# Intelligent Radio Resource Allocation for Human-Robot Collaboration

YE FENG [ID][1] (Student Member, IEEE), LIHUA RUAN [ID][2] (Member, IEEE),
AMPALAVANAPILLAI NIRMALATHAS [ID][1] (Senior Member, IEEE),
AND ELAINE WONG [ID][1] (Senior Member, IEEE)

[1]Department of Electrical and Electronic Engineering, The University of Melbourne, Melbourne, VIC 3010, Australia

[2]School of Science and Engineering, The Chinese University of Hong Kong, Shenzhen 518172, China

CORRESPONDING AUTHOR: E. WONG (e-mail: ewon@unimelb.edu.au)

**ABSTRACT** The recent surge in human-controlled robotics and haptic devices research is expediting a paradigm shift in today's communication networks towards human-to-robot (H2R) centric technologies that support Industrial Internet of Things (IIoT) and Industry 5.0. In both IIoT and Industry 5.0, human skills are extended through collaboration with robots that are geographically separated from the human. Depending on the dynamicity of the actual use case, human-to-robot communications necessitate low-latency networking. While Long Term Evolution (LTE) cellular technology has been successful in fulfilling the bandwidth demands of massively-connected sensors and devices of Industry 4.0, it is insufficient to meet the low latency demands of the future Industry 5.0 where dynamic interactions between humans and robots are paramount. In reducing the latency caused by radio resource contention in wireless H2R communications, in this work, we propose a novel approach that exploits an Attention-based Recurrent Neural Network (Att-RNN) to improve the Semi-Persistent Scheduling (SPS) resource allocation scheme adopted by LTE and new radio (NR) standards developed for the fifth generation (5G) mobile networks. We conduct a series of real haptic experiments to collect H2R traffic traces to train, test and evaluate the accuracy of Att-RNN in predicting H2R traffic. Then, with extensive simulations based on the empirical H2R traffic traces, we show that our proposed Att-RNN SPS scheme outperforms classic SPS and other existing resource allocation schemes in terms of reduced latency and improved resource allocation efficiency, thus making Att-RNN SPS a suitable candidate in future Industry 5.0 deployments.

**INDEX TERMS** Haptic communications, remote human-to-robot collaboration, low-latency, long term evolution, prediction methods, machine learning, 5G.

## I. INTRODUCTION

DURING the last half of the past decade, our society has witnessed an unprecedented advancement of Internet-of-Things (IoT) technology and its corresponding disruptive services, which in turn has led to an eco-system of IoT-based smart cities, smart healthcare systems, and smart factories. While NB-IoT [1] and Cat-M1 [2] have been widely deployed to support IoT networks which predominantly carry machine-to-machine (M2M) type traffic, these Long Term Evolution (LTE)-based networks may not be able to adequately support future human-to-robot/machine (H2R/M) traffic which incurs additional reliability and latency performance demands on the network [3], [4]. In [5],

we presented a solution to achieve ultra-low latency uplink transmission protocol in wireless access networks by exploiting an existing uplink medium access control protocol called Semi-Persistent Scheduling (SPS) as defined by LTE and 5G specifications. While our results show promise, it highlighted a limitation whereby the accuracy of our traffic prediction algorithm based on the statistical auto-regressive (AR) model and in turn the latency performance, degrade when supporting bursty traffic. According to [6], [7], the traffic distribution of haptic telerobot and haptic VR applications are characterized by the generalized Pareto distribution with bursty traffic arrivals. As these applications are key technology enablers of Industry 5.0 as illustrated by the roadmap shown
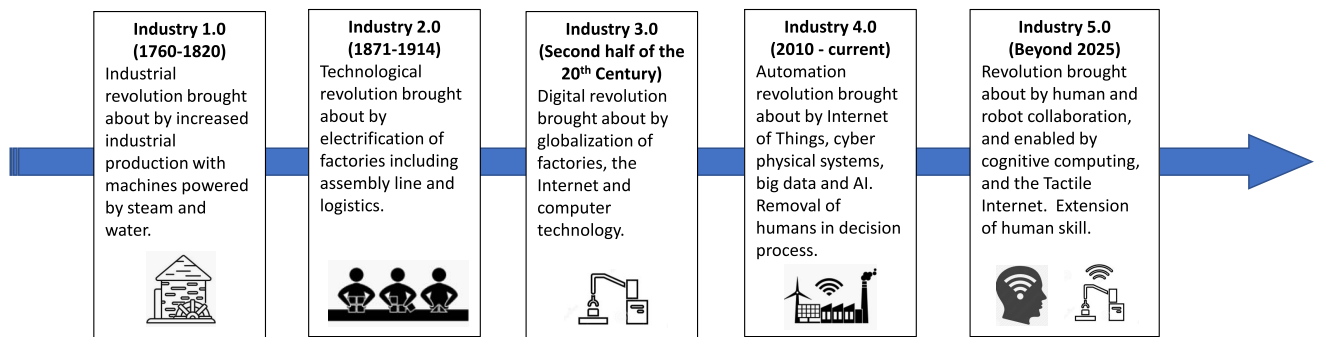
**FIGURE 1.** The roadmap towards Industry 5.0.

in Fig. 1, investigation into advanced algorithms that yield more accurate prediction is needed to improve resource block pair (RBP) allocation in predictive SPS schemes that support future H2R applications. This motivates our research into applying more advanced techniques with high prediction accuracy in support of bursty traffic.

In this article, we harness the use of machine learning intelligence to address the shortcoming in [5] with particular focus on supporting human-robot/machine collaboration for Industry 5.0. With the advancement of the processing power of contemporary computers and the availability of large dataset, tools from the machine learning paradigm have been proven effective for many purposes including regression, classification and decision making processes [8], [9]. Contrary to traditional computational algorithms which need to be explicitly programmed, machine learning algorithms can learn and improve themselves from the input data and are expected to understand the structure/feature of a given dataset. Three basic machine learning paradigms exist, namely supervised and unsupervised machine learning and reinforcement learning. The supervised machine learning takes labeled data as its input and applies statistical analysis, Bayesian models etc. to perform regression or classification on unlabeled data. When learning from the data, machine learning algorithms improve themselves by evaluating the deviation between the predicted and given labels. Examples of regression-type machine leaning are weather forecasting, future stock prediction and natural language processing, while classification-type machine learning involves facial recognition, optical character recognition and future event prediction. On the other hand, unsupervised machine learning algorithms are provided with unlabeled or partially labelled data, and they need to learn the hidden feature, pattern and commonalities among the data so that categorization and labeling can be performed [10]. Finally, reinforcement learning works by interacting with the known/unknown environment. It deploys an agent to take action on every encountered situation within the *environment*, and this agent receives rewards/penalties for performing correctly/incorrectly. By deploying reinforcement learning techniques, the agent is expected to learn the dynamics of the environment and causality between an

action and its resulting reward/penalty, in which case the performance of the agent can be improved by choosing the correct action. Different to supervised learning, where the ground-truth label of input is provided, the agent of a reinforcement learning algorithm only receives an immediate reward for its action, which does not necessarily represent the overall worthiness of that action. In other words, an action resulting in the best immediate reward might not be the correct choice to achieve the highest long-term gain.

In this paper, we harness supervised machine learning (ML) to make accurate predictions of H2R traffic for resource pre-allocation in SPS schemes of LTE and 5G networks. Below are a list of major contributions of this article:

- Proposal of a Feasible RBP Allocator (FRA) to determine the positions of scheduled resources in the predictive SPS scheme. Specifically, the FRA features radio resource collision avoidance amongst SPS-connected User Equipment (UE) and improves allocation efficiency by maintaining accessibility of unused resources for non-H2R UEs.
- Proposal of a novel Attention-based Recurrent Neural Network (Att-RNN) prediction model that uses empirical H2R traffic characteristics for prediction. In particular, the sequential-format input and output from the predictive SPS model is firstly processed with dimension expansion and logarithmic normalization to allow high-dimension feature detection and to remove data skewness. Then, data is fed into an encoder-decoder architecture such that temporal dependencies within the input can be captured. An additive attention mechanism is implemented to let the proposed model selectively focus on the most relevant input states for improved prediction accuracy. Finally, at the prediction stage, an initial input with learnable parameters is introduced at the decoder, which enables the model to perform with extra flexibility.
- Comparisons of the accuracy of the proposed Att-RNN model with that of existing statistical AR and ML-based models, including feed-forward artificial neural network

(ANN), Long Short-Term Memory (LSTM), and several other attention-based ML, in predicting H2R traffic. Normalized mean square error (NMSE) and mean absolute error (MAE) performances of each model are presented, highlighting the accurate prediction achieved by the proposed Att-RNN model.

- Proposal of an SPS scheme based on the Att-RNN prediction model, i.e., Att-RNN SPS scheme, that effectively pre-allocates resources for H2R traffic, and thereby achieving the latency metric demanded by haptic feedback delivery. System-level simulation results reveal that our proposed Att-RNN SPS scheme can improve latency performance and scheduling efficiency in support of bursty H2R traffic, as compared to existing AR-based predictive SPS scheme.

Evidently, this study differs from [5] through complementing the design of predictive SPS scheme via a novel resource block allocator that determines the positions of scheduled resources and through evaluating latency performance and prediction accuracy of the proposed Att-RNN SPS scheme in a more complex and realistic scenario when both H2R and non-H2R traffic exist.

The rest of this article is structured as follows. In Section II, medium access control (MAC) layer protocols of mobile cellular networks and related work on latency reduction are briefly discussed. Section III describes the general working principle of the Att-RNN SPS scheme and how positions of scheduled resources are determined by the novel FRA. Section IV presents an overview on recurrent neural network and its application in communication networks, followed by the detailed structure of Att-RNN prediction model for the proposed SPS scheme, following an overview on recurrent neural network and its application in communication networks. The implementation of the Att-RNN model for predictive SPS scheme and its prediction performance are discussed in Section V. In Section VI, the performance of the Att-RNN model is evaluated against that of existing schemes. Finally, our findings arising from this work are summarized in Section VII.

