

A Hybrid Computing Solution and Resource Scheduling Strategy for Edge Computing in Smart Manufacturing

Xiaomin Li, Jiafu Wan , *Member, IEEE*, Hong-Ning Dai , *Senior Member, IEEE*, Muhammad Imran , Min Xia , and Antonio Celesti 

Abstract—At present, smart manufacturing computing framework has faced many challenges such as the lack of an effective framework of fusing computing historical heritages and resource scheduling strategy to guarantee the low-latency requirement. In this paper, we propose a hybrid computing framework and design an intelligent resource scheduling strategy to fulfill the real-time requirement in smart manufacturing with edge computing support. First, a four-layer computing system in a smart manufacturing environment is provided to support the artificial intelligence task operation with the network perspective. Then, a two-phase algorithm for scheduling the computing resources in the edge layer is designed based on greedy and threshold strategies with latency constraints. Finally, a prototype platform was developed. We conducted experiments on the prototype to evaluate the performance of the proposed framework with a comparison of the traditionally-used methods. The proposed strategies have demonstrated the excellent real-time, satisfaction degree (SD), and energy consump-

tion performance of computing services in smart manufacturing with edge computing.

Index Terms—Edge computing, industry 4.0, resource scheduling, smart manufacturing.

I. INTRODUCTION

RECENTLY, thanks for the great progress of the information and communication technologies in manufacturing domains, Industrial Internet of Things (IIoT), cyber-physical system, and other smart frameworks and systems have been constructed and implemented to increase the flexibility and enhance economic efficiency [1]–[5]. In this scenario, an increasing attention from enterprises and academia has been devoted to inventing and extending new computing technology or framework, e.g., Industry 4.0, smart factory and intelligent manufacturing [6]–[10]. However, traditional manufacturing systems that are lacking of efficiency, in terms of a slow computing speed on complicated task, are not suitable for complicated manufacturing process and dealing with big data analysis, especially for artificial intelligence (AI) task. Hence, an extended framework on top of traditional manufacturing systems by introducing computing resources with different levels of computing capabilities, such as cloud and edge computing platforms, can meet fundamental requirement for resolving the aforementioned problems. However, the cloud computing platforms generally are far from the industrial devices; it consequently increases the latency, leads to the lag in data transmission, and cannot guarantee the real-time performance.

Real-time feature of data flow of the whole industrial system, directly affecting the production efficiency and normal operation of the system, plays a critical role in a smart factory and Industry 4.0 [11]. Therefore, reducing the overhead of data processing and transmission in computing-extensive tasks (e.g., AI and deep learning tasks) is another aspect to guarantee the real-time performance. The current latency-constraints methods concentrate on the network optimizations, data fusion, computing-task simplification, and computing resource decentralization. Edge/fog computing frameworks, close to producing equipment consequently leading to the decrement of the latency for data communications between servers and machines, are good candidate strategies for smart manufacturing to tackle the above problems [12]. In pioneering works [13]–[16], the authors

Manuscript received December 1, 2018; revised January 23, 2019; accepted February 11, 2019. Date of publication February 18, 2019; date of current version July 3, 2019. This work was supported in part by the National Key Research and Development Project of China under Grant 2017YFE0101000, in part by the Joint Fund of the National Natural Science Foundation of China and Guangdong Province under Grant U1801264, in part by the Key Program of Natural Science Foundation of Guangdong Province under Grant 2017B030311008, and in part by the Science and Technology Program of Guangzhou, China, under Grant 201802030005. The work of M. Imran was supported by the Deanship of Scientific Research, King Saud University through research group number RG-1435-051. Paper no. TII-18-3244.R1. (*Corresponding author: Jiafu Wan.*)

X. Li is with the School of Mechanical Engineering, Zhongkai University of Agriculture and Engineering, Guangzhou 510225, China (e-mail: lixiaomin@zhku.edu.cn).

J. Wan is with the School of Mechanical and Automotive Engineering, South China University of Technology, Guangzhou 510641, China (e-mail: mejwan@scut.edu.cn).

H.-N. Dai is with the College of Information Technology, Macau University of Science and Technology, Macau, China (e-mail: hndai@iee.org).

M. Imran is with the College of Computer and Information Sciences, King Saud University, Riyadh 11451, Saudi Arabia (e-mail: dr.m.imran@iee.org).

M. Xia is with the Department of Mechanical Engineering, University of British Columbia, Vancouver, BC V1V 1V7, Canada (e-mail: minxia@mech.ubc.ca).

A. Celesti is with the MIFT Department, University of Messina, Messina 98122, Italy (e-mail: acelesti@unime.it).

Color versions of one or more of the figures in this paper are available online at <http://ieeexplore.ieee.org>.

Digital Object Identifier 10.1109/TII.2019.2899679

presented the edge computing frameworks/strategies. Furthermore, the studies [17], [18] point out that computing-resource scheduling is an important factor to achieve the real-time performance and realize the fusion among of the edge computing, AI and smart manufacturing in a manufacturing environment. Therefore, strategies of scheduling computing resources are a significant perspective for ensuring real-time AI task. This paper will focus on optimizing the solution to this problem and improve the efficiency of resources in a real-time manner. In particular, we will develop a better strategy to improving the efficacy of resources in a real-time manner, reducing the power consumption, and enhancing the reliability of the whole system.

Recent research achievements, such as the fog-assisted manufacturing system [19], and the edge (fog) CPS (ECPS) [20] have provided a preliminary basis to enable the design of resource scheduling for smart manufacturing and Industry 4.0. Nevertheless, we should be aware that computing-resource scheduling to ensure the real-time requirement in the context of Industry 4.0 still faces many challenges. These challenges can be summarized as two fundamental problems: 1) how to propose a computing system to handle and integrate the historical heritage of computing resources and 2) how to construct some novel and efficient strategies and algorithms to ensure the real-time performance.

This article explores resource scheduling strategy for manufacturing in edge computing to provide real-time computing services, from the perspective of the implementation of Industry 4.0. In summary, there are three main contributions of this paper.

- 1) From the AI task operation perspective, a four-level computing system architecture for smart manufacturing is designed for industrial environments, which contributes to integration and fully utilization of different computing resources.
- 2) A two-phase scheduling strategy for the computing resources strategy in edge computing is provided to meet the performance requirement of different AI tasks with consideration of low-latency constraints.
- 3) The proposed scheduling method and the traditional algorithms are compared. The provided algorithm is implemented in a smart manufacturing prototype platform to validate its feasibility and effectiveness.

The remainder of this paper is organized as follows. Section II introduces related work about edge computing, real-time and computing resources scheduling in manufacturing. Section III gives the four-level architecture and the working process for manufacturing computing. Section IV proposes the methods used for computing resources scheduling for edge servers (ESs). In Section V, the experiments for the proposed algorithms are undertaken. Section VI concludes the paper.

II. RELATED WORK

The effective of computing resources control is a necessity to guarantee the low latency and continuous production in smart factories, consequently improving production efficiency and bringing economic benefits. Therefore, in this section, we briefly

outline existing efforts in aspects of edge computing, real-time methods, and resources scheduling strategies in manufacturing.

