

Task-Oriented Over-the-Air Computation for Multi-Device Edge AI

Dingzhu Wen, Xiang Jiao, Peixi Liu, Guangxu Zhu, Yuanming Shi, and Kaibin
Huang

Abstract

Edge inference refers to the use of artificial intelligent (AI) models at the network edge to provide mobile devices inference services and thereby enable intelligent services such as auto-driving and Metaverse towards 6G. However, departing from the classic paradigm of data-centric designs, the 6G networks for supporting edge AI features task-oriented techniques that focus on effective and efficient execution of AI task. Targeting end-to-end system performance, such techniques are sophisticated as they aim to seamlessly integrate sensing (data acquisition), communication (data transmission), and computation (data processing). Aligned with the paradigm shift, a task-oriented over-the-air computation (AirComp) scheme is proposed in this paper for multi-device split-inference system. In the considered system, local feature vectors, which are extracted from the real-time noisy sensory data on devices, are aggregated over-the-air by exploiting the waveform superposition in a multiuser channel. Then the aggregated features as received at a server are fed into an inference model with the result used for decision making or control of actuators. To design inference-oriented AirComp, the transmit precoders at edge devices and receive beamforming at edge server are jointly optimized to rein in the aggregation error and maximize the inference accuracy. The problem is made tractable by measuring the inference accuracy using a surrogate metric called discriminant gain, which measures the discernibility of two

D. Wen is with Network Intelligence Center, School of Information Science and Technology, ShanghaiTech University, Shanghai, China (e-mail: wendzh@shanghaitech.edu.cn), and was with Shenzhen Research Institute of Big Data, Shenzhen, China. (Corresponding author: G. Zhu)

Xiang Jiao and P. Liu are with State Key Laboratory of Advanced Optical Communication Systems and Networks, School of Electronics, Peking University, China, and Shenzhen Research Institute of Big Data, Shenzhen, China (e-mail: jiaoxiang@stu.pku.edu.cn, liupeixi@pku.edu.cn).

G. Zhu is with Shenzhen Research Institute of Big Data, Shenzhen, China (e-mail: gxzhu@sribd.cn).

Y. Shi is with Network Intelligence Center, School of Information Science and Technology, ShanghaiTech University, Shanghai, China (e-mail: shiym@shanghaitech.edu.cn).

K. Huang is with Department of Electrical and Electronic Engineering, The University of Hong Kong, Hong Kong (e-mail: huangkb@eee.hku.hk).

object classes in the application of object/event classification. It is discovered that the conventional AirComp beamforming design for minimizing the mean square error in generic AirComp with respect to the noiseless case may not lead to the optimal classification accuracy. The reason is due to the overlooking of the fact that feature dimensions have different sensitivity towards aggregation errors and are thus of different importance levels for classification. This issue is addressed in this work via a new task-oriented AirComp scheme designed by directly maximizing the derived discriminant gain. However, the resultant problem of joint transmit precoding and receive beamforming is nonconvex and difficult to solve due to the complicated form of discriminant gain and the coupling between the control variables. We overcome the difficulty using the successive convex approximation. The performance gain of the proposed task-oriented scheme over the conventional schemes is verified by extensive experiments targeting the application of human motion recognition.

I. INTRODUCTION

One main function of 6G networks is to provide artificial intelligent (AI) services such as auto-driving, eHealth, and Metaverse, at the network edge [1], [2], [3], [4], [5]. Existing data-driven services, i.e., transmitting multi-media data like voice, text, image, and video, focus on throughput maximization where only the communication process is considered. On the contrary, AI services at the network edge are goal oriented and aim to achieve the required accuracy and latency for completing a specific task (see, e.g., [6], [7], [8]). The task execution typically involves the tight integration of three processes, i.e., *sensing* for real-time data acquisition, *communication* for data transmission, and *computation* for decision making [9], [10]. Then to efficiently support edge-intelligence services drives an ongoing paradigm shift in wireless technologies from the traditional data-centric design toward the task-oriented design for 6G [11], [12], [13], [14]. On the other hand, edge-intelligence services rely on the deployment of trained AI models at the network edge for decision making and timely response to a dynamic environment [15], [16], [17]. This gives rise to an emerging research area, called edge inference. The design of edge inference faces two main challenges. On one hand, the task-oriented techniques for efficient edge inference must integrate sensing, communication, and computation, and thus their designs are sophisticated and cross-disciplinary. On the other hand, efficient edge inference has to overcome a communication bottleneck caused by the low-latency requirements of real-time AI services (e.g., human motion recognition in auto-driving) and the need to upload the sensory data from potentially many devices for aggregation to suppress sensing noise. One promising solution for these challenges is called over-the-air computation (AirComp) that leverages the waveform-superposition property

of a multi-access channel to realize over-the-air aggregation of analog modulated sensory data simultaneously transmitted by multiple devices. The communication-and-computing integration and the enabled simultaneous access promise to dramatically reduce multi-access latency and suppress communication overhead when there are many devices. In this work, we design task-oriented AirComp techniques to support communication-efficient edge inference.