## II. EXISTING LOW-LATENCY SCHEDULING SCHEMES FOR MOBILE CELLULAR NETWORKS

Currently, the communications community is striving to improve mobile cellular networks so that emerging new services can be supported with ultra-low latency. The extension of current LTE technology into the indoor network segment to ensure end-to-end connectivity may not be sufficient to support high-quality mobile real-time H2R applications [11]. In particular, the uplink latency caused by contention in wireless resources may severely affect H2R collaboration. Therefore, improving the uplink resource allocation scheme and the corresponding latency performance are our main focus. The uplink latency consists of processing latency $T_{proc}$, queueing latency $T_{queue}$, transmission latency $T_{trans}$ and propagation latency $T_{prop}$, where $T_{trans}$ depends on processing speed of communications devices
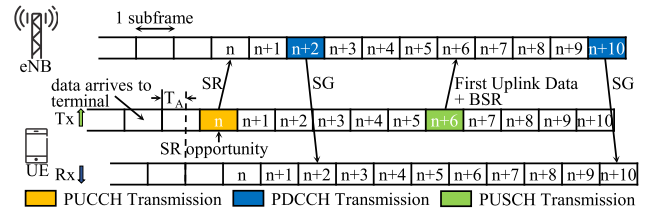


**FIGURE 2.** Timing diagram of Dynamic Scheduling (DS) Scheme in LTE.

and frame structures of LTE standards and $T_{prop}$ depends on transmission distance and speed of radio signal. In other words, $T_{trans}$ and $T_{prop}$ can only be reduced by redefining mobile cellular network standards and upgrading hardware. On the other hand, $T_{proc}$ is affected by the transmission protocols deployed by the communication standards and $T_{queue}$ is largely defined by scheduling functions. In this section, we discuss scheduling algorithms that bring down uplink latency.

### A. DYNAMIC SCHEDULING (DS) SCHEME

A conventional LTE network uses the DS scheme for uplink transmission. An overview of this operation is illustrated in Fig. 2. Once a UE is connected and synchronized to an E-UTRAN NodeB (eNB), the UE's uplink and downlink transmissions are then initialized with transmit timing advance [12]. The transmit timing advance–$T_A$ sets a negative offset between the start of a received downlink subframe and a transmitted uplink subframe, such that time alignment between the eNB and the UE can be achieved. In practice, $T_A$ is determined by propagation latency from UE to eNB. Once the UE detects the arrival of new uplink data, it transmits a Scheduling Request (SR) to the eNB on Physical Uplink Control Channel (PUCCH) at the next SR opportunity–subframe $n$. Upon successful reception of the SR, the eNB transmits a downlink control packet on Physical Downlink Control Channel (PDCCH) called Scheduling Grant (SG) at subframe $n + 2$, informing the scheduled RBP dedicated to the UE. The first uplink data frame is then transmitted by the UE on Physical Uplink Shared Channel (PUSCH), at 4 subframes after the UE received the SG. An important feature of the LTE uplink data frame is that it contains not only the uplink data but also a Buffer Status Report (BSR), indicating the size of leftover data in UE's buffer. By obtaining the buffer occupancy information, the eNB can dynamically schedules further RBPs and sends the second SG to the UE at subframe $n + 10$. The subsequent uplink data packets are therefore transmitted every 8 subframes until the UE's data buffer is cleared.

The study in [11] showed that instead of being limited by LTE bandwidth, the latency incurred by the signalling process of LTE networks dominates the network latency. Instead, existing research is focused on designing scheduling functions in the DS scheme [13]–[16], which only affects $T_{queue}$ and therefore has limited impact on the overall latency performance.
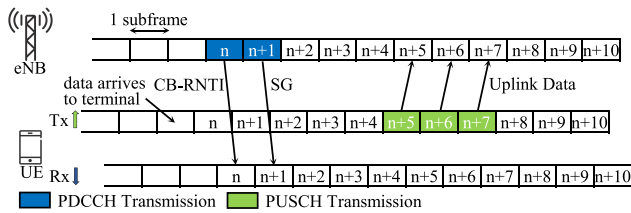
**FIGURE 3.** Timing diagram of Contention Based Scheme (CBS) in LTE.



**FIGURE 4.** Timing diagram of Semi-Persistent Scheduling (SPS) in LTE.

**TABLE 1.** Preliminary analysis of uplink transmission latency.

| Component | DS | CBS | SPS |
|---|---|---|---|
| Mean waiting time for transmission opportunity | 4 | 0.5 | 0.5 |
| PUSCH transmission | 1 | 1 | 1 |
| Minimum uplink latency | 5 | 1.5 | 1.5 |

### B. CONTENTION-BASED SCHEME

Contention Based Schemes (CBSs) have been proposed in [17], [18] to address the latency caused by $T_{proc}$. In general, CBS achieves statistical multiplexing by allowing multiple UEs to content for the same PUSCH resources for uplink transmission. First, an eNB assigns a group of RBPs that can be shared by multiple UEs using a contention-based (CB) grant and delivers this information via CB-Radio Network Temporary Identifiers (CB-RNTIs). By decoding CB grants, the UEs are allowed to send their uplink data directly using the shared PUSCH resources without waiting for SR or SG [19]. The procedure of CBS uplink transmission is illustrated in Fig. 3. Without the need to exchange SR and SG between UEs and eNB, CBS is able to conduct uplink transmissions with less latency. However, CBS suffers from data packet collision when multiple UEs transmit data using the same RBPs during the same subframe, in which case uplink packets cannot be correctly received and decoded.

Andreev *et al.* implemented CBS for Machine-Type Communications (MTC), which is characterized by small data transmissions [20]. Their solution to reducing data collision is a random backoff procedure followed by a packet retransmission mechanism. In [21], CBS was tested and evaluated with the implementation of consecutive Hybrid Automatic Repeat reQuest (HARQ) retransmission. The proposed HARQ retransmission approach is used to increase the reliability of CBS, and it is shown that CBS is more efficient in terms of resource scheduling than conventional SPS scheme and the proposed CBS is able to meet reliability requirement, given data collision is properly handled.

### C. SEMI-PERSISTENT SCHEDULING (SPS) SCHEME

The conventional DS scheme imposes a significant amount of control-signal overhead when transmitting small-sized packets. Additionally, instead of dynamic resource allocation, the regularly occurring pattern of constant bitrate (CBR) traffic motivates an uplink scheme to schedule resources periodically and deterministically. As such, the SPS scheme was proposed to address these two issues, and it is most notably implemented for VoIP [22]. In essence, an eNB periodically assigns predefined RBPs to an SPS-configured UE without consuming resource on PDCCH or PUCCH, except for the initial connection setup stage. The principle of operation of the SPS scheme with a 10-ms period is illustrated in Fig. 4. The SPS scheme follows a similar process as in the DS
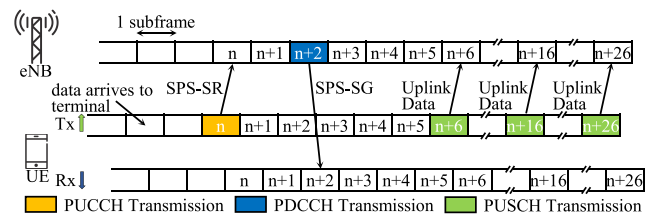
scheme until subframe $n + 6$. In particular, all SPS-related parameters including SPS periodicity, scheduling decision, selected Modulation and Coding Scheme (MCS) and a maximum number of SPS transmission are encapsulated in an SPS-SG. At subframes $n + 16$ and $n + 26$, a UE initiates uplink transmission using the granted periodic RBPs, and this continues until the maximum number of SPS transmission is reached or other termination mechanisms are triggered. Such a conventional SPS scheme may not well suit H2R communications since H2R traffic is bursty and UE transmissions cannot be timely adjusted once RBPs are scheduled. This leads to a trade-off between low/high latency and high/low resource waste in an over/under-scheduling problem. In particular, SPS might fail to transmit all the buffered data immediately when data traffic is bursty [23] or induce noticeable resource waste when data traffic reduces.

When the contention among UEs for uplink radio resources is light, the minimum uplink latency of DS, CBS and SPS schemes are presented in Table 1 based on [24]. It should be noted that the transmission latency is valid when sufficient radio resources are accessible. We can see that without exchanging SR and SG between UEs and eNB, CBS and SPS offers 3.5 ms less transmission latency as compared with the DS. When CBS is implemented, however, packet collision is likely to happen, which necessitates the implementation of collision reduction mechanisms and efficient retransmission techniques [21], [25]. This motivates our interest in designing an SPS scheme that allows flexible scheduling and ensures a contention-free transmission with low-latency performance.

### D. EXISTING LITERATURE ON LOW-LATENCY SPS SCHEMES

Avocanh *et al.* improved the SPS scheme with a MAC Layer method–an SPS scheme with Provisioning (SPS-P) that proactively allocates resources to consecutive future frames by prediction [26]. The prediction model used by the SPS-P is based on GBAR (Gamma-Beta Auto Regressive) algorithm, which is specialized for estimating short-term
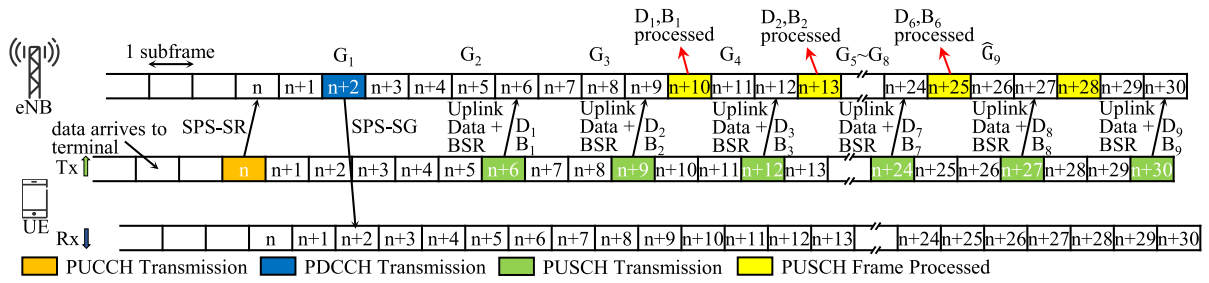
**FIGURE 5.** Timing diagram of the predictive SPS scheme with a 3 ms SPS periodicity [5].

fluctuation of Videotelephony data traffic. Compared to DS scheme, the SPS-P reduces packet loss by seven times and effectively improved the latency performance. Note that the GBAR algorithm in SPS-P is efficient for arrivals having variable packet sizes and constant arriving time. However, for bursty H2R arrivals, this scheme can not be directly applied due to inaccurate prediction for bursty traffic. Another drawback of the SPS-P is that the prediction window is limited, leaving it unsuitable for long-term traffic prediction. This, in turn, would require the SPS-P to initiate an entirely new SPS-P session after current prediction window, which will inevitably lead to increased latency for the overall Videotelephony session.