A. Edge/Fog Computing of Manufacturing

In [21], the authors proposed an architecture of edge computing for IoT-based manufacturing and analyze the function of edge computing in a manufacturing system. In [22], a manufacture inspection system for the smart industry was designed, while it adopted the deep learning models to find out the defects based on fog computing. In [23], the authors proposed a cyber-physical machine tool system based on fog computing-based. Meanwhile, this paper gave the definitions and functions for computer numerical control machines. In [24], the authors divided the data flows into ordinary and emergent streams according to different latency constraints, then adopted the edge computing the adaptive transmission strategies. In [25], a multi-tier multiaccess edge computing framework was provided and its role in the Industry 4.0 was investigated for manufacturing computing performance. In [26], an edge device capable of collecting, processing, storing, and analyzing data is constructed by using a single-board computer and a sensor. All these studies contribute to edge/fog computing in manufacturing. However, these literatures cannot consider the resource unbalance and differences among edge computing servers.

B. Real-Time Schemes in Manufacturing

In [27], the authors discussed the importance of real-time in the industrial system and pointed out that real-time is the most significant evaluation indicator for industrial automation applications. In [28], the wireless transmission characteristics of wireless networks were obtained and analyzed, then according to these characteristics, a real-time big data gathering algorithm for wireless networks is proposed for industrial operations. In [29], an industrial cyber-physical system based on the emerging fog computing paradigm was provided. In the system, machine learning models can be installed to support factory operation. In [30], the authors propose an innovative multi-microprogrammed control unit (MCU) system framework combining a field-programmable-gate-array-based hardware bridge and multiple scalable MCUs to realize an edge computing gateway to get low-latency performance of in industrial IoT. In [31], a ship inspection system based on fog computing was introduced. The system offers identifying and tracking of the pipe tasks, consequently decreasing latency. These literatures mostly focused on strategies of industry networks or adopted edge computing to ensuring real-time performance. However, there are few papers investigating resources scheduling to support real-time in intelligent manufacturing.

C. Resource Scheduling Strategies for Edge Computing

In [32], a novel model for allocating computing resources in an edge computing platform was proposed to allow service providers to establish resource sharing contracts with edge infrastructure. In [33], an intelligent agent at the edge computing was designed to develop a real-time adaptive policy for computational resource allocation for offloaded tasks of multiple

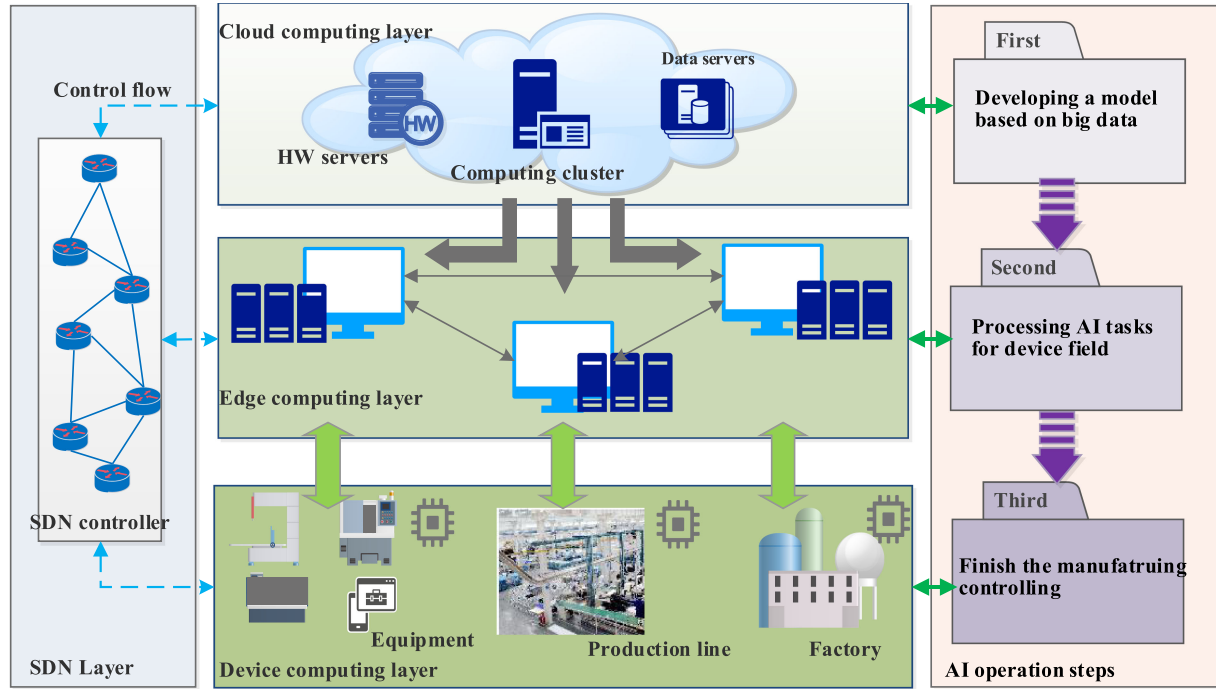


Fig. 1. Architecture for the hybrid computing system.

users in order to improve the system reliability. In [34], the authors formulated the computation offloading problem for mobile edge computing into the system cost minimization problem and present a distributed algorithm consisting of offloading strategy selection by taking into account the completion time and energy. In [35], the authors design a new optimization framework by using an extended Lyapunov technique. In [36], a resource allocation strategy for fog computing based on priced timed Petri nets, by in which the user can choose the satisfying resources autonomously. These studies abovementioned were not suitable for manufacturing, as the real-time constraint was not considered in these algorithms. Moreover, the literature seldom solves the problem by considering the cooperation of multiple computing servers.

III. SYSTEM ARCHITECTURE

This section presents the hybrid computing architecture in manufacturing and the working process of manufacturing computing.

A. Hybrid Computing Architecture in Manufacturing

In the traditional framework, there are two layers to complete the computing tasks: fog/edge and cloud. All the tasks are transmitted into a cloud or edge, on top of the traditional environment, the main drawback of this architecture lies in failing to fulfill the real-time requirement, especially many tasks queuing in ESs. After introducing smart networks nodes, agent devices with limited computing capability, cloud-computing servers and fog-computer servers, the traditional intelligent manufacturing system can be transferred to the hybrid computing architecture. Fig. 1 shows the system architecture for manufacturing comput-

ing using different computing resources for the computational task.

Obviously, cloud servers have the strengths of data storage and computing power; ESs are close to the industrial devices and equipment, thereby having benefits of real-time performance; device computing units can directly drive the mechanical structure; Software defined networking can simply provide the cooperation of different network devices. Hence, from the network perspective, all computing resources are integrated into the hybrid computing architecture to meet the latency requirement. This hybrid architecture essentially contains four parts from the task-node perspective (such as manufacturing devices).

- 1) Device computing layer.
- 2) Edge computing layer.
- 3) Cloud server.,
- 4) Software Defined Network (SDN) layer.

All these elements are collected by the industrial networks (i.e., wired/wireless network). In the cloud layer, the servers are mainly used to resolve the computing-extensive tasks, in which AI model are developed based on different information and big data. In edge computing layer, the edge computing servers are explored to finish the real-time, AI works. Additionally, in the devices computing layer, the devices are mainly responsible for finishing the sensing and controlling works. Besides, the SDN layer is used to control and coordinate different computing layer. There are differences between traditional maintenance and hybrid manufacturing computing architecture.

