

Robust Communication and Computation using Deep Learning via Joint Uncertainty Injection

Robert-Jeron Reifert*, Hayssam Dahrouj†, Alaa Alameer Ahmad‡, Haris Gacanin§ and Aydin Sezgin*

*Ruhr University Bochum, Germany, †University of Sharjah, United Arab Emirates,

‡Cariad SE, Wolfsburg, Germany, §RWTH Aachen University, Germany.

Abstract—The convergence of communication and computation, along with the integration of machine learning and artificial intelligence, stand as key empowering pillars for the sixth-generation of communication systems (6G). This paper considers a network of one base station serving a number of devices simultaneously using spatial multiplexing. The paper then presents an innovative deep learning-based approach to simultaneously manage the transmit and computing powers, alongside computation allocation, amidst uncertainties in both channel and computing states information. More specifically, the paper aims at proposing a robust solution that minimizes the worst-case delay across the served devices subject to computation and power constraints. The paper uses a deep neural network (DNN)-based solution that maps estimated channels and computation requirements to optimized resource allocations. During training, uncertainty samples are injected after the DNN output to jointly account for both communication and computation estimation errors. The DNN is then trained via backpropagation using the robust utility, thus implicitly learning the uncertainty distributions. Our results validate the enhanced robust delay performance of the joint uncertainty injection versus the classical DNN approach, especially in high channel and computational uncertainty regimes.

I. INTRODUCTION

The global mobile network data traffic is projected to surpass 550 exabytes per month over the next few years [1]. This surge will involve substantial amounts of remote data connectivity and processing, thereby highlighting the need for an integrated communication and computation platform. Concurrently, the rise of delay-sensitive applications such as extended reality (XR) and digital twins is bound to further spearhead the need for advanced resource management techniques toward the development of the sixth-generation of communication systems (6G) [2]. Notably, the International Telecommunication Union (ITU) envisions future 6G systems to natively incorporate machine learning (ML) and artificial intelligence (AI) as part of their design, deployment, and operation [3]. With ubiquitous intelligence, 6G would specifically be able to support the convergence of communication and computing, which is essential for XR video processing, as well as autonomous driving [4]. Deep learning techniques, particularly, promise to deliver 6G networks solutions with low

complexity, robustness, and near-optimality experience [4], e.g., provisioned wireless resource management [5], and robust optimization [6]. Along such lines, this paper proposes, and evaluates the benefit of, using a deep neural network (DNN) to generate robust solutions for a complex joint communication and computation delay minimization problem.

The problem addressed is related to the joint communication and computation resource management framework which is considered in the recent literature, e.g., [7]–[13]. To this end, reference [7] adopts a rate-splitting multiple access scheme and reference [8] focuses on XR-empowered systems, yet, both [7] and [8] assume perfect knowledge of channels and computing requirements, which is challenging to realize in practice. From a robust resource allocation perspective, addressing the absence of channel state information has been considered in several works, e.g., [9], [10]. From a robust computation perspective, reference [11] offers a framework to guarantee task deadlines under processing uncertainty. While references [9]–[11] address only one type of uncertainty, i.e., either communication uncertainty or computing uncertainty, reference [12] extends the scope to a joint perspective on computational intelligence in 6G. Further, [13] proposes a robust edge computing application, where video streams are recorded and processed at the network edge before being streamed to the devices.

Despite their numerical merits, [7]–[13] adopt classical numerical optimization approaches. The synergy between 6G and AI/ML, however, nowadays provides powerful alternatives that promise to enhance seamless network performance, particularly in complex tasks related to advanced radio resource management problems [5], [14]. For instance, the optimization of communication and computation resources through ML-based algorithms is explored in [15] using deep reinforcement learning, and in [16] using federated edge learning. As the available training data is never perfect, [6] proposes using an uncertainty injection method that employs a deep learning-based architecture to maximize the network minimum-rate subject to imperfect channel state information. The perspective of jointly robust communication and computation, especially the one that captures their respective uncertainties, remains, however, largely unexplored in the literature, and so this paper addresses such issues using a joint uncertainty injection scheme.

The paper proposes a robust deep learning framework to address the joint communication and computation

This work has been submitted to the IEEE for possible publication. Copyright may be transferred without notice, after which this version may no longer be accessible.

This work was supported in part by the German Federal Ministry of Education and Research (BMBF) in the course of the 6GEM Research Hub under grant 16KISK037.

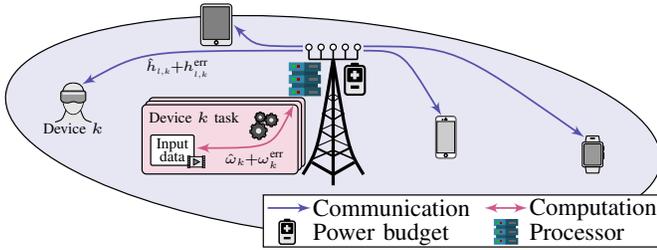


Fig. 1: Network setup with multi antenna BS and 4 devices.