For M2M communications, authors in [27], [28] designed an adaptive SPS scheme which utilizes the latest buffer occupancy information contained in the BSR to adjust the resource allocations for the next uplink SPS transmissions. Moreover, the proposed adaptive SPS scheme incorporates a channel condition measuring window, such that MCS selections that match current channel measurement can be updated, which provides certain diversity gains similar to the DS scheme. Simulation results of the proposed adaptive SPS scheme demonstrate that 85% packets have met delay budget when 50 gateways are deployed, roughly 3 times higher than the DS scheme. Moreover, the adaptive SPS scheme is shown to effectively reduce the uplink packet drop rate by almost 50% compared with the DS scheme. The proposed adaptive SPS scheme, however, does not elaborate how positions of scheduled resources are determined, making it infeasible for practical implementation. Moreover, the calculation of resource allocation does not consider new incoming uplink data after the BSR transmission. Lastly, the periodicity of adaptive SPS is over 10 ms, and hence ultra-low latency applications cannot be supported.

In [29], a soft resource reservation SPS scheme was proposed for teleoperation over mobile networks. The reservation strategy allows a UE to selectively send an SR to an eNB about whether new incoming data packets have been received at the UE buffer. Through this procedure, if the UE did not have any uplink data to send, the eNB will reschedule the reserved resources for other communications within the cell. Although this soft resource reservation SPS scheme does not appear to have any drawbacks as compared to the conventional DS scheme, it may incur a degradation

in latency performance when compared to the conventional periodic SPS scheme. Firstly, unlike the conventional SPS scheme, where no PUCCH resource is consumed after SPS establishment, the soft resource reservation SPS scheme must send uplink SR messages each time when transmission opportunity is needed. This may have a significant impact on accessible mobile connections of other communications within the cell. Secondly, an eNB needs to schedule resources for uplink transmissions 4 ms in advance and the processing time of PUCCH message is 3 ms. Hence, there will be at least a 7 ms gap between the SR packet and data packet transmission, introducing the possibility of a UE being unable to transmit in time even though it has new incoming uplink data. Results have shown that the resource reservation SPS scheme achieves an average uplink latency of 10 ms, outperforming the conventional DS scheme by 5 ms. However, no comparison has been made between the soft resource reservation SPS scheme and its preceding counterpart - the conventional SPS.

## III. PREDICTIVE SPS SCHEME FOR LATENCY-CRITICAL APPLICATIONS
### A. GENERAL WORKING PRINCIPLE OF PREDICTIVE SPS SCHEME
In the proposed Att-RNN SPS scheme, identical predictive schedulers are deployed in the MAC layers of eNB and corresponding UEs. As such, the predictive SPS scheme not only overcomes the disadvantage of unchangeable resource allocations of the conventional SPS scheme, but also preserves the advantage of zero control-signaling overhead as depicted in Fig. 5 which details the timing diagram corresponding to a newly-accepted predictive SPS session with 3 ms periodicity. More specifically, the flow diagram of the working principle of the Att-RNN SPS scheme is presented in Fig. 6, specifically illustrating $UE_j$ requesting an SPS session with an eNB. As the Att-RNN SPS scheme follows a similar procedure to the conventional SPS scheme, it is expected to be 3GPP-compliant.

### B. FEASIBLE RBP ALLOCATOR (FRA)
The FRA is designed to not only avoid resource collision among SPS-connected UEs, but also maintain the accessibility of unscheduled radio resources for DS scheme. Let $N_{\max,j}$ and $N_{n,j}$ denote the maximum number of schedulable
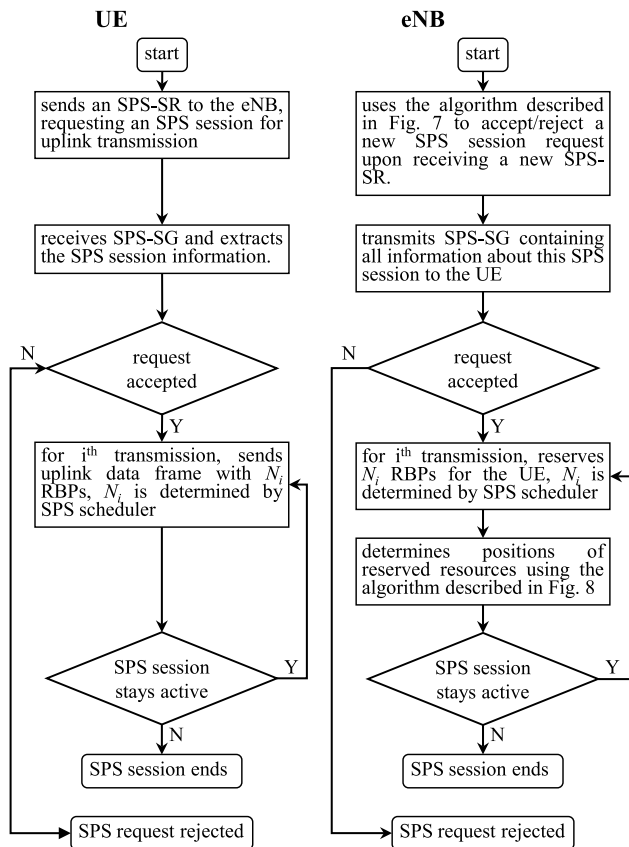
FIGURE 6. Flow diagram of working principle of predictive SPS scheme.



FIGURE 7. FRA algorithm for a new SPS session request.



FIGURE 8. FRA algorithm for an existing SPS session.

RBPs for $UE_j$ and the number of scheduled RBPs for $UE_j$ at subframe $n$. For an established SPS session between an eNB and $UE_j$, let $BR_j$, $p_j$ and $c_j$ denote the bitrate of upcoming H2R traffic, SPS periodicity and bits contained per RBP, respectively. Upon establishing a predictive SPS session, the eNB determines $N_{\max,j}$ based on:

$$N_{\max,j} = \left\lceil a_j \frac{BR_j \times p_j}{c_j} \right\rceil \qquad (1)$$

where $a_j$ is a multiplication factor determined by H2R traffic characteristics and is greater than 1.

The principle of FRA for eNB's MAC layer at subframe $n$ is described in Figs. 7 and 8. Radio resource allocation diagram for an example scenario where both $UE_1$ and $UE_2$ request predictive SPS sessions with preiodicities of 3 ms and 7 ms is illustrated in Fig. 9. Upon receiving SPS-SR from $UE_1$, the eNB runs the FRA algorithm described in Fig. 7 to accept/reject $UE_1$'s request. Specifically, if the amount of remaining contiguous RBPs is less than $N_{\max,1}$, $UE_1$'s SPS request is rejected and this situation is denoted as "Limitation of UE Admission". Otherwise, the eNB assigns $N_{\max,1}$ contiguous RBPs to $UE_1$, and $N_{\max,1}$ is the maximum amount of RBPs that $UE_1$ can access during its entire SPS session. As shown in Fig. 9, the eNB schedules and allocates $N_{n,1}$ RBPs for $UE_1$ at subframe $n$, while the remaining $(N_{\max,1} - N_{n,1})$ RBPs become available to DS-connected
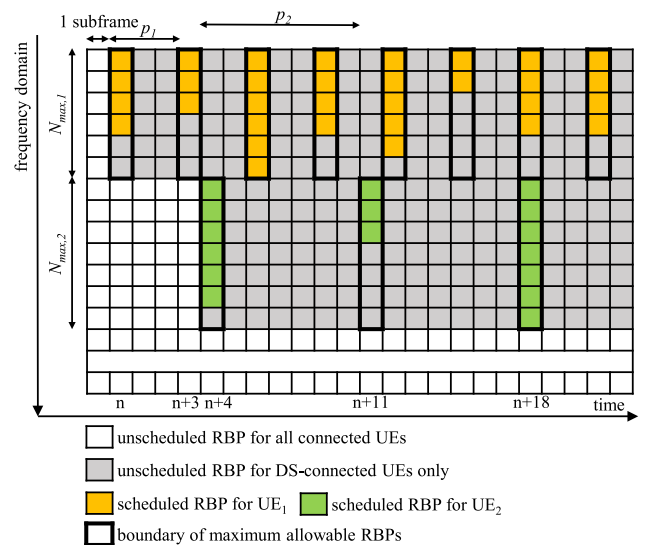


FIGURE 9. RBP allocation with FRA in an example scenario.

UEs. Subsequently after every 3 subframes, the eNB runs the FRA algorithm in Fig. 8 to schedule and allocate various RBPs to $UE_1$, as depicted in Fig. 9.

At subframe $n + 4$, the eNB establishes another predictive SPS session with $UE_2$ and it allocates RBPs to $UE_2$ every 7 subframes. Note that at subframe $n + 18$, when the eNB schedules and reserves RBPs for both the $7^{th}$ transmission of $UE_1$ and the $3^{rd}$ transmission of $UE_2$, Fig. 9 depicts that radio resource allocation of the Att-RNN SPS scheme does not cause collision, even when there exists concurrent uplink transmissions for multiple SPS-connected UEs. Furthermore, the unscheduled RBPs within $N_{max,1}$ and $N_{max,2}$ can still be allocated to DS-connected UEs. Hence the Att-RNN SPS scheme equipped with FRA can flexibly allocate radio resources among SPS-connected UEs and DS-connected UEs. As such, predictive SPS scheme achieves improved RBP utilization efficiency, as compared with conventional SPS scheme.

With regard to PUCCH usage, the Att-RNN SPS scheme differs from the soft resource reservation SPS scheme [29] in that no additional SR packet is required before data transmission, thus benefitting from a zero control-signaling overhead. For uplink latency, as explained in Section II-D, a minimum of 7 ms extra delay is introduced due to the processing of SR packets in the soft resource reservation SPS. Moreover, the predictive SPS scheme schedules RBPs based on network traffic prediction whereas the soft resource reservation SPS scheme schedules radio resources by informing eNB the amount of pending haptic data using SR.

## IV. PREDICTIVE SPS SCHEDULER BASED ON ATT-RNN
### A. OVERVIEW ON RECURRENT NEURAL NETWORK FOR SEQ2SEQ PROBLEM

Prediction of H2R traffic supported by the SPS scheme in nature is a sequence-to-sequence (seq2seq) problem. Relaying on a sequence of historical arrival data or pattern features in a time-series form, a sequence of future arrivals, either in short or long term, can be predicted. Classical seq2seq methods used in traffic prediction include statistical prediction algorithms such as AR, AR integrated moving average (ARIMA), and nonlinear AR exogenous (NARX).