B. Working Process of Manufacturing Computing

In this paper, we mainly focus on computing resource allocation in the proposed framework. The working process of manufacturing computing is briefly introduced. For the hybrid

computing system, all computing tasks are created on the field devices, including producing machines, wireless network nodes, and mobile elements. Tasks are random events which should usually be processed in real-time manner. For scheduling a task, there are three factors to be considered: compute capability, queueing time, and data communication latency. While according to the three factors, the latency of the task during the special time window can be computed, then in the hybrid computing system, in accordance with the real-time requirement, the AI task can be located at different computing layers.

For edge computing layer, the computing capability and queueing time are the significant factors to determine the task completion time. It is obvious that there are differences for one edge computing layer from computing power, storage power. Since the different servers may deal the task with different complexity, they may have different values of queueing time.

Actually, the historical legacy of computing resources and the low system latency are considered in the presented architecture. In particular, the computing framework can integrate the different level computing resources in a smart factory, such as device computing, ESs, and cloud servers. Therefore, the framework has outstanding performance in term of low latency with the comprehensive utilization of various computing resources, especially for device layer.

IV. RESOURCES SCHEDULING IN EDGE LAYER

This section mainly describes the resource scheduling in the edge layer. We first give the architecture of the edge layer, then present an algorithm for selecting edge computing server and a cooperation strategy for multiple ESs.

A. Architectures of the Edge Layer

After the above analysis, the manufacturing edge layer (MEL), cloud, device computing resources of devices (LCRD) are constructed and connected to the manufacturing computing system via SDN wired /wireless networks. However, cloud servers are typically far from the devices and the system has to spend more time in transmitting the task data between devices and the cloud server. Meanwhile, LCRD is limited by computing capability and responsible for dealing with the necessary tasks with the supporting of the local system normal operation. MEL that is close to the manufacturing equipment, plays the most important role in processing the real-time tasks. As shown in Fig. 2, the MEL consists of multiple ES clusters (ESCs). Every ES is heterogeneous in the capacity of computing, storage and task loads. In MEL, the ESs are connected via the high-bandwidth networks, such as wired links, optical fiber, etc. Therefore, ESs can form an ESC network with low delay. Therefore, every ES is deployed collocating with the devices to fulfill real-time computing tasks.

In order to achieve an efficient task process, ESs are placed in approximation to the devices. The tasks randomly are generated by the manufacturing equipment. They are then arranged and transmitted to the near and suitable ESs to ensure the real-time constraint. Obviously, there are two cases: Single ES can be qualified for the task and single ES cannot be qualified. Therefore, there are two strategies for computing resources

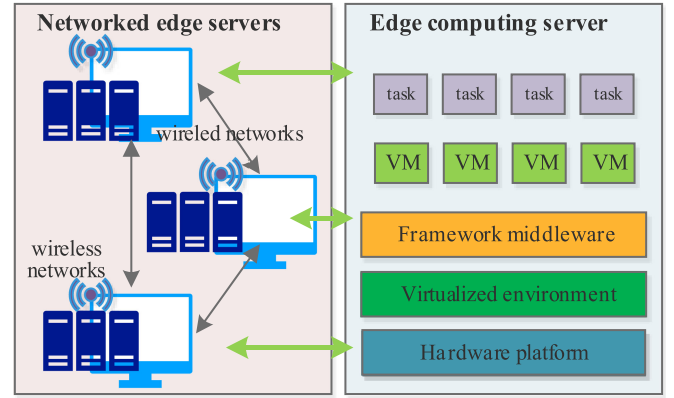


Fig. 2. Mechanism of tasks scheduling for edge computing layer.

scheduling: Selecting algorithms for ES (SAE) and cooperation of edge computing cluster (CEC) to fulfill real-time requirement. The former way can be used to meet the low real-time requirement of computing tasks. SAE scheduling algorithm undertakes to choose the suitable ES from the ES set (ESS) according to the task load, communication time and computing power. Moreover, the CEC is adopted for low-latency requirement, in which one ES cannot qualify to guarantee the low latency.

B. Algorithm for Selecting Edge Computing Server

MEL has the direct impact on the computing performances of the manufacturing task. It is indispensable to propose the scheduling algorithms for MEL and ESS. Scheduling for MEL contains two aspects: Selecting algorithms for ES (SAE) and cooperation of edge computing for low-latency task (CEC). Based on the requirement of specific application for manufacturing, the latency requirement of getting computing results depends on communication, computing and queueing time. In particular, the task x processing time in single ES es T_{task} can be formulated as

$$T_{\text{task}}(x, \text{esc}) = T_{\text{trans}}(x, es) + T_{\text{que}}(x, es) + T_{\text{process}}(x, es) + T_{\text{re}}(x, es) \quad (1)$$

where T_{trans} , T_{que} , T_{process} , T_{re} are the times of transmitting task to ES, queueing, processing, and receiving, respectively. Furthermore, assume that the data size of task and results are φ , \Im and the data rate is v , so T_{trans} , T_{re} can be described as

$$T_{\text{trans}}(x, es) = k_{x, \text{esc}} \frac{\varphi(x)}{v(es)} \quad (2)$$

$$T_{\text{re}}(x, es) = k_{x, \text{esc}} \frac{\Im(x)}{v(es)}. \quad (3)$$

Meanwhile, let X be set of tasks in ES, namely $X = \{x_1, x_2, \dots, x_{|X|}\}$. The set of computer instructions is denoted by $XN = \{xn_1, xn_2, \dots, xn_{|X|}\}$, which is used to deal with X . According to the [23], for the new tasks, the queueing time can be formulated as

$$T_{\text{que}}(x, es) = \sum_{i=1}^{|X|} \sum_{j=1}^{|xn_i|} \frac{IN_j}{V_{\text{process}}} \quad (4)$$

Algorithm 1: Pseudocode of Selecting single ES.

Initialization: Input task x , xn , v , v_{process} , X , XN , E , $t_{\text{requirement}}$, $Es = \emptyset$
Begin: $Es \leftarrow E$
for $i \leftarrow 1$ to $|E|$
 $T_{\text{com}}(e_i) \leftarrow k_{e_i} \cdot \frac{\varphi + \Im}{v}$
 // computing the communication time;
 $t_{\text{process}} \leftarrow \frac{IN \cdot xn}{V_{\text{process}}}$
 // computing the process time of task x ;
 for $j \leftarrow 1$ to $|X_{e_i}|$
 // computing the queuing and process time;
 $t_{\text{que}}(ij) \leftarrow \frac{xn_{ij} IN}{V_{\text{process}}}$
 if $(t_{\text{que}}(ij) \geq t_{\text{requirement}})$
 $Es \leftarrow Es \setminus e_i$
 //Selecting ES according with t_{que} ;
 else // Selecting ES with the total time of task x ;
 for $f \leftarrow 1$ to $|Es|$
 //computing the total time
 $T_{\text{task}}(x, e_f) \leftarrow T_{\text{com}}(e_f) + T_{\text{que}}(e_f) + T_{\text{process}}(e_f)$
 if $(T_{\text{task}}(x, e_f) > t_{\text{requirement}})$
 $Es \leftarrow Es \setminus e_f$
 //Selecting ES according with t_{que} ;
 else
 Break;
 End for
 End if
 End for
End for
Return Es

where IN_j and V_{process} are the j th instruction of the i th task and the process speed of the ES, respectively. In a similar way, we can get the equation for processing this task as follows:

$$T_{\text{process}}(x, es) = \sum_{j=1}^{|x|} \frac{IN_x}{V_{\text{process}}}. \quad (5)$$