resource management problem in the presence of channel and computing uncertainty. We consider a setup of one base station (BS) serving a number of devices using spatial multiplexing, i.e., via beamforming. Our focus is on minimizing the worst-case delay across all served devices by optimizing the transmit and computing powers, as well as the computing resources. Our proposed DNN-based approach takes estimated channels and computing requirements as inputs and produces optimized transmit and computing powers, along with computations allocations. The robustness of our solution is enhanced through a joint uncertainty injection framework, where the DNN's output undergoes further processing by injecting channel and computing error samples directly. Under such a framework, the DNN parameters are updated by leveraging the γ -th quantile value of worst-case delays under uncertainty injection utilizing backpropagation. Our results highlight the enhanced robust delay performance of the joint uncertainty injection, especially for high parameters uncertainties and power-limited regimes. Explicit comparisons with two different uncertainty injection schemes and a state-of-the-art DNN-based solution show how the proposed algorithm is highly flexible, generalizes the uncertainty injection schemes, and provides enhanced robust delays.

II. SYSTEM MODEL AND PROBLEM FORMULATION

In this work, we consider a single BS with L antennas and embedded computing ability serving K single-antenna devices. We denote the antenna set by $\mathcal{L} = \{1, \dots, L\}$, and the device set by $\mathcal{K} = \{1, \dots, K\}$. An example of the considered model is illustrated in Fig. 1. In fact, in many 6G-key applications, e.g., image and video processing or digital twinning, devices depend on strong remote processing capabilities at the network edge [12]. In such cases, each device experiences a delay consisting of two parts, namely the communication delay, i.e., the time loss on the physical communication medium, and the computation delay, i.e., the time for computing the device task.

A. Channel Setup and Communication Delay

Let $h_{l,k} \in \mathbb{C}$ be the channel coefficient from the BS's antenna l to device k , where $\mathbf{h}_k = [h_{1,k}, \dots, h_{L,k}]^T$ is the aggregate channel vector of device k . We assume that the BS has only access to an estimate of the channels. That is, the estimated channel from antenna l to device k is denoted by $\hat{h}_{l,k} \in \mathbb{C}$, with k 's aggregate estimated channel vector denoted by $\hat{\mathbf{h}}_k = [\hat{h}_{1,k}, \dots, \hat{h}_{L,k}]^T$. In practice, the channel uncertainty distribution is not known [10], [17]. Yet, as channel

estimations can be sampled in large amounts, a data-driven approach is able to *learn* the channel uncertainty implicitly. Without loss of generality we model the channel coefficients as

$$h_{l,k} = \hat{h}_{l,k} + h_{l,k}^{\text{err}}, \quad (1)$$

where $h_{l,k}^{\text{err}} \sim \mathcal{CN}(0, \sigma_h^2)$ is the channel estimation error.

The beamforming vector of device k 's signal $\mathbf{v}_k \in \mathbb{C}^L$ is defined by $\mathbf{v}_k = [v_{1,k}, \dots, v_{L,k}]^T$, where $\|\mathbf{v}_k\|_2 = 1$, adopting the regularized zero-forcing (RZF) strategy [6]. We also let p_k^{tx} be the BS's transmit power allocated to serve device k . The transmit power vector is then denoted by $\mathbf{p}^{\text{tx}} = [p_1, \dots, p_K]$.

The devices messages are encoded into $\xi_k, \forall k \in \mathcal{K}$, where $\xi_k \sim \mathcal{CN}(0, 1)$ are independent and identically distributed, circularly symmetric Gaussian random variables. Each device k receives its own desired signal, interference from other signals, and noise. The received signal at device k becomes

$$y_k = \mathbf{h}_k^H \mathbf{v}_k \sqrt{p_k^{\text{tx}}} \xi_k + \sum_{j \in \mathcal{K} \setminus \{k\}} \mathbf{h}_k^H \mathbf{v}_j \sqrt{p_j^{\text{tx}}} \xi_j + n_k, \quad (2)$$

where n_k is the additive white Gaussian noise, with zero mean and variance σ^2 . Based on (2), the achievable rate of device k is given by

$$r_k = W \log_2 \left(1 + \frac{|\mathbf{h}_k^H \mathbf{v}_k|^2 p_k^{\text{tx}}}{\sum_{j \in \mathcal{K} \setminus \{k\}} |\mathbf{h}_k^H \mathbf{v}_j|^2 p_j^{\text{tx}} + \sigma^2} \right), \quad (3)$$

where W is the system bandwidth. Device k 's communication delay is then calculated as D_k^{out}/r_k , where D_k^{out} is the output data size of the k -th device task computation. According to (3), r_k is a function of the channel values, which highlights the impact of the channel uncertainty on the communication delay.

B. Computing Framework and Computation Delay

From the computation perspective, device k requests to process a task, which consists of an input data size D_k^{in} , and a computation intensity ω_k , i.e., the required computation cycles per bit of input data is denoted by ω_k , for a total of $\omega_k D_k^{\text{in}}$ computation cycles for device k 's task. However, such a computational requirement ω_k may not be perfectly known at the BS, as only an estimation $\hat{\omega}_k$ is available. Specifically, computing environments are highly dynamic, i.e., they vary depending on system load and concurrent processes, and systems dependency, i.e., the inter-relation of input data characteristics. Various statistical models for ω_k exist in the literature, including uniform, normal, and Gamma distributions [18]. This paper addresses the computation uncertainty using a data-driven fashion, where we sample large amounts of uncertainty realizations and train the DNN accordingly. The computation intensity is modeled as

$$\omega_k = \hat{\omega}_k + \omega_k^{\text{err}}, \quad (4)$$

where $\omega_k^{\text{err}} \sim \mathcal{N}(0, \sigma_\omega^2)$ is the computing estimation error [18].