More recently, advanced machine learning techniques are being explored to improve traffic prediction accuracy. In that regard, the RNN encoder-decoder model is being considered as promising solution for seq2seq problems. Both the encoder and decoder consist of an RNN unit such as LSTM and in typical language processinng scenarios, or Gated Recurrent Unit (GRU) [30]. As for the encoder, each element in the source sentences is fed into the RNN one by one to generate a context vector. The decoder behaves differently at training and prediction stages. At the training stage, the decoder's RNN is fed with the target sentence and it is initialized with the context vector from the encoder. While at the prediction stage, the decoder is fed with a start-of-sentence (sos) token to generate the first prediction output. The following outputs are then produced by feeding the sos token as well as the already-generated outputs. Applied to language processing, such RNNs show more accurate

prediction, in terms of MSE and RMSE, than classic statistical and shallow ML models, such as feed-forward ANN [30], [31].

In [32], the authors showed that the performance of RNN in seq2seq can be effectively improved by incorporating attention-based learning. For an output at a given position in the output sequence, their attention-based model selectively puts more weights (hence attention) on more relevant input elements instead of the entire input sequence. On top of adjusting the weights, attention-based learning relaxes the need of an encoder to accurately learn the features of the entire input by modifying the network structure such that encoder-extracted feature is passed to decoder at every time-step. In processing language, results show that the attention-based model significantly outperforms the conventional encoder-decoder model (RNNencdec) for both short and long sentences.

### B. ATTENTION-BASED RECURRENT NEURAL NETWORK FOR PREDICTIVE SPS SCHEME

To this end, the Attention-based encoder-decoder model was proposed as a promising machine learning method for seq2seq tasks because it is able to capture the long-term temporal correlation of time series data. This also fits the demand of $k$-step ahead prediction for the Att-RNN SPS scheme as described in [5], as both the input and output are in a sequential format.

Practiced in empirical scenarios, attention-based learning have been shown to boost existing RNNs in achieving more accurate time-series prediction performance. Qin *et al.* proposed a Dual-stage Attention-based Neural Network (DA-RNN) for predicting the target time series based on history observation of multiple driving (exogenous) series [33]. Apart from the attention operation introduced in [32], their model performs an additional attention operation among different driving series. In other words, their model performs spatial attention at the encoder to select relevant driving series, and applies temporal attention at the decoder to select the relevant element of the output of the encoder. This design was verified by using room temperature as target series and 16 other monitoring data for about 40 days as driving series. Results show that DA-RNN outperforms not only traditional statistical time series prediction methods such as ARIMA and NARX, but also conventional machine learning models such as RNNencdec and Attention-based RNN network.

Attention-based RNN has previously been used to achieve prediction in natural language processing and computer vision [34]. Similarly, a multi-level attention network for geo-sensory data (GeoMAN) was proposed in [35]. The spatial and temporal attention mechanism was utilized in selecting relevant data series, and meanwhile, the external factors of dynamic environment were jointly considered to achieve the final sensor measurement prediction. The GeoMAN was tested in water and air quality prediction base on the data collected from various measuring sensors.
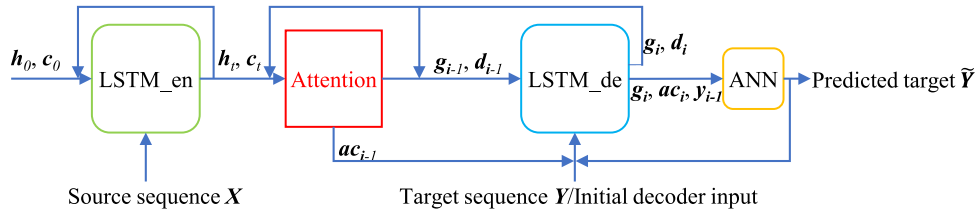
**FIGURE 10.** Recurrent neural network encoder-decoder architecture of Att-RNN prediction model.

**TABLE 2.** Notations for Attention-based prediction model.

| Symbol | Description | Symbol | Description |
|--------|-------------|--------|-------------|
| $X$ | Source sequence | $T$ | Length of source sequence |
| $n$ | Dimension of source sequence | $m$ | Dimension of hidden state of LSTM_en |
| $h_t$ | $t^{\text{th}}$ hidden state of LSTM_en | $c_t$ | $t^{\text{th}}$ cell state of LSTM_en |
| $Y$ | Target sequence | $l$ | Length of target sequence |
| $p$ | Dimension of target sequence | $q$ | Dimension of hidden state of LSTM_de |
| $g_i$ | $i^{\text{th}}$ hidden state of LSTM_de | $d_i$ | $i^{\text{th}}$ cell state of LSTM_de |
| $\tilde{Y}$ | Predicted target sequence | $ac_i$ | $i^{\text{th}}$ weighted compressed information from output of LSTM_en |

Results showed that the proposed GeoMAN achieved the best performance among considered alternatives such as ARIMA, LSTM, RNNencdec and DA-RNN.



**FIGURE 11.** Detailed structure of the LSTM unit at encoder.

### C. RECURRENT NEURAL NETWORK ENCODER-DECODER ARCHITECTURE

In this work, our proposed Att-RNN model for H2R traffic prediction employs an RNN encoder-decoder architecture which is illustrated in Fig. 10. The proposed network comprises an LSTM-based encoder and decoder, an attention block in between, and an ANN-based prediction block at the final stage. The notation used for explaining the procedure of our proposed Attention-based prediction model is listed in Table 2 and the functions of each block are explained in detail as follows:

#### 1) ENCODER USING LSTM

The encoder is used to extract compressed information from the source sequences to help the decoder generate accurate predictions. This is achieved by feeding elements of source sequences into an LSTM unit at the encoder and extracting generated hidden state ($h_t$) and cell state ($c_t$) as illustrated in Fig. 10. LSTM was first introduced by Hochreiter and Schmidhuber [36], and has a recurrent structure such that the generated information can be fed back as an input for the next computation. In this way, the long-term temporal dependencies in the source sequence can be captured. Moreover, an LSTM unit offers a non-linear transformation for the input signal without suffering from vanishing gradient problem. In comparison, conventional non-linear activation units based on sigmoid, tanh

or Rectified Linear Unit (ReLu) functions may encounter vanishing gradient in the backward propagation process, which in turn leads to non-optimized results or slow convergence speed. This issue has less effect on LSTM because it accumulates information flows over time, therefore allowing gradients to be unchanged. An improved version of LSTM contains a forget gate ($f$), an input gate ($i$), a cell gate and an output gate ($o$), by which hidden and cell states can be computed.

For a given source sequence $X = (x_1, x_2, \ldots, x_t, \ldots, x_T) \in \mathbb{R}^{n \times T}$ and target sequence $Y = (y_1, y_2, \ldots, y_i, \ldots, y_I) \in \mathbb{R}^{p \times I}$, the flow chart describing the computation process of a LSTM at encoder is presented in Fig. 11 and its calculation is shown as follows:

$$i_t = \sigma\left(W_i[h_{t-1}; x_t] + b_i\right), \tag{2}$$

$$f_t = \sigma\left(W_f[h_{t-1}; x_t] + b_f\right), \tag{3}$$

$$g_t = \tanh\left(W_g[h_{t-1}; x_t] + b_g\right), \tag{4}$$

$$o_t = \sigma\left(W_o[h_{t-1}; x_t] + b_o\right), \tag{5}$$

$$c_t = f_t \otimes c_{t-1} + i_t \otimes g_t, \tag{6}$$

$$h_t = o_t \otimes \tanh(c_t). \tag{7}$$

where $x_t \in \mathbb{R}^n$ is the $t^{\text{th}}$ element in the source sequence, $h_{t-1} \in \mathbb{R}^m$ is the hidden state from the previous LSTM calculation and $[h_{t-1}; x_t]$ represents the concatenation of these two vectors. $W_i$, $W_f$, $W_g$, $W_o \in \mathbb{R}^{m \times (n+m)}$ and $b_i$, $b_f$, $b_g$,
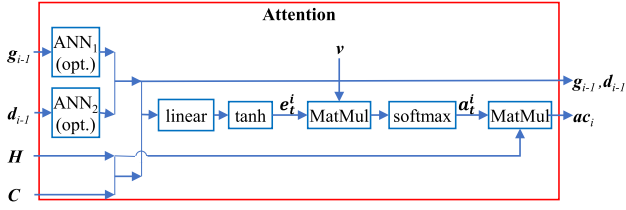
**FIGURE 12.** Additive attention mechanism implemented by the Att-RNN prediction model.



**FIGURE 13.** Data pre-processing procedure.

$b_o \in \mathbb{R}^m$ are multiplication matrices and bias terms of input, forget, cell and output gates. The symbol $\sigma$ and $\otimes$ denote sigmoid function and element-wise multiplication, respectively. Specifically, the multiplication factor $f_t$ determines the extent of forgetting the previous cell information and addition factor $o_t$ determines how much information of new cell state shall contribute to the calculation of new hidden state $h_t$. It should be noted that when feeding the first element of source sequence into the LSTM_en, the hidden and cell states denoted as $h_0$ and $c_0$ need to be initialized and provided.

### 2) ADDITIVE ATTENTION MECHANISM

A collection of hidden and cell states are available when encoding is complete and they are then processed by the attention mechanism. In our work, we employ an additive attention mechanism between the accumulated hidden and cell states from LSTM_en and the generated hidden and cell states from LSTM_de. The detailed operation is illustrated in Fig. 12, where $H = (h_1, h_2, \ldots, h_T)^\top \in \mathbb{R}^{T \times m}$ and $C = (c_1, c_2, \ldots, c_T)^\top \in \mathbb{R}^{T \times m}$ denote the accumulated hidden and cell state, respectively. The detailed mathematical derivation of the additive attention can be expressed as follows:

$$e_t^i = \tanh\left(W_a[h_t; c_t; g_{i-1}; d_{i-1}] + b_a\right), \tag{8}$$

$$a_t^i = \frac{\exp\left(v^\top e_t^i\right)}{\sum_{t=1}^{T} \exp\left(v^\top e_t^i\right)}, \tag{9}$$

$$ac_i = \sum_{t=1}^{T} a_t^i [h_t; c_t], \tag{10}$$

where $g_i \in \mathbb{R}^q$ and $d_i \in \mathbb{R}^q$ are the hidden and cell states calculated by LSTM_de. In particular, the energy score $e_t^i$ measures the importance of $t^{\text{th}}$ hidden and cell states from the encoder to the $i^{\text{th}}$ element of decoder, in order to make accurate target sequence prediction $\hat{y}_i$. The energy scores are then normalized by a softmax function [37]. Then, the weighted sum of the accumulated compressed information from LSTM_en can be calculated using a simple matrix multiplication operation. Such a procedure achieves the selection of the most relevant outputs from LSTM_en for a more accurate target prediction. $W_a \in \mathbb{R}^{q \times (2q+2m)}$, $b_a \in \mathbb{R}^q$ and $v \in \mathbb{R}^q$ are parameters of the attention mechanism which are learnt via a training process that will be described in Section V. It should be noted that at $i = 1$,
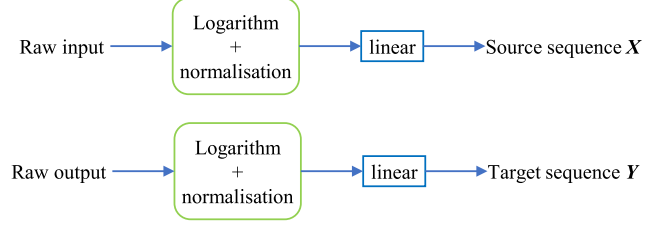
the hidden and cell states of LSTM_de are initialized with the last hidden and cell state of LSTM_en processed with feed-forward networks, which means $g_0 = \text{ANN}_1(h_T)$ and $d_0 = \text{ANN}_2(c_T)$.