According to formulations (1) to (5), the processing time of task x demoted by T_{task} can also be formulated as

$$T_{\text{task}}(x, es) = k_{x, \text{esc}} \frac{\varphi + \Im}{v} + \sum_{i=1}^{|X|} \frac{IN \cdot xn_i}{V_{\text{process}}} + \frac{IN \cdot xn}{V_{\text{process}}}. \quad (6)$$

To ensure the real-time requirement of processing task x , the edge computing server must be subject to the flowing as

$$T_{\text{task}}(x, es) \leq T_{\text{req}}(x). \quad (7)$$

Assume that there are multiple edge computing servers close to the device which contains the task x . For easily understanding, let E be the set of ESs, namely $E = \{e_1, e_2, \dots, e_{|E|}\}$.

So, we propose Algorithm 1 to fulfill the strategy of SAE. In this algorithm, the selecting a single ES strategy can mainly divide into the following steps: First, the system searches the all ESs, and construct the set of E . Then, according to (2)

–(4), we can get the communication time T_{com} and queuing time T_{process} of every ES in the set of E . Third, we evaluate the queuing time to determine whether it is larger than the deadline time of task x . Moreover, we update the candidate set Es of ES to processing the task. Finally, in the light of the total time for resolving the task, we update Es . Finally, the device randomly selects the ES from Es for the task x .

C. Cooperation Strategy of Networked Edge Computing to Achieve Low Latency

In Section IV-B, the low real-time requirement of processing algorithm is given. It is obvious that Algorithm 1 may not deal with the computing-extensive task as there is only one ES assigned for this task. Therefore, we propose a method to cooperate multiple edge computing servers to create ESC to fulfill the latency constraints of single ES.

Indeed, once the ESs are placed into smart factory, they are connected via industrial networks. Then, in the industrial system, the ESs are clustered with cloud servers via software defined network (SDN) controllers according to network distance between ESs and cloud servers. Hence, to achieve low latency, the latter is adopted in the novel framework. The main idea of the method is explained as follows: Selecting an ES as the main server for dividing task and merging the results; and choosing other ESs to cooperate to finish the task according to the latency.

Assume that the task x is be divided into N ($1 \leq N \leq |E|$) subtasks, which are executed in parallel at an ESC to ensure the real-time demands. We denote the set of subtask by $x = \{sx_0, sx_1, sx_2, \dots, sx_{N-1}\}$.

Let $Ec = \{ec_0, ec_1, ec_2, \dots, ec_{N-1}\}$ be the set for cooperating to process the task x . For the subtask $sx_i \in x$ ($0 \leq i \leq N-1$), the communication time can be given as

$$T_{\text{com}}(ec_0, ec_i) = \begin{cases} \frac{D_{\text{rough}}(sx_i) + D_{\text{result}}(sx_i)}{V_{ec_0, ec_i}} & \text{if } i \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

where $D_{\text{rough}}(sx_i)$ and $D_{\text{result}}(sx_i)$ are subtask rough and data size of results, respectively. The term of V_{ec_0, ec_i} is the average data rate between ec_0 and ec_i . Furthermore, in light of the literature [23], we can get the subtask processing time as

$$T_{\text{process}}(sx_i, ec_i) = \frac{IN_{sx_i}}{V_{\text{process}}(ec_i)} \quad (9)$$

where IN_{sx_i} is the subtask instruction number, and $V_{\text{process}}(ec_i)$ is the processing speed of the i th ES of Ec . It is worth mentioning that $|Ec| = N$, while we can get the formulation of $T_{\text{subtask}}(sx_i, ec_i)$ ($0 \leq i \leq N-1$) as shown in (10), according to formulations (1) and (6)

$$T_{\text{subtask}}(sx_i, ec_i) = T_{\text{com}}(ec_0, ec_i) + T_{\text{que}}(sx_i, ec_i) + T_{\text{process}}(sx_i, ec_i) \quad (10)$$

where $T_{\text{com}}(ec_0, ec_i)$ is the communication time between ec_0 and ec_i , $T_{\text{que}}(sx_i, ec_i)$ and $T_{\text{process}}(sx_i, ec_i)$ are queuing time and the process time for subtask sx_i in edge computing server ec_i , respectively.

Recall the fact that the main ES is responsible for dividing task and merging the results. Therefore, the running task time

in main ES is formulated as

$$T_{\text{main}}(x, sx_0, Ec) = T_{\text{divide}}(x) + \sum_{i=1}^{N-1} T_{\text{com}}(ec_0, ec_i) + \text{Max} \\ \times (T_{\text{subtask}}(sx_i, ec_i)) + T_{\text{merge}}(x, Ec) \quad (11)$$

where $T_{\text{divide}}(x)$ and $T_{\text{merge}}(x, Ec)$ are the dividing time and the merging-result time for task x in edge computing servers set Es , respectively. Therefore, the total time for task x running at Ec is decided by the communication time between main edge computing server and the device of task x and the running time $T_{\text{main}}(x, sx_0, Es)$ [as given in (11)]. It is formulated as

$$T_{\text{task}}(x, Es) = T_{\text{main}}(x, sx_0, Ec) + T_{\text{com}}(ec_0, \text{device}_x). \quad (12)$$

It is worth noting that $T_{\text{divide}}, T_{\text{merge}}, T_{\text{com}} \ll T_{\text{process}}$, so according to (11), the equation can be simplified as

$$T_{\text{task}}(x, Es) = \text{Max}(T_{\text{subtask}}(sx_i, ec_i)) + T_{\text{com}}(ec_0, \text{device}_x). \quad (13)$$

Equation (7) gives the time constraints for processing the task. Therefore, we can get

$$\begin{aligned} \text{Max}(T_{\text{subtask}}(sx_i, ec_i)) + T_{\text{com}}(ec_0, \text{device}_x) &\leq T_{\text{req}}(x) \\ \text{Max}(T_{\text{subtask}}(sx_i, ec_i)) &\leq T_{\text{req}}(x) - T_{\text{com}}(ec_0, \text{device}_x). \\ T_{\text{subtask}}(sx_i, ec_i) &\leq T_{\text{req}}(x) - T_{\text{com}}(ec_0, \text{device}_x) \end{aligned} \quad (14)$$

According to the (9) and (14), it is easy to get the task instruction number of the i th ES. It is described as

$$\begin{aligned} \text{IN}_{sx_i} &\leq (T_{\text{req}}(x) - T_{\text{com}}(ec_0, \text{device}_x) \\ &\quad - T_{\text{que}}(sx_i, ec_i))V_{\text{process}}(ec_i). \end{aligned} \quad (15)$$

Furthermore, task time $T_{\text{task}}(x, Es)$ is determined by the maximum T_{process} . Therefore, in light of the above discussion, we propose the strategy of cooperating edge computing servers for the extensive task as shown in Algorithm 2.