With an estimation of the required computation cycles for device k 's task, i.e., $\hat{\omega}_k D_k^{\text{in}}$, the BS allocates f_k^{co} in cycles/s to compute k 's task. The overall vector of computation allocations is denoted by $\mathbf{f}^{\text{co}} = [f_1^{\text{co}}, \dots, f_K^{\text{co}}]^T$. Due to the finite

computation resources at the BS, the maximum computation capacity constraint is given by

$$\sum_{k \in \mathcal{K}} f_k^{\text{co}} \leq F^{\text{max}}, \quad (5)$$

where F^{max} is the BS's maximum computation capacity. Due to the intertwined nature of communication and computation resources, especially in terms of the allocated power levels, we consider p^{co} as the total power allocated by the BS for the computation of all tasks. The relation between the computation power p^{co} and computation allocation f_k is given by

$$p^{\text{co}} = \tau \left(\sum_{k \in \mathcal{K}} f_k^{\text{co}} \right)^\mu, \quad (6)$$

where τ and μ are constants of the CPU model [8].

By further considering the communication delay discussed earlier, device k experiences a total delay t_k consisting of a computation and a communication delay as follows

$$t_k = \frac{\omega_k D_k^{\text{in}}}{f_k^{\text{co}}} + \frac{D_k^{\text{out}}}{r_k}. \quad (7)$$

In order to achieve fair service among all network devices, the paper aims at minimizing the worst-case delay t_{max} among all devices, i.e., $t_{\text{max}} = \max_{k \in \mathcal{K}} \{t_k\}$, computed via the ground-truth channels \mathbf{h}_k and computation intensities ω_k , $\forall k \in \mathcal{K}$.

C. Joint Communication and Computation Power

In this work, we consider the coupling of communication and computation power consumption (i.e., p^{tx} and p^{co} , respectively) at the BS, which leads to the overall power constraint

$$\sum_{k \in \mathcal{K}} p_k^{\text{tx}} + p^{\text{co}} \leq P^{\text{max}}, \quad (8)$$

with total power budgeted P^{max} . Equation (8) emphasizes the intertwined communication and computation powers, especially underlining the need for a careful selection of joint resource allocation under limited resources.

D. Robust Delay Utility

Based on (7), one can infer that the overall delay t_k of device k is a function of the channel vectors \mathbf{h}_k and the computation intensities ω_k . Both \mathbf{h}_k and ω_k follow some unknown probability distribution, such that the probability distributions of t_k and t_{max} are also unknown. Specifically, the BS only has access to the estimated $\hat{\mathbf{h}}_k$ and $\hat{\omega}_k$, $\forall k \in \mathcal{K}$.

To best account for the channel and computing uncertainty, the paper utilizes a DNN with inputs $\hat{\mathbf{h}}_k$, $\hat{\omega}_k$, $\forall k \in \mathcal{K}$, and inject samples of uncertainty realizations of both channels and computing requirements into the training phase just after the DNN output [6]. That is, for a given resource allocation, we inject estimation errors to impact the objective value, and afterwards train the network according to the new (robust) objective. Such uncertainty injection supplies a multitude of potential worst-case delays t_{max} for different realizations. This paper aims to minimize the γ -th quantile of the worst-case delay distribution, denoted by t_{max}^γ , where

$$\Pr \left[t_{\text{max}} > t_{\text{max}}^\gamma | \hat{\mathbf{h}}_k, \hat{\omega}_k, \forall k \in \mathcal{K} \right] \leq \gamma, \quad (9)$$

where t_{max}^γ can be interpreted as the value for which only a fraction γ of all realizations have higher worst-case delays.

E. Problem Formulation

The paper now aims at minimizing the robust worst-case delay t_{max}^γ among all devices. The corresponding optimization problem is formulated as follows:

$$\begin{aligned} \min_{\mathbf{p}^{\text{tx}}, p^{\text{co}}, \mathbf{f}^{\text{co}}} \quad & t_{\text{max}}^\gamma \\ \text{s.t.} \quad & (1), (4), (5), (8), (9), \\ & t_{\text{max}} \geq t_k, \quad \forall k \in \mathcal{K}, \end{aligned} \quad (10a)$$

where the optimization runs over the variables \mathbf{p}^{tx} , p^{co} , and \mathbf{f}^{co} , the computation capacity constraint is given in (5), and the power constraint is given in (8). Constraint (9) ensures the solution's robustness in terms of the γ -th quantile. Constraint (10a) expresses the worst delay among all devices as t_{max} . The intricacy of problem (10) stems from the mathematical coupling between the communication and computation variables and constraints, as well as from the unknown distribution of the worst-case delay t_{max} given the channel and computing states information uncertainties. The paper next proposes a data-driven solution approach based on a deep learning framework with joint uncertainty injection.