### 3) DECODER USING LSTM

Finally, the target sequence prediction is realized by the LSTM_de followed by a simple feed-forward ANN network. In particular, the LSTM_de generates hidden and cell states by taking the previous states, weighted compressed information from LSTM_en and previous target prediction as inputs. The detailed calculation can be expressed as follows:

$$i_i' = \sigma\left(W_i'[g_{i-1}; z_i] + b_i'\right), \tag{11}$$

$$f_i' = \sigma\left(W_f'[g_{i-1}; z_i] + b_f'\right), \tag{12}$$

$$g_i' = \tanh\left(W_g'[g_{i-1}; z_i] + b_g'\right), \tag{13}$$

$$o_i' = \sigma\left(W_o'[g_{i-1}; z_i] + b_o'\right), \tag{14}$$

$$d_i = f_i' \otimes d_{i-1} + i_i' \otimes g_i', \tag{15}$$

$$g_i = o_i \otimes \tanh(d_i). \tag{16}$$

where $z_i \in \mathbb{R}^{m+p}$ is the concatenation of the $i^{\text{th}}$ element in the target sequence and weighted compressed information. Note that $z_i = [ac_i; y_i]$ is used in the training stage and $z_i = [ac_i; \tilde{y}_i]$ is used in the evaluation stage, respectively. $W_i'$, $W_f'$, $W_g'$, $W_o' \in \mathbb{R}^{q \times (q+m+p)}$ and $b_i'$, $b_f'$, $b_g'$, $b_o' \in \mathbb{R}^q$ are multiplication matrices and bias terms of input, forget, cell and output gates of LSTM_de. The $i^{\text{th}}$ predicted target sequence is then derived by feeding the concatenation of $y_{i-1}$, $ac_i$ and $g_i$ into the ANN. Note that the ANN in this work implements ReLU function and it comprises one hidden layer of nine nodes and one output node.

### 4) DATA PRE-PROCESSING AND LEARNABLE INITIAL DECODER INPUT

Different to [33], [35] where multi-variant data is fed directly into a prediction model, our proposed Att-RNN prediction model pre-processes raw input and output data to form proper source and target sequences. We implement this to remove potential skewness in the raw data. As illustrated in Fig. 13, the logarithm values of raw data are first computed, and then followed by normalization. Subsequently, both normalized input and output data undergo a dimension expansion process denoted as a linear block, such that potential high-dimension pattern of the input data such as
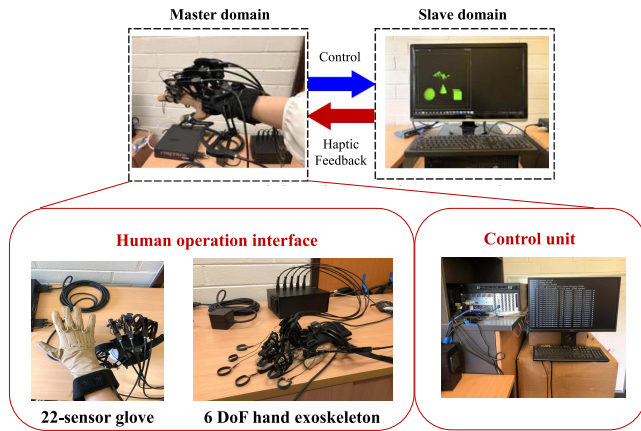
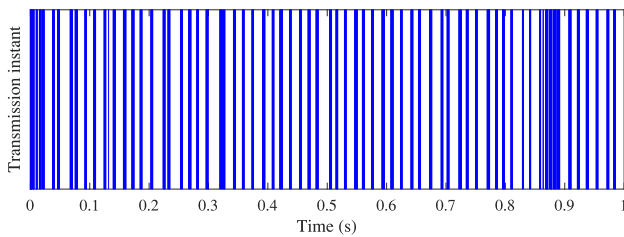**FIGURE 14.** Human-to-robot (H2R) experimental setup.



**FIGURE 15.** Experimental haptic feedback data arrivals.

**TABLE 3.** Statistics of data inter-arrival time (ms).

| Experiment indices | Statistics | | | |
|---|---|---|---|---|
| | Min | Mean | Max | Std |
| 1 | 0.0 | 1.151 | $2.530 \times 10^1$ | 3.337 |
| 2 | 0.0 | 1.450 | $1.957 \times 10^2$ | 5.533 |
| 3 | 0.0 | 1.117 | $7.639 \times 10^1$ | 3.215 |
| 4 | 0.0 | 1.146 | $3.318 \times 10^1$ | 3.273 |
| 5 | $6.000 \times 10^{-3}$ | 4.882 | $1.979 \times 10^2$ | 6.606 |

0 ms (multiple data packets were generated simultaneously) to 200 ms. Therefore, conventional periodic SPS schemes will not be able to support low latency transmission as demanded by H2R collaboration. In the following, we train the proposed Att-RNN prediction model using the empirical H2R traffic traces and validate its prediction performance.

### B. TRAINING AND TESTING DATA

As explained in Section III-A, for $UE_j$ with an SPS periodicity of $p_j$, an allocation decision $\widehat{G}_l$ for $l^{th}$ SPS transmission is determined based on available history observations, including transmitted data size $\{D_{l-k-M}, \ldots, D_{l-k}\}$, remaining buffer length $\{B_{l-k-M+1}, \ldots, B_{l-k}\}$ reported in BSRs and previously granted data size $\{G_{l-k}, \ldots, G_{l-1}\}$, where $M$ represents the selected windows size for prediction and $G_l$ can be directly calculated from $N_l$. These form the raw model of our proposed predictive scheduler. The available historic observations and prediction output for the $l^{th}$ transmission is summarized in Table 4, where $\{T_{l-k-M+1}, \ldots, T_{l-k}\}$ and $\{T_{l-k+1}, \ldots, T_l\}$ are the input and output sequences to the Att-RNN model. Note that in Table 4 "$\cdots$" stands for omission of elements and "-N/A-" denotes unavailable observations due to "Information Lagging". The elements of raw input data can be derived by:

$$T_i = D_i + B_i - B_{i-1} \qquad (17)$$

where $l - k - M + 1 \leq i \leq l - k$. It is important to note that there exists a 4 ms leading time between RBP grant and uplink data transmission and another 4 ms processing time for decoding an uplink data packet. This results in an "Information Lagging" such that for $l^{th}$ SPS transmission, a eNB has observed only the transmitted data and BSR up to $(l - k)^{th}$ order, where $k$ is determined by:

$$k = \left\lfloor \frac{7\ ms}{p_j} \right\rfloor + 1 \qquad (18)$$

Let us consider an eNB that has reached subframe $n + 26$ in Fig. 5. The eNB needs to schedule for the uplink transmission at subframe $n + 30$ (the $9^{th}$ transmission), but it had just finished processing the uplink data transmitted at subframe $n + 21$ (the $6^{th}$ transmission). Therefore, the eNB has observed only up to $D_6$, $B_6$ and $G_8$ and thus, the order of "Information Lagging" is 3.

In essence, the proposed Att-RNN SPS scheme determines the number of granted RBPs–$N_l$, based on the estimated length of the previous data buffer ($\widehat{B}_{l-1}$) and predicted size

long-term temporal dependencies could be captured. The existing attention-based RNNs as overviewed in [32] typically requires fixed size inputs at the decoder. The Att-RNN prediction model proposed in this work overcome such drawback by allowing inputs to have variable size and learnable parameters for the prediction.

## V. SUPERVISED TRAINING FOR ATTENTION-BASED PREDICTION MODEL

### A. TRACED DATA TRAFFIC FROM HUMAN-TO-ROBOT (H2R) EXPERIMENTS

To verify the proposed Att-RNN prediction model and to evaluate its performance as a predictive SPS scheme, we first collect haptic feedback traffic traces from H2R experiments which we then use in our simulations as data generated by the application layer. Fig. 14 presents our experimental platform, which comprises a human master haptic control device that is made up of 22-sensor haptic glove with position tracker and hand exoskeleton, and a robot device in a virtual environment. Based on this platform, real-time human operations and haptic response from the virtual environment can be collected. An example of the experimental haptic feedback data flow with a time window of 1 second is illustrated in Fig. 15. In addition, the statistics of inter-arrival time of haptic feedback data flows from 5 separate H2R experiments are shown in Table 3. These results reveal that the data flow generated by the haptic interaction application of the H2R experiment exhibits randomness and bursty characteristics. In particular, the inter-arrival time between two adjacent haptic feedback experimental data could range from

**TABLE 4.** Available observations and prediction target of SPS scheduler for $l^{\text{th}}$ transmission.

| Index | $l-k-M$ | $l-k-M+1$ | $\cdots$ | $l-k-2$ | $l-k-1$ | $l-k$ | $\cdots$ | $l-1$ | $l$ |
|---|---|---|---|---|---|---|---|---|---|
| New | $T_{l-k-M}$ | $T_{l-k-M+1}$ | $\cdots$ | $T_{l-k-2}$ | $T_{l-k-1}$ | $T_{l-k}$ | -N/A- | -N/A- | $\widehat{T_l}$ |
| Transmitted Data | $D_{l-k-M}$ | $D_{l-k-M+1}$ | $\cdots$ | $D_{l-k-2}$ | $D_{l-k-1}$ | $D_{l-k}$ | -N/A- | -N/A- | Not needed |
| BSR | $B_{l-k-M}$ | $B_{l-k-M+1}$ | $\cdots$ | $B_{l-k-2}$ | $B_{l-k-1}$ | $B_{l-k}$ | -N/A- | $\widehat{B}_{l-1}$ | Not needed |
| Granted Data Size | $G_{l-k-M}$ | $G_{l-k-M+1}$ | $\cdots$ | $G_{l-k-2}$ | $G_{l-k-1}$ | $G_{l-k}$ | $\cdots$ | $G_{l-1}$ | $\widehat{G_l}$ |

