. Vitturi, C. Zunino, and T. Sauter, "Industrial communication systems and their future challenges: Next-generation Ethernet, IIoT, and 5G," *Proc. IEEE*, vol. 107, no. 6, pp. 944–961, Jun. 2019.

[99] E. Sisinni, A. Saifullah, S. Han, U. Jennehag, and M. Gidlund, "Industrial Internet of Things: Challenges, opportunities, and directions," *IEEE Trans. Ind. Informat.*, vol. 14, no. 11, pp. 4724–4734, Nov. 2018.

[100] M. Polese, M. Giordani, T. Zugno, A. Roy, S. Goyal, D. Castor, and M. Zorzi, "Integrated access and backhaul in 5G mmWave networks: Potential and challenges," *IEEE Commun. Mag.*, vol. 58, no. 3, pp. 62–68, Mar. 2020.

[101] K.-K.-R. Choo, S. Gritzalis, and J. H. Park, "Cryptographic solutions for industrial Internet-of-things: Research challenges and opportunities," *IEEE Trans. Ind. Informat.*, vol. 14, no. 8, pp. 3567–3569, Aug. 2018.

[102] J. Feng, L. T. Yang, Q. Zhu, and K.-K.-R. Choo, "Privacy-preserving tensor decomposition over encrypted data in a federated cloud environment," *IEEE Trans. Dependable Secure Comput.*, vol. 17, no. 4, pp. 857–868, Jul. 2020.

[103] A. Rostami, "Private 5G networks for vertical industries: Deployment and operation models," in *Proc. IEEE 2nd 5G World Forum (5GWF)*, Sep. 2019, pp. 433–439.

**ARTHUR SOUSA DE SENA** (Student Member, IEEE) received the B.Sc. degree in computer engineering and the M.Sc. degree in teleinformatics engineering from the Federal University of Ceará, Brazil, in 2017 and 2019, respectively. He is currently pursuing the Ph.D. degree with the School of Energy Systems, LUT University, Finland. From 2014 to 2015, he studied computer engineering as an Exchange Student at the Illinois Institute of Technology, USA. He is also a Researcher with the Cyber-Physical Systems Group, LUT. His research interests include signal processing, mobile communications systems, non-orthogonal multiple access techniques, intelligent metasurfaces, and massive MIMO.
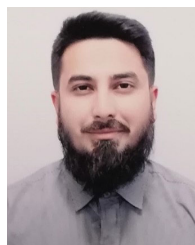
**DANIEL GUTIERREZ-ROJAS** (Graduate Student Member, IEEE) received the B.Sc. degree in electrical engineering from the University of Antioquia, Colombia, in 2016, and the M.Sc. degree in protection of power systems from the University of São Paulo, Brazil, in 2017. He is currently pursuing the Ph.D. degree with the School of Energy Systems, LUT University, Finland. From 2017 to 2019, he has worked as a Security of Operation and Fault analyst for Colombia's National electrical operator. His research interests include predictive maintenance, power systems, microgrids, mobile communication systems, and electrical protection systems.

**DICK CARRILLO MELGAREJO** (Graduate Student Member, IEEE) received the B.Eng. degree (Hons.) in electronics and electrical engineering from San Marcos National University, Lima, Perú, in 2004, and the M.Sc. degree in electrical engineering from the Pontifical Catholic University of Rio de Janeiro, Rio de Janeiro, Brazil, in 2008. He is currently pursuing the Ph.D. degree in electrical engineering with the Lappeenranta–Lahti University of Technology. From 2008 to 2010, he contributed to WIMAX (IEEE 802.16m) standardization. From 2010 to 2018, he has worked with the design and implementation of cognitive radio networks and projects based on 3GPP technologies. Since 2018, he has been a Researcher with the Lappeenranta–Lahti University of Technology. His research interests include mobile technologies beyond 5G, energy harvesting, intelligent meta-surfaces, and cell-free mMIMO.

**ARUN NARAYANAN** (Member, IEEE) received the B.E. degree in electrical engineering from the Visvesvaraya National Institute of Technology, Nagpur, India, in 2002, the M.Sc. degree in energy technology from the Lappeenranta University of Technology (LUT), Lappeenranta, Finland, in 2013, and the Ph.D. degree from the School of Energy Systems, LUT University, in 2019.

He is currently a Postdoctoral Researcher with the Research Group "Cyber-Physical Systems Group," LUT University. His research interests include renewable energy-based smart microgrids, electricity distribution and markets, demand-side management, energy management systems, and information and communications technology. He focuses on applying optimization, computational concepts, and artificial intelligence techniques to renewable electrical energy problems.

**HAFIZ MAJID HUSSAIN** (Member, IEEE) received the B.S. degree in electrical engineering from the National University of Computer and Emerging Sciences, in 2014, and the M.S. degree in electrical engineering from the University of Engineering and Technology Taxila, Pakistan, in 2017. He is currently pursuing the Ph.D. degree in electrical engineering with the Lappeenranta University of Technology, Finland. He is also the part of the research group "Cyber-Physical Systems Group," and a project EnergyNet that focuses on building the energy internet as a large-scale IoT-based cyber-physical systems. His research interests include demand response applications, energy resource optimization in the smart grid, and information security technologies.

**MEHAR ULLAH** (Graduate Student Member, IEEE) received the B.S. degree in information technology from Iqra National University, Pakistan, and the master's degree in software engineering from the Lappeenranta-Lahti University of Technology (LUT), Finland, where he is currently pursuing the Ph.D. degree. His main research interests include the IoT and cyber-physical systems especially for industrial applications.

**SUZAN BAYHAN** received the Ph.D. degree in computer engineering from Bogazici University, in 2012. She has worked at the University of Helsinki, Aalto University, and TU Berlin, from 2012 to 2019. She is currently an Assistant Professor with the University of Twente and an Adjunct Professor (Docent) with the University of Helsinki. Her current research interests include spectrum sharing and coexistence of wireless networks, WiFi and LTE radio resource management, and edge computing. She received the Best Paper Award at ACM ICN 2015 and the Best Demo Award at IEEE INFOCOM 2020.

**PEDRO H. J. NARDELLI** (Senior Member, IEEE) received the B.S. and M.Sc. degrees in electrical engineering from the State University of Campinas, Brazil, in 2006 and 2008, respectively, and the Ph.D. degree from the University of Oulu, Finland and the State University of Campinas, following a dual degree agreement, in 2013. He is currently an Assistant Professor (tenure track) in IoT in energy systems with LUT University, Finland, and holds a position of an Academy of Finland Research Fellow with a project called Building the Energy Internet as a Large-Scale IoT-Based Cyber-Physical System that manages the energy inventory of distribution grids as discretized packets via machine-type communications (EnergyNet). He leads the Cyber-Physical Systems Group, LUT, where he is also the Project Coordinator of the CHIST-ERA European consortium Framework for the Identification of Rare Events via Machine Learning and IoT Networks (FIREMAN). He is also an Adjunct Professor with the University of Oulu in the topic of "communications strategies and information processing in energy systems." His research interest includes wireless communications particularly applied in industrial automation and energy systems. He received the Best Paper Award of IEEE PES Innovative Smart Grid Technologies Latin America 2019 in the track "Big Data and Internet of Things." More information: https://sites.google.com/view/nardelli/

• • •