III. DEEP LEARNING VIA JOINT UNCERTAINTY INJECTION

To tackle the intricate optimization problem (10), we harness the promising abilities of deep learning in terms of automatic learning, feature representation, and versatility. More specifically, we first fix the beamforming vectors \mathbf{v}_k using: $\mathbf{V} = \hat{\mathbf{H}}(\hat{\mathbf{H}}\hat{\mathbf{H}}^H + \alpha\mathbf{I})^{-1}$, where $\mathbf{V} = [\mathbf{v}_1, \dots, \mathbf{v}_K]$, $\hat{\mathbf{H}} = [\hat{\mathbf{h}}_1, \dots, \hat{\mathbf{h}}_K]$, α is the RZF factor, and \mathbf{I} is a $K \times K$ identity matrix. Each \mathbf{v}_k is normalized to unit norm. We then utilize a DNN to directly map problem (10)'s parameters to a resource allocation solution, which is then further processed by injecting uncertainties for optimizing a robust objective. In a typical DNN, neurons are organized into layers, including input, hidden, and output layers. These layers are densely interconnected, with each connection between neurons having an associated weight that is continuously adjusted during the training process. Additionally, activation functions are applied to the neurons within these layers, facilitating nonlinear transformations of their inputs. In general, depth (number of layers) and width (number of neurons in each layer) of the DNN are key factors that influence its capacity to learn and represent data.

In the context of this paper, we propose a DNN architecture, which takes the estimated channels and estimated computing requirements as inputs, i.e., $\hat{\mathbf{h}}_k$ and $\hat{\omega}_k$, $\forall k \in \mathcal{K}$. Thus, the input layer consists of $2LK + K$ neurons, which represent the $2LK$ real and imaginary channel coefficients and the K distinct devices' tasks. After the input layer, we attach $\mathcal{F}_{\text{depth}}$ fully-connected hidden layers with $\mathcal{F}_{\text{width}}$ neurons each, and Rectified Linear Unit (ReLU) activation functions. The output layer consists of $2K + 1$ neurons, where the first $K + 1$ entries represent the K devices' allocated power and the BS's computation power, and the remaining K entries represent the computation capacities. The DNN then returns the output vector $\mathbf{x} = [x_1, \dots, x_{2K+1}]^T$, representing the optimized resource allocation \mathbf{p}^{tx} , p^{co} , and \mathbf{f}^{co} , after some mathematical manipulations. We let \mathcal{F}_Θ denote the proposed DNN, such that problem (10) is solved by

$$\mathbf{x} = \mathcal{F}_\Theta(\hat{\mathbf{h}}, \hat{\omega}), \quad (11)$$

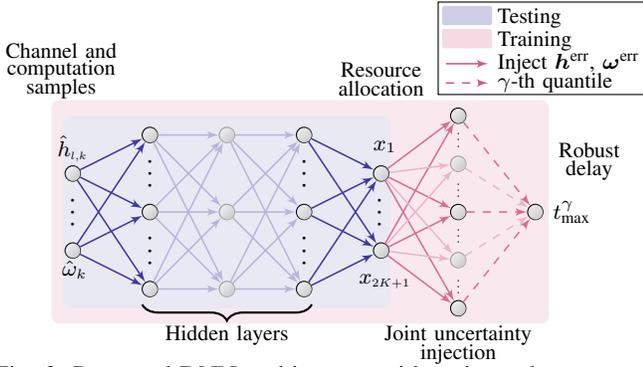


Fig. 2: Proposed DNN architecture with estimated parameters input, resource allocation output, and uncertainty injection.

where $\hat{\mathbf{h}} = [\hat{h}_1, \dots, \hat{h}_K]^T$ and $\hat{\boldsymbol{\omega}} = [\hat{\omega}_1, \dots, \hat{\omega}_K]^T$. Note that Θ refers to the DNN model parameters, e.g., weights, biases, etc.¹ Fig. 2 illustrates a representation of the proposed DNN-based solution. In the context of minimizing the robust worst-case delay among devices, we next explore how to map the DNN output to the resource allocation, how to ensure the robustness of the solution, i.e., in terms of the γ -th quantile, and how to train the DNN according to the robust objective.

A. DNN Output to Resource Allocation Mapping

A major complication of the joint communication and computation resource management problem (10) is the interdependent computation capacity (5) and power (8) constraints. We, thus, first use the softmax function² on parts of the DNN output to compute the transmit and computation powers as

$$[(\mathbf{p}^{\text{tx}})^T, \mathbf{p}^{\text{co}}]^T = P^{\text{max}} \cdot \text{softmax}([x_1, \dots, x_{K+1}]^T). \quad (12)$$

By using the softmax function, we obtain $\sum_{i=1}^{K+1} x_i = 1$, which ensures the full utilization of the complete power budget at the BS. Next, we transform the computation power into cycles using (6), so as to obtain the maximum number of available computation cycles limited by the power allocation, denoted by F^{pow} . Based on (6), F^{pow} can be written as:

$$F^{\text{pow}} = \sqrt[\mu]{\frac{p^{\text{co}}}{\tau}}. \quad (13)$$

At this point, there are two limiting factors for the computation allocation (i.e., \mathbf{f}^{co}), namely, the BS's maximum computation capacity F^{max} in (5), and the above power-limited available computation cycles F^{pow} . We then use the smallest among the two limiting factors, i.e., $\min(F^{\text{max}}, F^{\text{pow}})$. To determine the computation allocation \mathbf{f}^{co} , we use the softmax function on the remaining part of the output layer as follows

$$\mathbf{f}^{\text{co}} = \min(F^{\text{max}}, F^{\text{pow}}) \cdot \text{softmax}([x_{K+2}, \dots, x_{2K+1}]^T), \quad (14)$$

which utilizes all available computation cycles, either limited by the capacity or by the power.