**TABLE 5.** Number of samples used in off-line training and testing.

| normalized traffic load | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
|---|---|---|---|---|---|---|---|---|---|---|
| 2 ms periodicity | 249650 | 249650 | 218663 | 163994 | 131190 | 109320 | 93700 | 81985 | 72872 | 65585 |
| 3 ms periodicity | 166435 | 166432 | 145780 | 109330 | 87460 | 72880 | 62470 | 54660 | 48585 | 43723 |

**TABLE 6.** NMSE of prediction for 2 ms periodicity.

| | NMSE | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Traffic load | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
| AR | 1.5741 | 1.3013 | 1.6288 | 1.6314 | 1.4575 | 1.5346 | 2.0609 | 1.5062 | 1.2798 | 1.5827 |
| ANN | 0.6467 | 0.9262 | 0.9876 | 1.0522 | 0.9450 | 0.8776 | 0.9809 | 0.9538 | 0.8592 | 0.9565 |
| DA-RNN | 1.1370 | 1.2423 | 1.1205 | 1.0258 | 1.3155 | 1.1065 | 1.1765 | 1.0096 | 1.1469 | 0.9615 |
| DP-Attn | 1.5144 | 1.7567 | 1.1584 | 0.9976 | 1.2948 | 0.9188 | 1.0215 | 1.2374 | 1.1531 | 1.1522 |
| S-Attn | 0.9433 | 0.9983 | 0.9986 | 0.9949 | 1.0015 | 0.9910 | 1.0015 | 0.9968 | 0.9894 | 1.0037 |
| **Proposed model** | **0.5794** | **0.5508** | **0.7352** | **0.7638** | **0.5982** | **0.6110** | **0.8508** | **0.6645** | **0.5626** | **0.7391** |

of new data that arrive between previous and next SPS transmission ($\widehat{T_l}$) based on:

$$N_l = \begin{cases} N_{\max}, & \text{if } \left\lceil \frac{\widehat{G_l}}{c} \right\rceil > N_{\max} \\ \left\lceil \frac{\widehat{G_l}}{c} \right\rceil, & \text{otherwise.} \end{cases} \quad (19)$$

where $\widehat{G_l} = \widehat{B}_{l-1} + \widehat{T_l}$ is the buffer occupancy predicted by the proposed Att-RNN SPS scheme and $c$ represents the amount of data contained within one RBP. To obtain the training and testing data for the Att-RNN prediction model, we first apply traced H2R arrivals obtained via the experiments described in Section V-A to the regular SPS scheme with a periodicity of 2 and 3 ms, and then extract relevant samples including the amount of transmitted data, BSR and scheduled data for each uplink transmission. Hence, the newly arrived data $T$ during an SPS period can be calculated by (17). The total size of experimental data samples applied to the regular SPS scheme with a periodicity of 2 and 3 ms for each normalized network load, is provided in Table 5. Note that 80% of total samples are used for training and the remaining for testing.

## C. PREDICTION ACCURACY DURING OFF-LINE TRAINING

In this section, the extracted time series of newly arrived data $N$ during an SPS period is used for off-line supervised training for our proposed Att-RNN prediction model. Note that the selected windows size $M$ determines the length of source sequence $X$. On one hand, a larger window size provides more information about the network traffic that in turn, could improve prediction accuracy of the proposed SPS scheduler. On the other hand, a larger window size requires the scheduler to wait for a longer time before the prediction can start, thus leading to degraded scheduling accuracy within the waiting phase. In this study, we choose 30 as the selected windows size, such that the resulting waiting time of the SPS scheduler is less than 100 ms with the selected SPS periodicity.

We then verify the proposed prediction model by examining the H2R arrival prediction accuracy measured by MAE and NMSE. The MAE of newly-arrived data is determined by $\frac{1}{n}\sum_{l=1}^{n}|T_l - T_l'|$, where $n$ denotes number of samples, and $T$ and $T'$ denote actual and predicted size of newly-arrived data, respectively. MAE measures the magnitude of overall error in forecasting but its value depends on the scale of examined data. On the other hand, the NMSE of newly-arrived data is determined by $\frac{1}{\sigma^2 n}\sum_{l=1}^{n}(T_l - T_l')^2$, where $\sigma$ is the sample standard deviation of $\{T_1, T_2, \ldots, T_n\}$ [38]. NMSE measures against the average squared error which is then normalized based on deviation of examined data group. Hence in this study, NMSE is useful for assessing the accuracy of a prediction model across data groups with different periodicities. With MAE and NMSE, we compare the accuracy of the proposed model with other alternative models including AR, simple feed-forward ANN, DA-RNN [33], Dot-Product Attention-based network (DP-Attn) [39] and Self-Attention-based network (S-Attn) [37]. It should be noted that our specific prediction task contains only one driving sequence, and therefore the spatial attention computation in the encoding stage of DA-RNN [33] is removed when it is deployed for our traffic data. The NMSE and MAE for different periodicities and varying normalized taffic load are listed in Tables 6–9, respectively.

From the results, we can clearly see that our proposed Att-RNN prediction model generates the least prediction error among all considered models and thus highlights

**TABLE 7.** NMSE of prediction for 3 ms periodicity.

| | NMSE | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Traffic load | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
| AR | 1.4034 | 1.6552 | 1.6224 | 1.5601 | 1.6317 | 1.2441 | 1.5859 | 1.4595 | 1.5830 | 1.5364 |
| ANN | 0.5611 | 1.0134 | 0.9875 | 0.8909 | 0.9701 | 0.8443 | 0.9638 | 0.8243 | 0.9382 | 0.8914 |
| DA-RNN | 0.7247 | 1.1291 | 1.0516 | 1.0363 | 1.0134 | 0.9581 | 1.0701 | 0.9681 | 0.9682 | 0.9954 |
| DP-Attn | 1.2310 | 1.4359 | 0.8925 | 1.0854 | 1.0526 | 1.0456 | 2.5316 | 1.6783 | 1.0717 | 0.9483 |
| S-Attn | 0.9506 | 1.0029 | 0.9904 | 0.9870 | 0.9928 | 0.9814 | 0.9849 | 0.9848 | 0.9844 | 0.9789 |
| **Proposed model** | **0.5211** | **0.7468** | **0.5992** | **0.6023** | **0.7516** | **0.5949** | **0.7201** | **0.5528** | **0.6547** | **0.6792** |

**TABLE 8.** MAE of prediction for 2 ms periodicity.

| | MAE | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Traffic load | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
| AR | 0.9408 | 1.0033 | 0.8467 | 0.8427 | 0.9875 | 0.9268 | 0.8420 | 0.9299 | 0.9666 | 0.8924 |
| ANN | 0.6027 | 0.8770 | 0.7188 | 0.6474 | 0.8363 | 0.7599 | 0.5945 | 0.8168 | 0.7445 | 0.6664 |
| DA-RNN | 0.8125 | 0.9466 | 0.8342 | 0.6872 | 0.9349 | 0.8210 | 0.7139 | 0.8382 | 0.8932 | 0.6578 |
| DP-Attn | 0.9890 | 1.0472 | 0.8223 | 0.7426 | 0.9518 | 0.7372 | 0.6198 | 0.9041 | 0.8952 | 0.7649 |
| S-Attn | 0.8447 | 0.9260 | 0.7106 | 0.7561 | 0.9118 | 0.8238 | 0.5977 | 0.8483 | 0.8585 | 0.6737 |
| **Proposed model** | **0.5428** | **0.5518** | **0.5416** | **0.5790** | **0.5869** | **0.5630** | **0.5548** | **0.6043** | **0.5623** | **0.5786** |

**TABLE 9.** MAE of prediction for 3 ms periodicity.

| | MAE | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| Traffic load | 0.1 | 0.2 | 0.3 | 0.4 | 0.5 | 0.6 | 0.7 | 0.8 | 0.9 | 1.0 |
| AR | 0.9859 | 0.8530 | 0.9232 | 0.9446 | 0.8697 | 0.9514 | 0.8466 | 0.9312 | 0.8839 | 0.9123 |
| ANN | 0.5845 | 0.7904 | 0.8154 | 0.7802 | 0.7133 | 0.7695 | 0.5762 | 0.7547 | 0.6353 | 0.7502 |
| DA-RNN | 0.6619 | 0.8042 | 0.8332 | 0.7746 | 0.7041 | 0.8423 | 0.6910 | 0.8304 | 0.6649 | 0.7710 |
| DP-Attn | 0.8629 | 0.9728 | 0.7563 | 0.7299 | 0.6704 | 0.8349 | 1.1894 | 0.9992 | 0.7341 | 0.7940 |
| S-Attn | 0.9063 | 0.7198 | 0.8541 | 0.8178 | 0.7671 | 0.8467 | 0.6348 | 0.8555 | 0.6535 | 0.8215 |
| **Proposed model** | **0.5302** | **0.5463** | **0.5588** | **0.5503** | **0.5766** | **0.6044** | **0.5039** | **0.5496** | **0.5639** | **0.6055** |

its potential in improving latency performance and RBP allocation efficiency, as will be evident in the next section.

## VI. SIMULATION RESULTS AND DISCUSSION

We verify the performance enhancement of the SPS scheme arising from implementing our trained Attention-based prediction model by simulating an LTE network that supports both low-latency H2R UEs and latency-tolerant non-H2R UEs. The implementation is based on feeding the traced H2R arrivals into an LTE simulation program–SimuLTE [40], which is a system-level simulator developed on the well-known discrete event simulation platform OMNeT++ and it is widely used to simulate LTE and LTE-Advanced systems. Prediction via our proposed Att-RNN SPS scheduler is performed online and in real-time while both H2R and non-H2R traffic is generated and delivered across the LTE network. The LTE network parameters and traffic characteristics are summarized in Table 10. In Table 10, a normalized traffic load of 1 represents the maximum aggregate data offered to all UEs and this uplink capacity is chosen to be 36 Mbps over 500 radio resource block pairs. In comparison, the amount of received data–data plane utilization, is dependent on the number of connected UEs and the packet drop rate.