Recent years have witnessed the advancements of edge inference on different fronts. Split inference is arguably the most popular edge-inference architecture, which divides an AI model into two parts: one deployed on resource-limited devices for feature extraction [via, e.g., principal component analysis (PCA) or using a convolutional neural network], and the other at an edge server for completing the remaining computation-intensive inference task. Thereby, the avoidance of direct data uploading helps preserve data privacy and the offloading of computation-intensive task to the server overcomes the devices' resource constraints. These advantages motivate us to adopt split inference in this work. One main research focus on edge split inference is to balance the trade-off between the computation and communication overhead on edge device via, e.g., compressing the feature map of the split layer [18], [19], [20], a two-step pruning strategy [21], progressive feature transmission [22], setting early existing point [23], [24], and joint source and channel coding using deep neural networks [25]. However, the aforementioned designs focus on the case of a single edge device. This pertains to scenarios where the device either senses the source in a narrow view to obtain highly accurate sensory data by focusing a single angle, or obtains a noise-corrupted wide-view sensory data for wide angle object detection by, e.g., scanning from angle to angle [26]. To address the incomplete feature space caused by the narrow-view sensing, a multi-device cooperated multi-view edge inference scheme is proposed in [10] to maximize the inference accuracy via the design of task-oriented sensing, computation, and communication integration. However, the issue of suppressing the noise of wide-view sensory data for inference accuracy enhancement remains unresolved and is the theme of this paper.

In this work, a multi-device edge inference system is considered. Each device obtains a noise-corrupted version of the ground-true wide-view sensory data and extracts from it a noisy local feature vector using simple linear operations like PCA. To suppress the sensing noise, we adopt a common approach, which averages out the noise via a weighted sum of all local feature vectors [27], [28]. To this end, the technique of AirComp can be employed to enhance the communication efficiency due to its capability in supporting fast data aggregation from a large number of devices [29]. Specifically, in AirComp, signals from all devices are allowed to transmit simultaneously

over the same frequency band. At the receiver, the functional value of the aggregated signals is directly calculated using the waveform superposition property of wireless channels, instead of first decoding the individual data stream from each device. There has been comprehensive research for the efficient implementation of AirComp, including the design of beamforming in multi-input-multi-output (MIMO) system (see, e.g., [30], [31], [32], [33], [34]), power control for combating the non-uniform channel fading (see, e.g., [35]), the investigation of tradeoff between computation and energy efficiency (see, e.g., [36]), the design of unmanned aerial vehicle (UAV) assisted AirComp (see, e.g., [37]), etc. In view of its low-latency merit in wireless data aggregation, AirComp has been a promising technique widely exploited in federated edge learning for communication efficiency enhancement (see, e.g., [38], [39], [40], [41], [42]). Most recently, researchers have proposed an AirComp based edge inference system, where the same inference task is performed in multiple devices and a server aggregates the local inference results and makes a final decision based on majority voting [43]. However, such existing design builds on the traditional on-device inference, which causes huge computation overhead at the resource-limited edge devices. It remains an uncharted area to implement edge split inference using AirComp, and thus motivates the current work.

In the considered multi-device edge inference system, the server is equipped with multiple antennas and all devices are equipped with one single antenna. The server aggregates all local feature vectors in a low-latency manner via AirComp to attain a denoised feature vector for the subsequent inference task. In such a system, the transmit precoding and receive beamforming need to be jointly designed to rein in the aggregation error caused by AirComp and maximize the inference accuracy. It is noteworthy that the traditional AirComp design criterion, i.e., minimum mean square error (MMSE) used in existing literature, is no longer effective for edge split inference systems. To be specific, the schemes based on MMSE minimize the average distortion between aggregated data by AirComp and the ideally aggregated one without any corruption by channel fading and noise. However, in the context of edge inference, the MSE measure fails to respect the fact that some feature dimensions are more sensitive to the aggregation error than the others when the ultimate inference accuracy is concerned. As an example, a classification task is shown in Fig. 1, whose feature vector has two elements (dimensions). It is observed that feature element 2 is more tolerant to distortion than element 1 in terms of violating the inference accuracy. Obviously, in the case of MMSE, the non-uniform importance levels at different feature elements are ignored and thus may lead to poor performance. To address