In Algorithm 2, the steps can be mainly described as followings: First, according to the referenced equations, we compute subtask instruction number $\text{IN}(E)$ in the constraints of $T_{\text{req}}(x)$ for every ES in the set of E . Then, we sort $\text{IN}(E)$ in the descending order (i.e., from largest to smallest), and create the sorted ES set Es' . Third, we sum the subtask instruction number, temp_sum and evaluate whether the temp_sum meets the requirement of task x . Finally, the main ES divides the task x and finishes the processing task x by Es , and returns the result of the task.

V. ANALYSIS AND EXPERIMENT

This section mainly evaluates the performance of the proposed framework based on an implementation of a prototype in a realistic manufacturing platform.

A. Prototyping Platform and Experiment Installment

For analyzing the proposed strategies of scheduling edge computing resources, a prototyping platform with edge computing servers and industrial IoTs is constructed. As demonstrated in Fig. 3, the prototyping platform is composed of four parts: De-

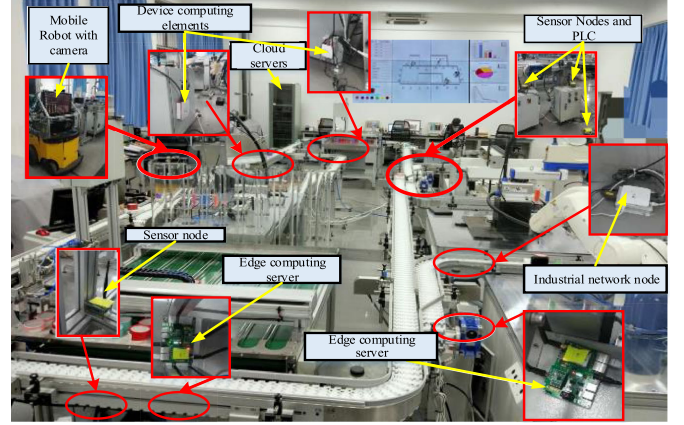


Fig. 3. Prototyping Platform for smart manufacturing with edge computing.

Algorithm 2: Pseudocode of CES.

```

Initialization
for  $k \leftarrow 1$  to  $|E|$ 
     $\text{IN}(e_k) \leftarrow (T_{\text{req}}(x) - T_{\text{com}}(ec_0, \text{device}_x) - T_{\text{que}}(e_k))V_{\text{process}}(e_k)$ 
    //getting the subtask instruction number in the constraints of  $T_{\text{req}}(x)$ 
     $\text{IN}(E) \leftarrow \text{IN}(e_k)$ 
     $Es' \leftarrow \text{sort}(\text{IN}(E))$ 
    //sort () is the function for sorting the E according with the IN ( $e_k$ );
     $\text{main\_ESC} \leftarrow \text{Max}(\text{IN}(E))$ 
    //selecting the main ES
End for
for  $i \leftarrow 1$  to  $|Es'|$ 
     $\text{Temp\_sum} = \text{IN}(es'_i) + \text{Temp\_sum}$ 
    if ( $\text{Temp\_sum} < xn$ )
         $Es \leftarrow Es/es'_i$ 
    else
        break;
    End if
End for
divided_task(x, Es) // divide the task x according with Es
processing_subtask() // processing subtask in selecting edge server
Return_subtask_result()
//returning the subtask result from different edge server
RES  $\leftarrow$  merge_subresult()
//the main edge server merging the results
Return RES

```

vice field; edge computing layer; industrial private cloud servers; and industrial networks (with SDN).

The device field contains multiple types of equipment mobile robots, products processing machines, conveyor belts, and manipulators. To construct the edge computing layer, multiple single-board computers connected with wireless/wired local area network and Bluetooth (i.e., Raspberry Pi 3 Model B) are adopted to serve as ESs. Each single-board computer is

equipped with a quad-core 1.2 GHz Broadcom BCM2837 64bit central processing unit and 1GB random-access memory. For linking with other ESs, we exploit the Ethernet with 100 Mbps bandwidth and adopt servers switch hubs. In addition, wireless communications (Wi-Fi or Bluetooth) are responsible for connection to the devices. From the software aspects, every ES is installed with operating systems (Linux) and OpenCV framework to support image processing tasks for mobile robots or machines. Based on the XenServer and Hadoop cloud ecosystem, we construct a private cloud platform supporting big data processing platform. Meanwhile, the wired ethernet and wireless communication system are used to data interaction between the different layers of smart manufacturing. In a word, an industrial manufacturing prototype platform with ESs has been constructed.

Then based on the prototype platform, an experiment is constructed to evaluate the performance of the proposed. In the experiment, a mobile robot equipped with an industrial camera moved along with a fixed trajectory to monitor the operation of equipment of the constructed industry 4.0 platform. The robot periodically captured pictures and then to calculate the images boundary for the next work process. Then, the running state of the equipment is judged by a neural network. So, to assess the performance of different methods, these pictures are sent to different computing layer via industrial networks. Then, the related parameters are measured during the experiment.

In this paper, a task for processing an image with the average size of 10–20 Mb was executed at the edge computing server (in which typical OpenCV algorithms were executed). The different ESs communicate with each other via Ethernet with bandwidth 100 Mbps. An ES connects with the mobile robot with bandwidth 1–54 Mbps. Moreover, the data is transmitted to the cloud server in four hops via industrial networks. Assume that the communication energy consumption is 1J/s and the computing energy consumption rate is 0.8J/s.

Meanwhile, we use the computing latency, SD and the energy consumption as the evaluation metrics to assess the provided strategies via selecting algorithms for ES (SAE) and cooperating edge computing servers (CEC). We evaluate the performance of the proposed strategies with the comparison with traditional strategies such as cloud server computing, ordinary edge computing (OEC) without scheduling.

In particular, computing latency (CL) is the time between devices transmission data task from devices to the cloud server and receiving results at device. SD is the QoS lever performance metric. SD is expressed as the ratio between requirement computing service time and the computing latency. The higher SD implies the better performance.

B. Analysis and Results

Fig. 4 shows the average of computing latency of different methods with their best performance. In particular, Fig. 4(a) shows that with an increment of the task data size, the computing latency increases for all these methods. It is obvious that the proposal (SAE_CEC) outperforms other strategies, as SAE_CEC will select the better ES to complete the task. More-

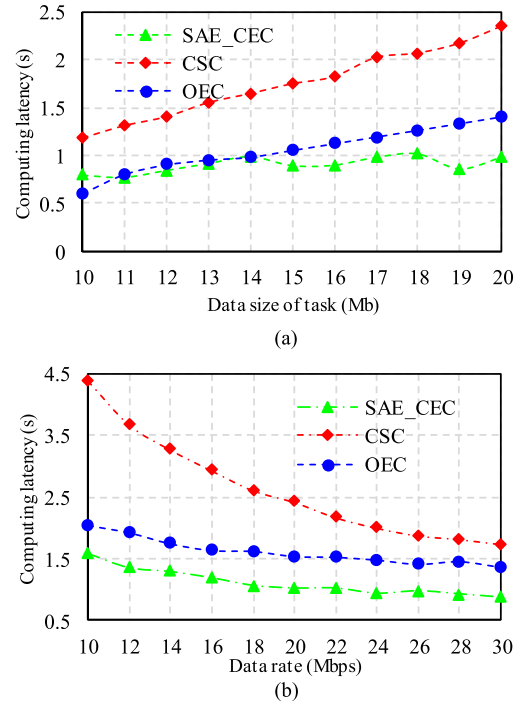


Fig. 4. Comparison of computing latency. (a) Computing latency in different data size of task. (b) Computing latency in different data rate.

over, with the increase the data size, the single ES cannot meet the deadline, SAE_CEC will call Algorithm 2 and exploit multiple ESs together to process the computing-extensive task to ensure the real-time constraints. It is worth mentioning that OEC gets better performance than CSC, as OEC adopts edge computing servers to finish the tasks. Meanwhile, the computing servers far away the device and system have to spend more time in data transmission. Hence, the CSC gets the worst computing latency performance.