### A. LATENCY PERFORMANCE

Figure 16 compares the average uplink MAC layer latency performance of H2R UEs incurred by conventional DS,

**TABLE 10.** Simulation parameters.

| Parameter | Value | Notation |
|---|---|---|
| Simulation time | 100 s | |
| Number of Cells | 1 | |
| Number of H2R UEs | 10 | $M$ |
| Number of non-H2R UEs | 20 | |
| Bandwidth | 200 MHz | |
| Uplink RBPs | 500 | |
| Theoretical uplink capacity | 36.0 Mbps | $C$ |
| Normalized traffic load | 0.1–1 | $L$ |
| Non-H2R traffic load ratio | 60% | $\alpha$ |
| Non-H2R traffic | Pareto distribution | |
| Non-H2R packet payload | Uniform(20,1500) B | |
| H2R traffic load ratio | 40% | $\beta$ |
| Haptic packet payload | 12 B | |
| Scheduling algorithm for DS scheme | Proportional fair | |
| SR periodicity | 1 ms | |
| Tested SPS period | 2, 3 ms | $p$ |
| UE MAC buffer size | 1 MB | |
| Bytes per RBP | 9 | $c$ |
| LTE multiplexing mode | FDD | |
| Tx mode | Single antenna | |
| Maximum UE Tx power | 398 mW | |
| Pathloss model | ITU-R, Rural Macro | |

conventional SPS, adaptive SPS in [27], [28], AR-based predictive SPS in [5], and our proposed Att-RNN SPS. Results shows that the achievable uplink latency of the
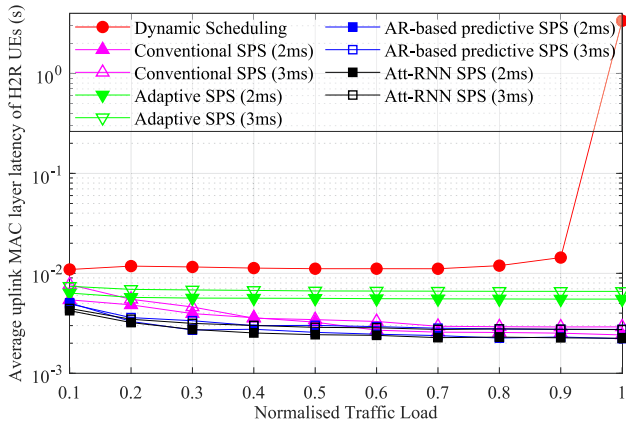
**FIGURE 16.** Average uplink MAC layer latency of H2R UEs as a function of normalized traffic load.
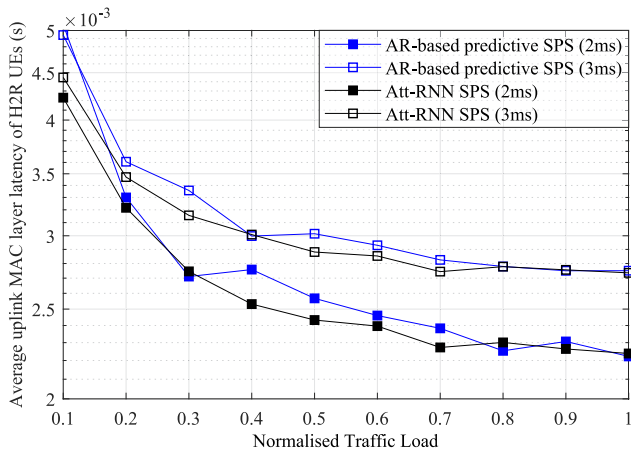


**FIGURE 17.** Detailed latency performance comparison as a function of normalized traffic load.



**FIGURE 18.** Overall data plane utilization as a function of normalized traffic load.



**FIGURE 19.** Number of connected H2R UEs as a function of normalized traffic load.

conventional DS protocol is around 10 ms latency when the normalized traffic load is under 0.9, and this rises drastically to 3 seconds when the network is near saturation. On the other hand, SPS schemes provide relatively stable latency performance at under 7 ms. It is worth noting that the resulting average uplink latency of Adaptive SPS is higher than that of other SPS schemes, reflecting the limitation of Adaptive SPS to predict bursty traffic. A detailed comparison of uplink latencies between the AR-based predictive SPS and Att-RNN SPS schemes is presented in Fig. 17. Note that the average uplink MAC layer latency of both AR-based predictive SPS and the proposed scheme is seen to increase as traffic loads are reduced. This is predominantly due to higher prediction errors, which will be explained in more detail in Section VI-C.

## B. DATA PLANE UTILIZATION AND NUMBER OF CONNECTED H2R UES

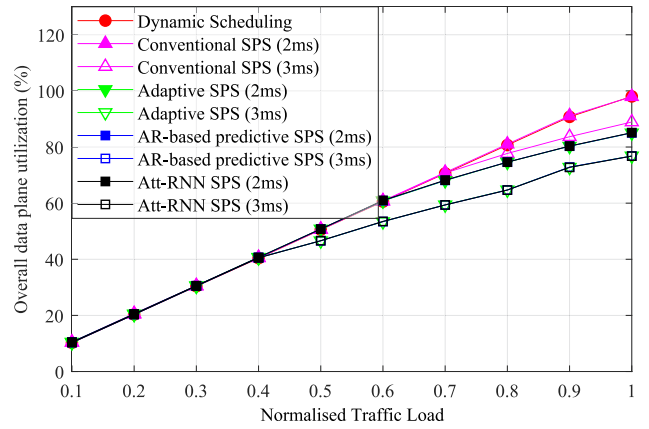Fig. 18 compares the data plane utilization representing overall throughput and Fig. 19 shows the number of connected (admitted) H2R UEs incurred by the investigated uplink MAC schemes. From Fig. 18, the utilization curves of all schemes increase linearly with the traffic load up to 0.4. However, the conventional SPS, Adaptive SPS, AR-based predictive SPS and Att-RNN SPS produce less utilization when the traffic load is above 0.5. This results from the decreased number of connected H2R UEs shown in Fig. 19. As described in Step 2 of Section III-A, the admission control implemented by SPS scheme could lead to "Limitation of UE Admission" when ($N_{max}$) increases with normalized traffic load, reflecting the known trade-off between throughput and latency in wireless communications [41].

## C. SCHEDULING AND PREDICTION PERFORMANCE

The ability of different SPS schemes to allocate sufficient RBPs for uplink transmission is measured by the NMSE of RBP allocation, whereby a lower NMSE indicates less deviation between actual and ideal allocations and hence a higher prediction accuracy and allocation efficiency. The simulation results in Fig. 20 show that Att-RNN SPS provides the best allocation accuracy, which explains the improvement of latency performance presented in Fig. 16.

The allocation efficiency reflected by buffer occupancy during simulation is shown in Fig. 21. In general, the
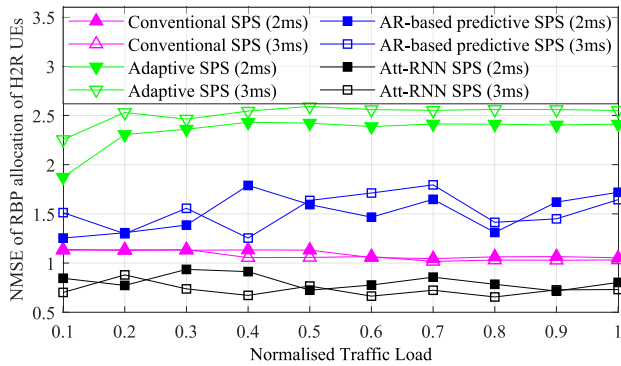
**FIGURE 20.** NMSE of RBP allocation for H2R UEs as a function of total normalized traffic load.
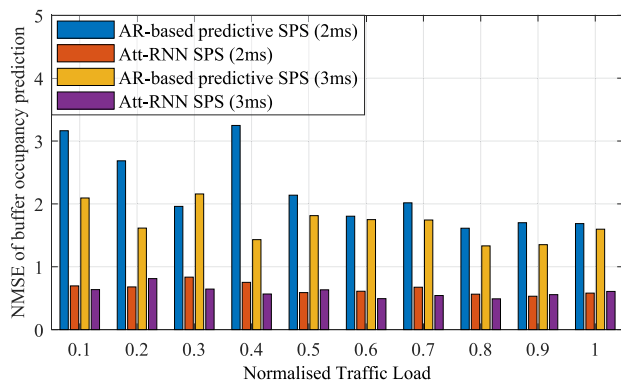


**FIGURE 21.** NMSE of predicted buffer occupancy as a function of total normalized traffic load.

prediction performed by Att-RNN SPS is more accurate than AR-based predictive SPS. However, when the operation condition denoted as (load, periodicity (ms)) is (0.2, 3), (0.3, 2), (1.0, 2) and (1.0, 3), the prediction accuracies of these two models are comparable, and this conforms to comparable latency performances as depicted in Fig. 17. Furthermore, the results from Figs. 17 and 21 indicate an increasing trend of prediction error as traffic load is decreased. A possible reason for this can be attributed to the increasing sparsity of bursty data arrivals when traffic is lightly loaded, in which case temporal dependencies is harder to obtain due to the limited size of raw input data.

## VII. CONCLUSION

In this article, we presented a novel scheduler for the predictive SPS scheme based on an Attention-based Recurrent Neural Network for H2R traffic prediction, in order to achieve enhanced latency performance for H2R collaboration over LTE-based industrial networks. Specifically, we interpreted the "Information Lagging" issue arising from predictive SPS schemes into a seq2seq problem and design the Att-RNN model for this task. In addition, we proposed a new feasible RBP allocator that ensures collision-free radio resource allocation and enhances allocation efficiency by retaining accessibility of unused resources. Further, we

studied the characteristics of traced data traffic from real H2R haptic feedback experiments. Simulation results showed that for all the prediction methods that are compared, their H2R arrival prediction accuracy may degrade under light traffic loads, when the temporal dependency is relatively weak among arrivals in a time-series form. Nonetheless, the Att-RNN SPS outperforms the existing schemes in latency performance, achieving < 5ms latency over all network load range. In summary, our proposed Att-RNN SPS improves the SPS scheme adopted by LTE and NR standards developed for 5G, and hence better supports key applications such as haptic telerobot and haptic AR/VR in advanced industrial deployments of the future.

## REFERENCES

[1] M. Chen, Y. Miao, Y. Hao, and K. Hwang, "Narrow band Internet of Things," *IEEE Access*, vol. 5, pp. 20557–20577, 2017.

[2] O. Liberg, M. Sundberg, Y.-P. E. Wang, J. Bergman, J. Sachs, and G. Wikström, "Chapter 5—LTE-M," in *Cellular Internet of Things*, 2nd ed. Cambridge, MA, USA: Academic Press, 2020, pp. 155–254.

[3] E. Wong, M. Pubudini Imali Dias, and L. Ruan, "Predictive resource allocation for tactile Internet capable passive optical LANs," *J. Lightw. Technol.*, vol. 35, no. 13, pp. 2629–2641, Jul. 2017.