¹In equation (11), the output \mathbf{x} is a function of the DNN parameters Θ . Such DNN \mathcal{F}_{Θ} and its parameters Θ are sequentially updated via backpropagation as discussed in Section III-C later.

²The softmax function of an input vector \mathbf{x} is an output vector, the i -th entry of which is $e^{x_i} / (\sum_{j=1}^n e^{x_j})$, where n is the dimension of \mathbf{x} .

B. Joint Uncertainty Injection Scheme

Given the resource allocation \mathbf{p}^{tx} , \mathbf{p}^{co} , and \mathbf{f}^{co} found in (12) and (14), one can readily compute a worst-case delay for *one given* channel and computing estimation. However, problem (10) aims at providing a robust t_{max}^{γ} accounting for the uncertainties $\mathbf{h}^{\text{err}} = [h_{1,1}^{\text{err}}, \dots, h_{L,K}^{\text{err}}]^T$ and $\boldsymbol{\omega}^{\text{err}} = [\omega_1^{\text{err}}, \dots, \omega_K^{\text{err}}]^T$, given in (1) and (4), respectively. Thus, we next explain the uncertainty injection approach [6] and present the joint uncertainty injection scheme adopted in this paper's context.

The basic idea of uncertainty injection is to estimate the statistical objective numerically, using a large number of realizations or samples of the uncertainty. Thereafter, such realizations are *injected* after the DNN output to compute a large number of objective values. These objectives are then used to obtain an empirical estimate of the robust objective. Such estimate is then used to derive gradients and update the DNN parameters via backpropagation [6].

In the context of this paper, however, the problem is subject to both channel and computing state information uncertainties. The adopted objective is also a function of the worst-case delay, as presented in (9), (10). Therefore, we sample N realizations of the channel and computing uncertainties $\mathbf{h}_n^{\text{err}}$, $\boldsymbol{\omega}_n^{\text{err}}$, $n \in \{1, \dots, N\}$. These samples are used to compute N worst-case delays $t_{\text{max},n}$, $\forall n \in \{1, \dots, N\}$. With these delay-values derived from the uncertainty samples, we sort all N delays and select the $(\gamma \cdot N)$ -th highest value, which forms an empirical estimate of the γ -th quantile worst-case delay t_{max}^{γ} , also referred to as $\hat{t}_{\text{max}}^{\gamma}$.

C. Unsupervised Learning via Backpropagation

An integral part of the proposed solution is the DNN training procedure, especially in accounting for the uncertainties to ensure a robust solution. More specifically, by taking the $(\gamma \cdot N)$ -th highest value (i.e., $\hat{t}_{\text{max}}^{\gamma}$) of all computed worst-case delays $t_{\text{max},n}$, $\forall n \in \{1, \dots, N\}$, gradients are derived from $\hat{t}_{\text{max}}^{\gamma}$ to update the DNN parameters. To be more specific, we use the gradient $\partial \hat{t}_{\text{max}}^{\gamma} / \partial \Theta$, which is computed using Pytorch [19]. In that way, the DNN is trained in an unsupervised fashion, implicitly learning the distribution of channel and computing uncertainties. The numerical merit of our proposed method is discussed in the simulations section of the paper, as shown next.

IV. SIMULATION RESULTS AND DISCUSSION

In the following, we numerically evaluate the performance of the proposed DNN-based solution with joint uncertainty injection. The network consists of one BS with $L = 8$ antennas serving $K = 6$ devices. Other parameters are summarized in Tab. I. The noise is set to -75 dBm/Hz. We also assume Rayleigh fading $h_{l,k} \sim \mathcal{CN}(0, 1)$, $\forall (l, k) \in (\mathcal{L}, \mathcal{K})$, utilize minimum mean-square estimation for the estimated channels $\hat{\mathbf{h}}$, and set the beamforming vectors \mathbf{v}_k , $\forall k \in \mathcal{K}$, according to RZF [6]. As the beamforming vectors are determined a priori, we slightly modify the DNN's input layer, and only input effective channels, i.e., $\hat{h}_{i,j}^{\text{eff}} = \hat{\mathbf{h}}_i^H \mathbf{v}_j$, $\forall i, j \in \mathcal{K}$, thereby reducing the input layer to $2K^2 + K$ nodes. The ground-truth computation intensities are drawn according to a Gamma distribution, i.e., $\omega_k \sim \Gamma(2, 200)$, $\forall k \in \mathcal{K}$ [11].