Fig. 4(b) demonstrates the results the latency of different methods with the same task data size 10 Mb in varied data rate between devices and server. Similarly, SAE_CEC shows a better performance than other algorithms. Because SAE_CEC method can dynamically select the better ESs according to the deadline of the task. Furthermore, with the increase of data rate the latency of CSC decreases dramatically, as data rate has a great impact on CSC.

In addition, the SD is useful to evaluate the computing services with the consideration of the latency constraints. In particular, Fig. 5(a) gives the result of SD with different deadline requirements and the same data rate and the data size of the task. It is shown that when the data size of the task is less than 14Mb, SDs of three methods are more than 100%. However, when the data size of the task is larger than 14 Mb, only SD of SAE_CEC can get 100%. That is to say, the proposal can meet the real-time requirement for processing the task. Fig. 5(b) shows the SD in different communication bandwidths and data amount of the task with different deadlines (I: 2s, II: 3s, III: 3s). Fig. 5(b) shows the similar results to Fig. 5(a). In summary, the proposals can adapt to different environments and different real-time constraints. The traditional methods cannot be adopted in

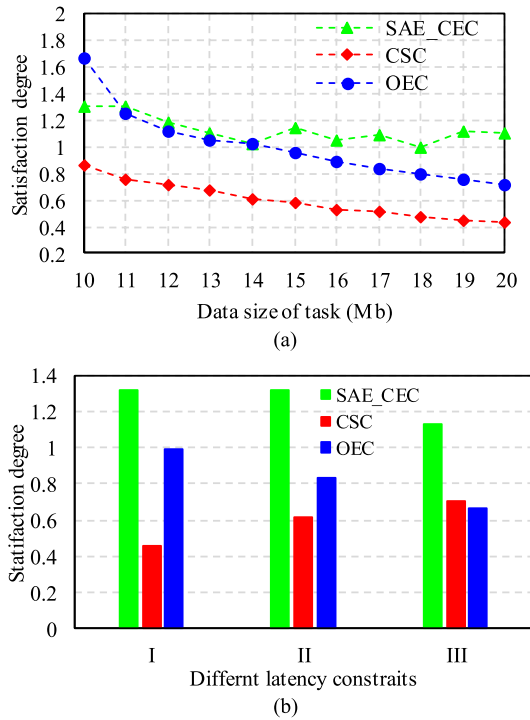


Fig. 5. Comparison of SD. (a) SD in different data size of task. (b) SD in different time constraints.

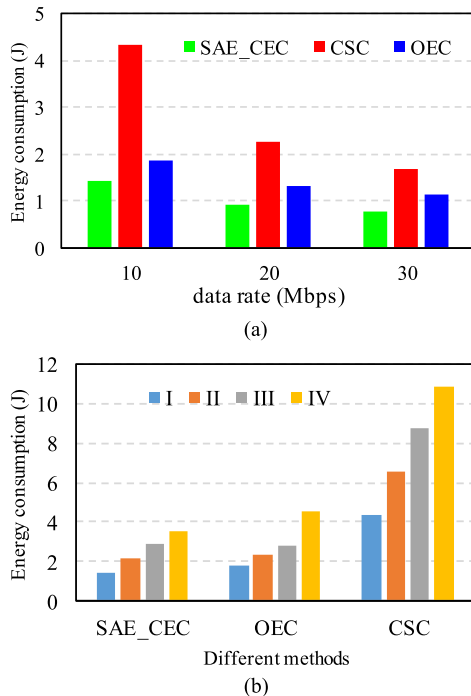


Fig. 6. Comparison of energy consumption. (a) Energy consumption in different data rate. (b) Energy consumption in different task data size.

the industrial maneuvering system due to the failure of fulfilling the real-time requirement.

Moreover, Fig. 6 shows the average energy consumption of different methods in different data rate and task data size. Fig. 6(a) gives the results of energy consumption for differ-

ent strategies. It is demonstrated that with the increment of data rate, the energy consumption will reduce, as all the methods spend less energy in communication. Meanwhile, the proposed scheme outperforms other strategies in terms of energy consumption because of ES being close to device. In particular, when the data rate is 10 Mbps, the proposal can reduce more than 50% energy consumption. Fig. 6(b) shows the energy consumption in different task data size [I: 10, II: 15, III: 20, VI: 20 (Mb)] with same data rate. Similarly, when the data size increases, the energy consumption will be added. Fig. 6(b) also shows the advantage of our method in energy consumption.

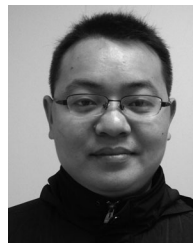
VI. CONCLUSIONS

In this paper, we have focused on resources scheduling to ensure the real-time requirement of smart manufacturing with consideration of integrating different computing resources. First, according to the feature of smart factories, AI task operations and network perspective, we provide the four-layer architecture with integration of historical heritage of computing resources. Then, we focus on the edge computing layer and propose a two-phase scheduling strategy to allocate the computing resources to meet the latency constraint. In the first phase, different factors are considered to select the edge computing server for supporting the computing services for the task with the low real-time constraint. Moreover, the second phase is explored for resolving cooperation of multiple ESs to construct a task of ESC operating the lower latency computing services. Finally, for verifying the feasibility of our proposal, a prototype platform is implemented. In particular, we conduct experiments to compare the traditionally-used methods with the proposed computing-resource scheduling strategy. The proposed computational resources allocation strategies have ensured the real-time for smart manufacturing with edge computing. In summary, the proposed frameworks and computing resources scheduling can accelerate the implementation of Industry 4.0 and smart factory.