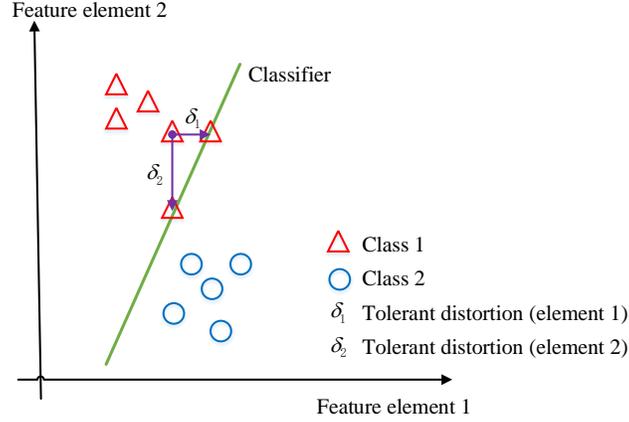


Fig. 1. Various distortion tolerance of different feature elements in classification tasks: For a distortion level δ_1 obtained under the MMSE criterion, incorrect inference occurs if it is on element 1, but the inference result is correct if it is on element 2.

this issue, the best approach is directly maximizing the inference accuracy in the AirComp design. However, it is not tractable to measure the instantaneous inference accuracy prior to the execution of the inference task. Alternatively, this work adopts an approximate but tractable metric, namely discriminant gain, as the surrogate accuracy measure for classification tasks. The metric is originally proposed in [22] building on the well-known KL divergence [44]. Specifically, for arbitrary two classes in the Euclidean feature space, discriminant gain is the distance of their centroids normalized by their covariance. With a larger distance, the two classes are better separated, which implies a higher inference accuracy. However, the joint design of transmit precoding and receive beamforming in AirComp under the criterion of discriminant gain maximization still faces the challenges due to the complicated form of the objective function, and the coupling between the control variables.

To address the challenges above, the solution framework for inference-task-oriented AirComp is proposed in this paper. The detailed contributions are summarized below.

- **Multi-device Over-the-air Inference Systems:** An AirComp based multi-device edge split inference system is established, where the feature vector used for inference at the server is estimated by aggregating all noisy local feature vectors via AirComp. In each time slot, two real feature elements are linearly analog modulated into a complex scalar symbol; feature vectors of different devices are transmitted simultaneously as blocks of symbols. Under the system settings, the impact of the sensing noise and channel noise on the inference accuracy is theoretically characterized by the derived discriminant gain in closed-form.

- **Task-oriented AirComp Transceiver Design for Edge Inference:** Based on the derived discriminant gain, a joint transmit precoding and receive beamforming design problem for the over-the-air inference system is formulated as a problem of maximizing the discriminant gain. To tackle the challenging non-convex problem, the method of variables transformation is first applied to convert it to an equivalent difference of convex (d.c.) form, which is further solved via using the technique of successive convex approximation (SCA) (see e.g., [45]) to yield a sub-optimal solution.
 - It is noteworthy that the sub-optimal solution meet all the Karush – Kuhn – Tucker (KKT) conditions of the original problem, from which one can derive the insight that the optimal beam steered by a device to the edge server has the power inversely proportional to the sensing noise incurred by the device. This further suggests that optimal joint beamforming design should favor those devices with good sensing quality (i.e., small sensing noise).
- **Performance Evaluation:** Extensive simulations using the wireless sensing simulator proposed in [46] have been performed while taking into account the specific task of wide-view human motion recognition with two inference models, i.e., support vector machine (SVM) and multi-layer perception (MLP) neural network, respectively. It is demonstrated that, for both models, maximizing the discriminant gain is effective in maximizing the inference accuracy. Furthermore, it is demonstrated that the proposed scheme significantly outperforms the benchmarking scheme (designed using the criteria of MMSE) in terms of inference accuracy.