General	Communication	Computation
$P^{\max} = 38$ dBm	$W = 1$ MHz	$F^{\max} = 4.6 \cdot 10^9$ cycles/s
$N = 1000$ samples	$D_k^{\text{out}} = 7.5 \cdot 10^4$ bits	$D_k^{\text{in}} = 5 \cdot 10^4$ bits
$\gamma = 0.05$	$\alpha = 0.2$ RZF factor	$\tau = 10^{-28}$, $\mu = 3$

Tab. I: Simulation parameters (unless otherwise mentioned).

As for the DNN parameters and training procedure, we use $\mathcal{F}_{\text{depth}} = 10$ fully-connected hidden layers with $\mathcal{F}_{\text{width}} = 400$ neurons each, ReLU activation, train 500 epochs with 50 minibatches, each consisting of 1000 independent network realizations. Validation and test sizes are set to 2000 each, where each realization experiences uncertainty injection to obtain the robust objectives.

To benchmark the proposed joint uncertainty injection (referred to as *Joint UI*), we utilize two modified schemes, *Comm UI*, a method which only accounts for communication channel uncertainty, and *Comp UI*, a method which only accounts for uncertainty in the computing requirements. Further, we consider the *Regular DNN*, a scheme that does not account for the uncertainties at all, and is trained according to the estimated inputs only.

A. Impacts of the Uncertainties

In the first simulation set, we observe the impact of both channel and computing requirement uncertainties on the robust delay \hat{t}_{\max}^{γ} . To this end, Fig. 3 depicts the robust delay versus the channel estimation error variance σ_h^2 , while maintaining a fixed computation estimation error variance of $\sigma_{\omega}^2 = 6400$. Initially, it is noteworthy that all uncertainty injection schemes demonstrate improvements over the baseline *Regular DNN*, with gains ranging between 8 and 20 ms. Notably, the *Joint UI* scheme consistently achieves the lowest robust delay across all levels of channel uncertainty variances. Fig. 3 highlights the generalizability and flexibility of *Joint UI* in different uncertainty scenarios. Specifically, lower values of σ_h^2 mark a region where computing uncertainty dominates, while channel uncertainty has a lesser impact. Here, *Comp UI* closely rivals *Joint UI*, while *Comm UI* performs relatively worse. In contrast, for higher channel uncertainties (e.g., $\sigma_h^2 = 0.06$), *Comm UI* and *Joint UI* exhibit similar performances, while *Comp UI* shows a delay increase of over 6 ms. A noteworthy comparison between *Regular DNN* and *Joint UI* can be drawn at $\hat{t}_{\max}^{\gamma} = 80$ ms. *Regular DNN* achieves $\hat{t}_{\max}^{\gamma} = 80$ ms at approximately $\sigma_h^2 = 0.01$, whereas *Joint UI* achieves lower delays, down to around $\sigma_h^2 = 0.053$. Consequently, when the robust delay is considered as a measure of delay guarantee, it becomes apparent that *Joint UI* can guarantee better delays than *Regular DNN*, even with significantly increased uncertainty variance. These findings underscore two significant insights: The importance of uncertainty injection in minimizing robust delays in general, and the effectiveness of the proposed *Joint UI* scheme across various uncertainty scenarios in particular.

Fig. 4 reaffirms the advantages of *Joint UI* by showing the robust delay as a function of the computation estimation error variance σ_{ω}^2 , with $\sigma_h^2 = 0.035$. In Fig. 4, the high computation estimation error region ($\sigma_{\omega}^2 > 6400$) offers further insights. By comparing the slope of *Regular DNN* and *Comm UI* versus the slope of *Comp UI* and *Joint UI*, it is observed that both

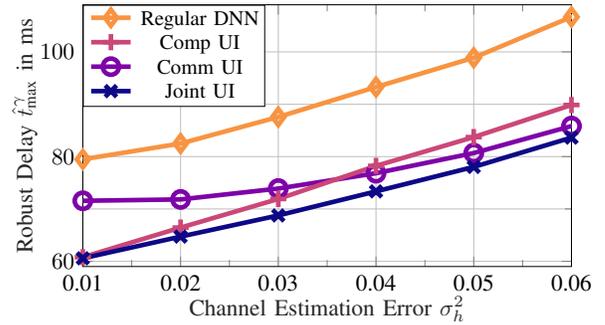


Fig. 3: Robust worst-case delay \hat{t}_{\max}^{γ} as a function of the channel uncertainty, i.e., σ_h^2 .

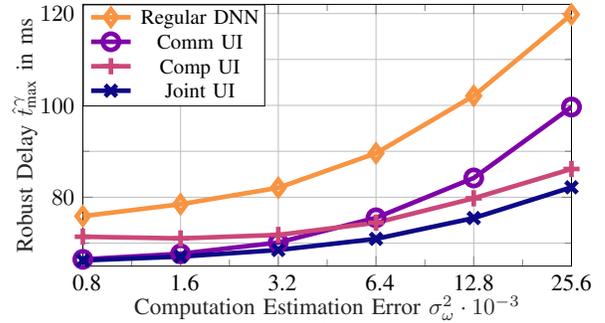


Fig. 4: Robust worst-case delay \hat{t}_{\max}^{γ} as a function of computing uncertainty, i.e., σ_{ω}^2 .