[4] L. Ruan, M. P. I. Dias, M. Maier, and E. Wong, "Understanding the traffic causality for low-latency human-to-machine applications," *IEEE Netw. Lett.*, vol. 1, no. 3, pp. 128–131, Sep. 2019.

[5] Y. Feng, A. Nirmalathas, and E. Wong, "A predictive semi-persistent scheduling scheme for low-latency applications in LTE and NR networks," in *Proc. IEEE Int. Conf. Commun. (ICC)*, May 2019, pp. 1–6.

[6] L. Ruan, M. P. I. Dias, and E. Wong, "Achieving low-latency human-to-machine (H2M) applications: An understanding of H2M traffic for AI-facilitated bandwidth allocation," *IEEE Internet Things J.*, vol. 8, no. 1, pp. 626–635, Jan. 2021.

[7] S. Mondal, L. Ruan, M. Maier, D. Larrabeiti, G. Das, and E. Wong, "Enabling remote human-to-machine applications with AI-enhanced servers over access networks," *IEEE Open J. Commun. Soc.*, vol. 1, pp. 889–899, Jul. 2020.

[8] C. M. Bishop, *Pattern Recognition and Machine Learning* (Information Science and Statistics). New York, NY, USA: Springer, 2006.

[9] L. Ruan, M. P. I. Dias, and E. Wong, "Enhancing latency performance through intelligent bandwidth allocation decisions: A survey and comparative study of machine learning techniques," *IEEE/OSA J. Opt. Commun. Netw.*, vol. 12, no. 4, pp. B20–B32, Apr. 2020.

[10] S. J. Russell and P. Norvig, *Artificial Intelligence: A Modern Approach* (Pearson Custom Library). Harlow, U.K.: Pearson, 2014.

[11] Z. Tan, Y. Li, Q. Li, Z. Zhang, Z. Li, and S. Lu, "Supporting mobile VR in LTE networks: How close are we?" *Proc. ACM Meas. Anal. Comput. Syst.*, vol. 2, no. 1, p. 8, Apr. 2018.

[12] E. Dahlman, S. Parkvall, and J. Sköld, "Chapter 5—Physical transmission resources," in *4G LTE-Advanced Pro and The Road to 5G*, 3rd ed., E. Dahlman, S. Parkvall, and J. Sköld, Eds. London, U.K.: Academic, 2016, pp. 1–5.

[13] A. S. Lioumpas and A. Alexiou, "Uplink scheduling for machine-to-machine communications in LTE-based cellular systems," in *Proc. IEEE GLOBECOM Workshops (GC Wkshps)*, Dec. 2011, pp. 353–357.

[14] A. Aijaz, "Toward human-in-the-loop mobile networks: A radio resource allocation perspective on haptic communications," *IEEE Trans. Wireless Commun.*, vol. 17, no. 7, pp. 4493–4508, Jul. 2018.

[15] S. A. AlQahtani, "Delay-aware resource allocation for M2M communications over LTE-A networks," *Arabian J. Sci. Eng.*, vol. 44, no. 4, pp. 3639–3653, Apr. 2019.

[16] A. M. Maia, D. Vieira, M. F. de Castro, and Y. Ghamri-Doudane, "A fair QoS-aware dynamic LTE scheduler for machine-to-machine communication," *Comput. Commun.*, vols. 89–90, pp. 75–86, Sep. 2016.

[17] "Contention based uplink transmission," Huawei, HiSilicon, Shenzhen, China, Rep. R2-154191, Oct. 2015.

[18] "Impacts of contention based uplink in RAN2," Ericsson, ST-Ericsson, Stockholm, Sweden, Rep. R2-100125, Jan. 2010.

[19] L. Zhang and J. Ma, *Grant-Free Multiple Access Scheme*. Cham, Switzerland: Springer, 2019, pp. 515–533.

[20] S. Andreev *et al.*, "Efficient small data access for machine-type communications in LTE," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Jun. 2013, pp. 3569–3574.

[21] C. Wang, Y. Chen, Y. Wu, and L. Zhang, "Performance evaluation of grant-free transmission for uplink URLLC services," in *Proc. IEEE 85th Veh. Technol. Conf. (VTC Spring)*, Jun. 2017, pp. 1–6.

[22] D. Jiang, H. Wang, E. Malkamaki, and E. Tuomaala, "Principle and performance of semi-persistent scheduling for VoIP in LTE system," in *Proc. Int. Conf. Wireless Commun. Netw. Mobile Comput.*, Sep. 2007, pp. 2861–2864.

[23] J. Seo and V. C. M. Leung, "Performance modeling and stability of semi-persistent scheduling with initial random access in LTE," *IEEE Trans. Wireless Commun.*, vol. 11, no. 12, pp. 4446–4456, Dec. 2012.

[24] "Study on latency reduction techniques for LTE," 3GPP, Sophia Antipolis, France, Rep. TR 36.881, Jun. 2016.

[25] T. Jacobsen *et al.*, "System level analysis of uplink grant-free transmission for URLLC," in *Proc. IEEE Globecom Workshops (GC Wkshps)*, Dec. 2017, pp. 1–6.

[26] J. T. S. Avocanh, M. Abdennebi, J. Ben-Othman, and G. Piro, "A semi-persistent scheduling scheme for videotelephony traffics in the uplink of LTE networks," in *Proc. 17th ACM Int. Conf. Model. Anal. Simulat. Wireless Mobile Syst.*, 2014, pp. 321–325.

[27] N. Afrin, J. Brown, and J. Y. Khan, "An adaptive buffer based semi-persistent scheduling scheme for machine-to-machine communications over LTE," in *Proc. 8th Int. Conf. Next Gener. Mobile Apps Services Technol.*, Sep. 2014, pp. 260–265.

[28] N. Afrin, J. Brown, and J. Y. Khan, "Design of a buffer and channel adaptive LTE semi-persistent scheduler for M2M communications," in *Proc. IEEE Int. Conf. Commun. (ICC)*, Jun. 2015, pp. 5821–5826.

[29] M. Condoluci, T. Mahmoodi, E. Steinbach, and M. Dohler, "Soft resource reservation for low-delayed teleoperation over mobile networks," *IEEE Access*, vol. 5, pp. 10445–10455, May 2017.

[30] I. Sutskever, O. Vinyals, and Q. V. Le, "Sequence to sequence learning with neural networks," 2014, *arxiv:1409.3215*.

[31] K. Cho, B. van Merrienboer, Ç. Gülçehre, F. Bougares, H. Schwenk, and Y. Bengio, "Learning phrase representations using RNN encoder-decoder for statistical machine translation," 2014, *arxiv:1406.1078*.

[32] D. Bahdanau, K. Cho, and Y. Bengio, "Neural machine translation by jointly learning to align and translate," 2015, *arXiv:1409.0473*.

[33] Y. Qin, D. Song, H. Chen, W. Cheng, G. Jiang, and G. W. Cottrell, "A dual-stage attention-based recurrent neural network for time series prediction," in *Proc. 26th Int. Joint Conf. Artif. Intell.*, 2017, pp. 2627–2633. [Online]. Available: https://doi.org/10.24963/ijcai.2017/366

[34] B. Wang, K. Liu, and J. Zhao, "Inner attention based recurrent neural networks for answer selection," in *Proc. 54th Annu. Meeting Assoc. Comput. Linguist.*, vol. 3, 2016, pp. 1288–1297.

[35] Y. Liang, S. Ke, J. Zhang, X. Yi, and Y. Zheng, "GeoMAN: Multi-level attention networks for geo-sensory time series prediction," in *Proc. 27th Int. Joint Conf. Artif. Intell.*, 2018, pp. 3428–3434. [Online]. Available: https://doi.org/10.24963/ijcai.2018/476

[36] S. Hochreiter and J. Schmidhuber, "Long short-term memory," *Neural Comput.*, vol. 9, no. 8, pp. 1735–1780, 1997. [Online]. Available: https://doi.org/10.1162/neco.1997.9.8.1735

[37] A. Vaswani *et al.*, "Attention is all you need," 2017, *arXiv:1706.03762*.

[38] R. Adhikari and R. K. Agrawal, "An introductory study on time series modeling and forecasting," 2013, *arXiv:1302.6613*.

[39] M. Luong, H. Pham, and C. D. Manning, "Effective approaches to attention-based neural machine translation," 2015, *arXiv:1508.04025*.

[40] A. Virdis *et al.*, "Simulating LTE/LTE-advanced networks with simuLTE," in *Advances in Intelligent Systems and Computing*, vol. 402. Cham, Switzerland: Springer, 2015, pp. 83–105.

[41] A. E. Gamal, J. Mammen, B. Prabhakar, and D. Shah, "Throughput-delay trade-off in wireless networks," in *Proc. IEEE INFOCOM*, vol. 1, Mar. 2004, p. 475.

**YE FENG** (Student Member, IEEE) received the Ph.D. degree from the University of Melbourne, Melbourne, VIC, Australia, in 2020. He is working for developing embedded IoT device and his research interests also include network latency reduction, teletraffic engineering, and machine learning.

**LIHUA RUAN** (Member, IEEE) received the Ph.D. degree from the University of Melbourne, Melbourne, VIC, Australia, in 2020. She is currently a Postdoctoral Fellow of the Chinese University of Hong Kong, Shenzhen. She has coauthored more than 30 papers and served as a reviewer for multiple IEEE Communication Society journals. She research interests include low-latency communications and networking, and machine learning-assisted low-latency applications. She was a recipient of the Melbourne University's John Collier Scholarship and the John Melvin Scholarship for the Best Ph.D. Thesis in Engineering and IT.

**AMPALAVANAPILLAI NIRMALATHAS** (Senior Member, IEEE) received the Ph.D. degree in electrical and electronic engineering from the University of Melbourne. He is currently a Professor and the Deputy Dean (Research) of the Faculty of Engineering and Information Technology. His current research interests include energy efficient telecommunications, access networks, optical-wireless network integration, mobile-access edge computing, photonic reservoir computing, Internet of Things, and broadband wireless systems and devices.

**ELAINE WONG** (Senior Member, IEEE) received the Ph.D. degree in electrical and electronic engineering from the University of Melbourne, Melbourne, VIC, Australia, in 2002, where she is currently a Professor and an Associate Dean (Diversity and Inclusion) of the Faculty of Engineering and Information Technology. She has coauthored more than 170 journal and conference publications. Her research interests include energy-efficient optical and wireless networks, optical-wireless integration, broadband applications of vertical-cavity surface-emitting lasers, wireless sensor body area networks, and emerging optical and wireless technologies for 6G. She has served on the editorial board for the *Journal of Lightwave Technology* and the *Journal of Optical Communications and Networking*. She is a Fellow of Optica (formerly Optical Society).