II. SYSTEM MODEL

A. Network and Sensing Model

Consider an edge inference system where there is one server equipped with a multi-antenna access point (AP) and K single-antenna sensing devices (e.g., radar sensors and cameras), as shown in Fig. 2. The server aims at aggregating the noisy local feature vectors, which are extracted from the real-time noise-corrupted sensory data, on all devices to form a global denoised feature vector for completing the remaining inference task. Specifically, the noise-corrupted sensory data obtained by device k is given as

$$\mathbf{z}_k = \mathbf{z} + \mathbf{e}_k, \quad (1)$$

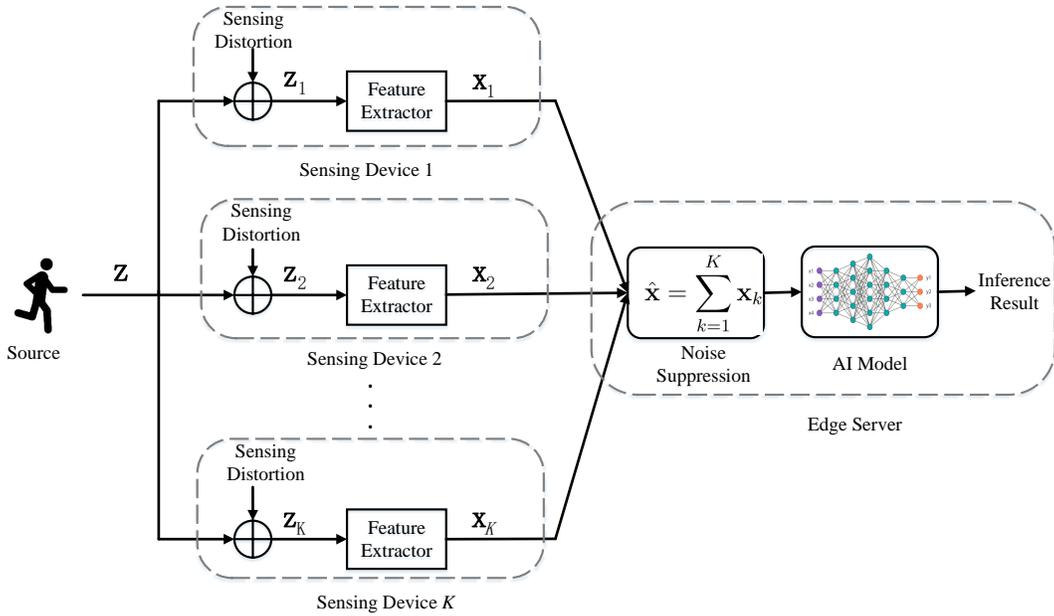


Fig. 2. Model of Over-the-air Computation Based Edge Inference Systems.

where $\mathbf{z} = [z_1, z_2, \dots, z_S]^T$ is the ground-true sensory data of the source, $\mathbf{z}_k = [z_{k,1}, z_{k,2}, \dots, z_{k,S}]^T$ is the local observation of device k , \mathbf{e}_k is the sensing noise, and S is the dimension of the raw sensory data. According to [27], [28], different elements of the sensing noise vector follow identical and independent zero-mean Gaussian distributions:

$$\mathbf{e}_k \sim \mathcal{N}(\mathbf{0}, \epsilon_k^2 \mathbf{I}), \quad (2)$$

where $\mathcal{N}(\cdot, \cdot)$ is the Gaussian distribution, ϵ_k^2 is the sensing noise power, and $\mathbf{I} \in \mathbb{R}^{S \times S}$ is the identical matrix.

The server and the sensing devices communicate via wireless links. Time-division multiple access is adopted. The channels are assumed to be static in each time slot and varying among different slots. The channel gain of device k is denoted as $\mathbf{h}_k \in \mathbb{C}^N$, with N being the number of receive antennas at the server and \mathbb{C}^N being a complex vector space with the dimension of N . Moreover, the server is assumed to work as the coordinator and has the ability to acquire the channel gains of all devices' uplink links.

B. Feature Generation and Distribution

In this part, the feature generation procedure is first introduced, followed by the description of the feature distribution.

1) *Feature Generation*: As transmitting the raw sensory data with large dimensions causes large communication overhead as well as violates the data privacy, an alternative approach is to move the feature extraction part (e.g., PCA and convolutional operations) of an AI model on devices. In this work, PCA is adopted to extract a latent low-dimensional feature sub-space from the raw sensory data on each device. The detailed procedure is described as follows.