REFERENCES

- [1] J. Wan, S. Tang, D. Li, M. Imran, C. Zhang, C. Liu, and Z. Pang, "Reconfigurable smart factory for drug packing in healthcare industry 4.0," *IEEE Trans. Ind. Inform.*, vol. 15, no. 1, pp. 507–516, 2019.
- [2] B. Chen, J. Wan, L. Shu, P. Li, M. Mukherjee, and B. Yin, "Smart factory of industry 4.0: Key technologies, application case, and challenges," *IEEE Access*, vol. 6, pp. 6505–6519, 2018.
- [3] M. Wollschlaeger, T. Sauter, and J. Jasperneite, "The future of industrial communication: Automation networks in the era of the internet of things and industry 4.0," *IEEE Ind. Electron. Mag.*, vol. 11, no. 1, pp. 17–27, Mar. 2017.
- [4] J. He, G. Jia, G. Han, H. Wang, and X. Yang, "Locality-aware replacement algorithm in flash memory to optimize cloud computing for smart factory of industry 4.0," *IEEE Access*, vol. 5, pp. 16252–16262, 2017.
- [5] W. Zhang, Y. Liu, G. Han, Y. Feng, and Y. Zhao, "An energy efficient and QoS aware routing algorithm based on data classification for industrial wireless sensor networks," *IEEE Access*, vol. 6, pp. 46495–46504, 2018.
- [6] Y. Bi, G. Han, C. Lin, Q. Deng, L. Guo, and F. Li, "Mobility support for fog computing: An SDN approach," *IEEE Commun. Mag.*, vol. 56, no. 5, pp. 53–59, May 2018.
- [7] X. Li, D. Li, J. Wan, A. V. Vasilakos, C. F. Lai, and S. Wang, "A review of industrial wireless networks in the context of industry 4.0," *Wireless Netw.*, vol. 23, no. 1, pp. 23–41, 2017.

- [8] J. Wan, B. Yin, D. Li, A. Celesti, F. Tao, and Q. Hua, "An ontology-based resource reconfiguration method for manufacturing cyber-physical systems," *IEEE/ASME Trans. Mechatronics*, vol. 23, no. 6, pp. 2537–2546, Oct. 2018.
- [9] P. K. Illa and N. Padhi, "Practical guide to smart factory transition using IoT, big data and edge analytics," *IEEE Access*, vol. 6, pp. 55162–55170, 2018.
- [10] G. Han, M. Guizani, J. Lloret, S. Chan, L. Wan, and W. Guibene, "Emerging trends, issues, and challenges in big data and its implementation toward future smart cities: Part 2," *IEEE Commun. Mag.*, vol. 56, no. 2, pp. 76–77, Feb. 2018.
- [11] E. Shellshear, R. Berlin, and J. S. Carlson, "Maximizing smart factory systems by incrementally updating point clouds," *IEEE Comput. Graph. Appl.*, vol. 35, no. 2, pp. 62–69, Mar./Apr. 2015.
- [12] D. Georgakopoulos, P. P. Jayaraman, M. Fazio, M. Villari, and R. Ranjan, "Internet of things and edge cloud computing roadmap for manufacturing," *IEEE Cloud Comput.*, vol. 3, no. 4, pp. 66–73, Jul./Aug. 2016.
- [13] M. Marjanović, A. Antonić, and I. P. Žarko, "Edge Computing Architecture for Mobile Crowdsensing," *IEEE Access*, vol. 6, pp. 10662–10674, 2018.
- [14] T. Taleb, K. Samdanis, B. Mada, H. Flinck, S. Dutta, and D. Sabella, "On multi-access edge computing: A survey of the emerging 5g network edge cloud architecture and orchestration," *IEEE Commun. Surveys Tuts.*, vol. 19, no. 3, pp. 1657–1681, Jul.–Sep. 2017.
- [15] 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./Feb. 2018.
- [16] Y. Sun, S. Zhou, and J. Xu, "EMM: Energy-aware mobility management for mobile edge computing in ultra dense networks," *IEEE J. Sel. Areas Commun.*, vol. 35, no. 11, pp. 2637–2646, Nov. 2017.
- [17] J. Xu, B. Palanisamy, H. Ludwig, and Q. Wang, "Zenith: Utility-aware resource allocation for edge computing," in *Proc. IEEE Int. Conf. Edge Comput.*, 2017, pp. 47–54.
- [18] L. Yin, J. Luo, and H. Luo, "Tasks scheduling and resource allocation in fog computing based on containers for smart manufacturing," *IEEE Trans. Ind. Inform.*, vol. 14, no. 10, pp. 4712–4721, Oct. 2018.
- [19] T. M. Fernándezcaramés, P. Fragalamas, M. Suárezalbela, and M. Villarmonesinos, "A fog computing and cloudlet based augmented reality system for the industry 4.0 shipyard," *Sensors*, vol. 18, no. 6, 2018.
- [20] M. Aazam, S. Zeadally, and K. A. Harras, "Deploying fog computing in industrial internet of things and industry 4.0," *IEEE Trans. Ind. Inform.*, vol. 14, no. 10, pp. 4674–4682, Oct. 2018.
- [21] B. Chen, J. Wan, A. Celesti, D. Li, H. Abbas, and Q. Zhang, "Edge computing in IoT-based manufacturing," *IEEE Commun. Mag.*, vol. 56, no. 9, pp. 103–109, Sep. 2018.
- [22] L. Li, K. Ota, and M. Dong, "Deep learning for smart industry: Efficient manufacture inspection system with fog computing," *IEEE Trans. Ind. Inform.*, vol. 14, no. 10, pp. 4665–4673, Oct. 2018.
- [23] Z. Zhou, J. Hu, Q. Liu, P. Lou, J. Yan, and W. Li, "Fog computing-based cyber-physical machine tool system," *IEEE Access*, vol. 6, pp. 44580–44590, 2018.
- [24] X. Li, D. Li, J. Wan, C. Liu, and M. Imran, "Adaptive transmission optimization in sdn-based industrial internet of things with edge computing," *IEEE Internet Things J.*, vol. 5, no. 3, pp. 1351–1360, Jun. 2018.
- [25] N. Dao, Y. Lee, S. Cho, E. Kim, K. Chung, and C. Keum, "Multi-tier multi-access edge computing: The role for the fourth industrial revolution," *Proc. Int. Conf. Inf. Commun. Technol. Convergence*, 2017, pp. 1280–1282.
- [26] D. Park, S. Kim, Y. An, and J. Y. Jung, "Lired: A light-weight real-time fault detection system for edge computing using LSTM recurrent neural networks," *Sensors*, vol. 18, no. 7, 2018.
- [27] V. C. Gungor, and G. P. Hancke, "Industrial wireless sensor networks: challenges, design principles, and technical approaches," *IEEE Trans. Ind. Electron.*, vol. 56, no. 10, pp. 4258–4265, Oct. 2009.
- [28] X. Ding, Y. Tian, and Y. Yu, "A real-time big data gathering algorithm based on indoor wireless sensor networks for risk analysis of industrial operations," *IEEE Trans. Ind. Inform.*, vol. 12, no. 3, pp. 1232–1242, Jun. 2016.
- [29] P. O'Donovan, C. Gallagher, K. Bruton, and D. T. J. O. Sullivan, "A fog computing industrial cyber-physical system for embedded low-latency machine learning industry 4.0 applications," *Manuf. Lett.*, vol. 15, pp. 139–142, 2018.
- [30] C. Chen, M. Lin, and C. Liu, "Edge computing gateway of the industrial internet of things using multiple collaborative microcontrollers," *IEEE Netw.*, vol. 32, no. 1, pp. 24–32, Jan./Feb. 2018.
- [31] F. Tiago, P. Fragalamas, S. Manuel, and D. Manuel, "A fog computing based cyber-physical system for the automation of pipe-related tasks in the industry 4.0 shipyard," *Sensors*, vol. 18, no. 6, 2018.
- [32] J. Xu, B. Palanisamy, H. Ludwig, and Q. Wang, "Zenith: Utility-aware resource allocation for edge computing," *Proc. IEEE Int. Conf. Edge Comput.*, 2017, pp. 47–54.
- [33] T. Yang, Y. Hu, M. C. Gursoy, A. Schmeink, and R. Mathar, "Deep reinforcement learning based resource allocation in low latency edge computing networks," in *Proc. 15th Int. Symp. Wireless Commun. Syst.*, 2018, pp. 1–5.
- [34] H. Yu, Q. Wang, and S. Guo, "Energy-efficient task offloading and resource scheduling for mobile edge computing," in *Proc. IEEE Int. Conf. Netw. Architecture Storage*, 2018, pp. 1–4.
- [35] X. Wang, K. Wang, S. Wu, S. Di, H. Jin, K. Yang, and S. Ou, "Dynamic resource scheduling in mobile edge cloud with cloud radio access network," *IEEE Trans. Parallel Distrib. Syst.*, vol. 29, no. 11, pp. 2429–2445, 2018.
- [36] L. Ni, J. Zhang, C. Jiang, C. Yan, and K. Yu, "Resource allocation strategy in fog computing based on priced timed petri nets," *IEEE Internet Things J.*, vol. 4, no. 5, pp. 1216–1228, 2017.