Comp UI and *Joint UI* exhibit a relatively small increase in performance degradation in response to high computing uncertainty. In fact, at $\sigma_{\omega}^2 = 25600$, the proposed *Joint UI* achieves a gain of approximately 37 ms over *Regular DNN*, 17 ms over *Comm UI*, and 4 ms over *Comp UI*.

All the examined uncertainty injection schemes in Fig. 3 and Fig. 4 notably improve the network robust delay performance, albeit with varying degrees of effectiveness for *Comp UI* and *Comm UI* across different uncertainty values. *Joint UI*, particularly, showcases superior flexibility, generalizes both *Comp UI* and *Comm UI*, and achieves significantly reduced robust delays. These results highlight *Joint UI*'s potential in spearheading the prospective joint communication and computation resource management schemes in future 6G networks.

B. Joint Power Constraint

With $\sigma_h^2 = 0.05$ and $\sigma_{\omega}^2 = 6400$, Fig. 5 illustrates the robust delay \hat{t}_{\max}^{γ} as a function of P^{\max} . Across all schemes, \hat{t}_{\max}^{γ} decreases with increasing power budget. With more power available at the BS, resources can be allocated more effectively, leading to improved delay performance. Notably, the proposed *Joint UI*, *Comm UI*, and *Comp UI* schemes consistently outperform the *Regular DNN* scheme, which neglects channel and computing uncertainties. For instance, in the low-power region ($P^{\max} < 32$ dBm), *Joint UI* surpasses *Regular DNN* by over 40 ms, while this gain diminishes to around 15 ms in the high-power region. Particularly in the low-power regime, *Comm UI* exhibits significant performance degradation compared to *Joint UI*, as it fails to consider computing uncertainty. In contrast, *Comp UI* closely tracks *Joint UI* with a nearly constant gap, highlighting the dominance of

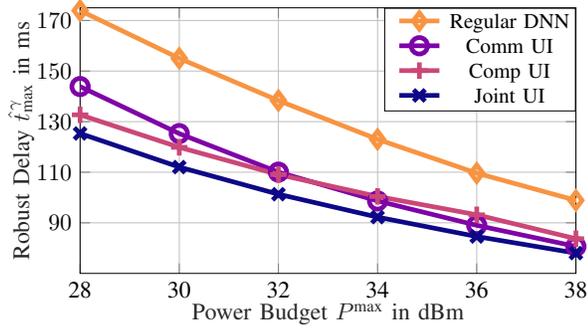


Fig. 5: Robust worst-case delay \hat{t}_{\max}^{γ} vs. power budget P^{\max} .

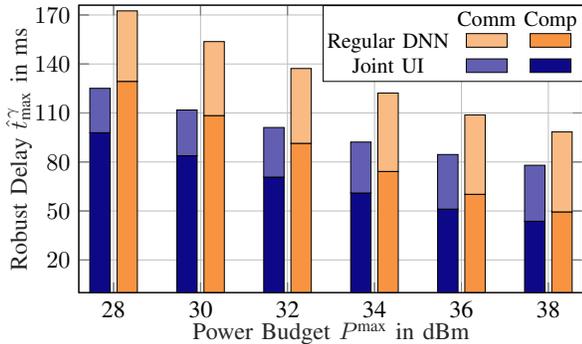


Fig. 6: Robust worst-case communication (*Comm*) and computation (*Comp*) delay vs. power budget P^{\max} .

computation delay in low-power regimes. Once again, Fig. 5 emphasizes the necessity of accounting for both channel and computing uncertainties to optimize delay performance, particularly in the more constrained low-power regime.

To best illustrate the communication and computation delays of *Joint UI* and *Regular DNN*, Fig. 6 depicts both the worst-case communication (D_k^{out}/r_k) and computation ($\omega_k D_k^{\text{in}}/f_k^{\text{co}}$) delays against the power budget P^{\max} . Firstly, from a computing perspective, Fig. 6 illustrates how an increasing power budget leads to reduced worst-case computation delays. On the other hand, the impact of P^{\max} on the communication delays is less pronounced. This discrepancy is reasonable, given that computation delays are, for this parameter choice, significantly higher than communication delays, thereby dominating the overall delay metric \hat{t}_{\max}^{γ} . The imbalance between communication and computation delays is particularly evident in the low-power region. For instance, for *Joint UI* at 28 dBm, the computation delay is 98 ms, whereas the communication delay is only 27 ms. However, as more power becomes available (as depicted on the right-hand side of Fig. 6), the *Joint UI* scheme allocates additional resources favoring computation delay reduction. Consequently, computation delays approach communication delays as P^{\max} increases.

V. CONCLUSION

Automating the decision making via ML algorithms is among the most prominent key enablers of future seamless networks. This paper introduces a novel ML-aided approach to address the convergence of communication and computation in 6G systems, while taking uncertainties in channel estimation and computing requirements into consideration. Our proposed

robust solution effectively manages transmit and computing powers, alongside computation allocation, using a DNN-based approach. By injecting uncertainty samples during training and optimizing resource allocations based on robust utility, our approach minimizes the worst-case delay among multiple devices. Our findings demonstrate the superior performance of joint uncertainty injection for various power budgets and uncertainty levels, thereby unveiling the proposed solution potential at promoting robust communication and computation resource management schemes in future 6G systems.