Xiaomin Li received the B.S. degree in computer science and engineering from the Air Force Engineering University, Xi'an, China, in 2004, the M.S. degree in agricultural mechanization engineering from the South China Agricultural University, Guangzhou, China, in 2014, and the Ph.D. degree in mechatronic engineering from the South China University of Technology, Guangzhou, China, in 2018.

He is currently an Associate Professor with the Zhongkai University of Agriculture and Engineering, Guangzhou, China. His research interests include cyber-physical systems, smart manufacturing, big data, and wireless network.



Jiafu Wan (M'11) received the B.S. degree in automation from North University of China, Taiyuan, China, in 2000, the M.S. degree in control theory and control engineering from Guangdong University of Technology, Guangzhou, China, in 2003, and the Ph.D. degree in mechatronic engineering from South China University of Technology, Guangzhou, China, in 2008.

He is currently a Professor with the School of Mechanical and Automotive Engineering, South China University of Technology. He has directed

18 research projects, including the National Key Research and Development Project, and the Joint Fund of the National Natural Science Foundation of China and Guangdong Province. He has authored or coauthored more than 150 scientific papers, including 90+ SCI-indexed papers, 40+ IEEE Trans./Journal papers, 18 ESI Highly Cited Papers and 4 ESI Hot Papers. According to Google Scholar Citations, his published work has been cited more than 6700 times (H-index = 40). His SCI other citations, sum of times cited without self-citations, reached 1700 times (H-index = 26) according to Web of Science Core Collection. His research interests include cyber-physical systems, intelligent manufacturing, big data analytics, Industry 4.0, smart factory, and cloud robotics.



Hong-Ning Dai (SM'16) received the B.Eng. and M.Eng. degrees in computer science and engineering from South China University of Technology, Guangzhou, China, in 2000 and 2003, respectively, and the Ph.D. degree in computer science and engineering from the Department of Computer Science and Engineering, Chinese University of Hong Kong, Hong Kong, in 2008.

He is currently with the Faculty of Information Technology, Macau University of Science and Technology, as an Associate Professor. He has authored or coauthored more than 80 peer-reviewed papers in refereed journals and conferences. He is also a holder of 1 U.S. patent and 1 Australia innovation patent. His current research interests include big data analytics and Internet of Things. He holds visiting positions at the Hong Kong University of Science and Technology, the University of New South Wales, Sun Yat-sen University, University of Electronic Science and Technology of China and Hong Kong Applied Science and Technology Research Institute, respectively.

Dr. Dai was the recipient of the Bank of China (BOC) Excellent Research Award of Macau University of Science and Technology in 2015. He was a Guest Editor for the IEEE TRANSACTIONS ON INDUSTRIAL INFORMATICS and an Editor for *International Journal of Wireless and Mobile Communication for Industrial Systems*. He is a Professional Member of the Association for Computing Machinery.



Muhammad Imran received the MCS degree in computer science from National University of Modern Languages, Islamabad, Pakistan, in 2005, the M.S. degree with majors in communication and networks from Federal Urdu University of Arts, Science and Technology, Karachi, Pakistan, in 2007, and the Ph.D. degree in information technology from the University Teknologi PETRONAS, Seri Iskandar, Malaysia in 2011.

He is currently an Assistant Professor with the College of Computer and Information Sciences, King Saud University, Riyadh, Saudi Arabia. His research is financially supported by several grants and he has completed a number of international collaborative research projects with reputable universities. He has authored or coauthored more than 150 research articles in top international conferences and journals. He served/serving as a Guest Editor for more than a dozen special issues in the top international journals. He has been involved in more than 75 conferences and workshops in various capacities such as a chair, co-chair and technical program committee member. European Alliance for Innovation (EAI) has appointed him as a Co-Editor in Chief for EAI Transactions on Pervasive Health and Technology. His research interest includes Internet of Things, big data analytics, intelligent transportation systems, cloud and edge computing, and security and privacy.

Dr. Imran was an Associate Editor for reputable international journals, such as IEEE COMMUNICATIONS MAGAZINE, FUTURE GENERATION COMPUTER SYSTEMS, IEEE ACCESS, *Wireless Communication and Mobile Computing Journal*, *Ad Hoc and Sensor Wireless Networks Journal*, *IET Wireless Sensor Systems*, *International Journal of Autonomous and Adaptive Communication Systems*.



Min Xia received B.S. degree in industrial engineering from Southeast University, Nanjing, China in 2009, the M.S. degree in precision machinery and instrumentations from the University of Science and Technology of China, Hefei, China, in 2012 and the Ph.D. degree in mechanical engineering from the University of British Columbia, Vancouver, BC, Canada in 2017.

He is currently a Postdoctoral Research Fellow with the University of British Columbia, Vancouver, BC, Canada. His research interests include smart manufacturing, machine diagnostics and prognostics, deep neural networks, wireless sensor network, and sensor fusion.



Antonio Celesti received the B.S. and M.S. degrees in computer science in 2006 and 2008, respectively, from the University of Messina, Messina, Italy, where he received the Ph.D. degree in advanced technologies for information engineering in 2012.

He is a Professor of database II (advanced database) and a Fellow with the Scientific Research Organizational Unit of University of Messina. He was a Collaborator for the University of Messina in several international research

projects including EU FP7 RESERVOIR (2008–2011), EU FP7 VISION CLOUD, EU FP7 frontierCities (2010–2013), EU Horizon 2020 BEACON (2015–2017) and he was a Principal Investigator with the EU Horizon 2020 frontierCities2 - FiTech - CASMOB (2017–2018). He is responsible for the University of Messina of the Italian Ministry Healthcare founded project entitled “Do Severe acquired brain injury patients benefit from Telerehabilitation? A Cost-effectiveness analysis study”. He is co-author of more than 140 scientific documents published in international journals, books, and conference proceedings. His main research interests includes distributed systems, cloud computing, edge computing, Internet of Things, artificial intelligence, big data, and system federation and security.