REFERENCES

- [1] "Ericsson mobility report November 2023," Ericsson, Tech. Rep., Jun. 2023. [Online]. Available: <https://www.ericsson.com/en/reports-and-papers/mobility-report/reports/november-2023>
- [2] L. U. Khan, W. Saad, D. Niyato, Z. Han, and C. S. Hong, "Digital-twin-enabled 6G: Vision, architectural trends, and future directions," *IEEE Commun. Mag.*, vol. 60, no. 1, pp. 74–80, 2022.
- [3] *IMT-2030 Framework: Recommendation ITU-R M.2160-0*, International Telecommunication Union (ITU), ITU-R Working Party 5D, Nov 2023.
- [4] N. Kato, B. Mao, F. Tang, Y. Kawamoto, and J. Liu, "Ten challenges in advancing machine learning technologies toward 6G," *IEEE Wirel. Commun.*, vol. 27, no. 3, pp. 96–103, 2020.
- [5] A. G. Onalan and S. Coleri, "Deep learning based low complexity relay selection for wireless powered cooperative communication networks," in *Proc. BalkanCom*, 2023, pp. 1–6.
- [6] W. Cui and W. Yu, "Uncertainty injection: A deep learning method for robust optimization," *IEEE Trans. Wirel. Commun.*, vol. 22, no. 11, pp. 7201–7213, 2023.
- [7] R.-J. Reifert, H. Dahrouj, A. A. Ahmad, A. Sezgin, T. Y. Al-Naffouri, B. Shihada, and M.-S. Alouini, "Rate-splitting and common message decoding in hybrid cloud/mobile edge computing networks," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 5, pp. 1566–1583, 2023.
- [8] R.-J. Reifert, H. Dahrouj, and A. Sezgin, "Extended reality via cooperative NOMA in hybrid cloud/mobile-edge computing networks," *IEEE Internet Things J.*, pp. 1–1, 2023.
- [9] M. Hasan, E. Hossain, and D. I. Kim, "Resource allocation under channel uncertainties for relay-aided device-to-device communication underlaying LTE-A cellular networks," *IEEE Trans. Wirel. Commun.*, vol. 13, no. 4, pp. 2322–2338, 2014.
- [10] A. J. Anandkumar, A. Anandkumar, S. Lambotharan, and J. A. Chambers, "Robust rate maximization game under bounded channel uncertainty," *IEEE Trans. Veh. Technol.*, vol. 60, no. 9, pp. 4471–4486, 2011.
- [11] S. Li, C. Li, Y. Huang, B. A. Jalaian, Y. T. Hou, and W. Lou, "Enhancing resilience in mobile edge computing under processing uncertainty," *IEEE J. Sel. Areas Commun.*, vol. 41, no. 3, pp. 659–674, 2023.
- [12] B. Ji, Y. Wang, K. Song, C. Li, H. Wen, V. G. Menon, and S. Mumtaz, "A survey of computational intelligence for 6G: Key technologies, applications and trends," *IEEE Trans. Industr. Inform.*, vol. 17, no. 10, pp. 7145–7154, 2021.
- [13] Z. Liu, C. Zhan, Y. Cui, C. Wu, and H. Hu, "Robust edge computing in UAV systems via scalable computing and cooperative computing," *IEEE Wirel. Commun.*, vol. 28, no. 5, pp. 36–42, 2021.
- [14] S. Tamoor-ul Hassan, S. Samarakoon, M. Bennis, and M. Latva-aho, "Latency-aware radio resource optimization in learning-based cloud-aided small cell wireless networks," *IEEE Trans. Green Commun. Netw.*, vol. 8, no. 1, pp. 542–558, 2024.
- [15] C. Huang, G. Chen, P. Xiao, Y. Xiao, Z. Han, and J. A. Chambers, "Joint offloading and resource allocation for hybrid cloud and edge computing in SAGINs: A decision assisted hybrid action space deep reinforcement learning approach," *arXiv preprint arXiv:2401.01140*, 2024.
- [16] X. Mo and J. Xu, "Energy-efficient federated edge learning with joint communication and computation design," *J. Commun. Inf. Netw.*, vol. 6, no. 2, pp. 110–124, 2021.
- [17] A. M. Ahmed, A. A. Ahmad, S. Fortunati, A. Sezgin, M. S. Greco, and F. Gini, "A reinforcement learning based approach for multitarget detection in massive MIMO radar," *IEEE Trans. Aerosp. Electron. Syst.*, vol. 57, no. 5, pp. 2622–2636, 2021.
- [18] J. R. Lorch and A. J. Smith, "Improving dynamic voltage scaling algorithms with PACE," *ACM SIGMETRICS Performance Evaluation Review*, vol. 29, no. 1, pp. 50–61, 2001.

- [19] A. Paszke, S. Gross, F. Massa, A. Lerer, J. Bradbury, G. Chanan, T. Killeen, Z. Lin, N. Gimelshein, L. Antiga *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *Advances in neural information processing systems*, vol. 32, 2019.