- At the training stage, PCA is first performed by the server over the offline training dataset to extract the principal dimensions of each sample. The learning model is trained using the principal feature dimensions.
- At the inference stage, before the server aggregates the local observations from each device, the principal eigen-space is broadcast to each device. For each device, the local feature vector is extracted by projecting the sensory data into the principal eigen-space, and then transmitted.

Thereby, the extracted local feature vectors from the sensory data \mathbf{z}_k can be expressed as

$$\mathbf{x}_k = \mathbf{U}^T \mathbf{z}_k = \mathbf{U}^T \mathbf{z} + \mathbf{U}^T \mathbf{e}_k = \mathbf{x} + \mathbf{d}_k, \quad 1 \leq k \leq K, \quad (3)$$

where \mathbf{U} is a $S \times M$ real column unitary matrix representing the principal eigen-space of PCA, M is the dimension of the principal feature eigen-space,

$$\mathbf{x} = \mathbf{U}^T \mathbf{z} = [x_1, x_2, \dots, x_M]^T, \quad (4)$$

is the ground-true feature vector, and

$$\mathbf{d}_k = \mathbf{U}^T \mathbf{e}_k, \quad 1 \leq k \leq K, \quad (5)$$

is the projected noise vector of device k . By substituting the distribution of \mathbf{e}_k in (2), the distribution of \mathbf{d}_k can be derived as

$$\mathbf{d}_k \sim \mathcal{N}(\mathbf{0}, \epsilon_k^2 \mathbf{I}), \quad 1 \leq k \leq K, \quad (6)$$

where the variance is derived from $\mathbb{E}(\mathbf{d}_k^T \mathbf{d}_k) = \mathbb{E}(\mathbf{U}^T \mathbf{e}_k \mathbf{e}_k^T \mathbf{U}) = \mathbf{U}^T \mathbb{E}(\mathbf{e}_k \mathbf{e}_k^T) \mathbf{U} = \epsilon_k^2 \mathbf{I}$.

2) *Feature Distribution*: By considering a classification inference task with L classes and following the setting in [22], the ground-true feature vector \mathbf{x} is assumed to follow a Gaussian mixture as

$$\mathcal{F}(\mathbf{x}) = \frac{1}{L} \sum_{\ell=1}^L \mathcal{F}_\ell(\mathbf{x}), \quad (7)$$

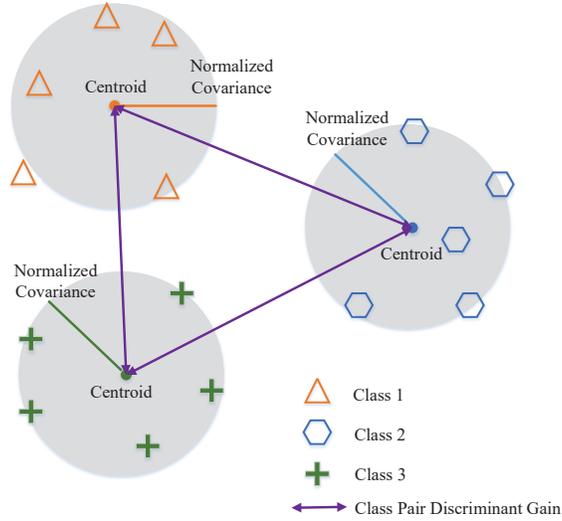


Fig. 3. Geometric interpretation of discriminant gain in the feature space.

where $\mathcal{F}_\ell(\mathbf{x})$ is the Gaussian distribution of \mathbf{x} in terms of the ℓ -th class. As PCA is performed, different feature elements are linearly independent. Thereby, $\mathcal{F}_\ell(\mathbf{x})$ can be written as

$$\mathcal{F}_\ell(\mathbf{x}) = \mathcal{N}(\boldsymbol{\mu}_\ell, \boldsymbol{\Sigma}), \quad 1 \leq \ell \leq L, \quad (8)$$

where $\boldsymbol{\mu}_\ell \in \mathbb{R}^M$ is the centroid of the ℓ -th class, given as

$$\boldsymbol{\mu}_\ell = [\mu_{\ell,1}, \mu_{\ell,2}, \dots, \mu_{\ell,M}]^T, \quad 1 \leq \ell \leq L, \quad (9)$$

and $\boldsymbol{\Sigma} \in \mathbb{R}^{M \times M}$ is a diagonal covariance matrix, given as

$$\boldsymbol{\Sigma} = \text{diag}\{\sigma_1^2, \sigma_2^2, \dots, \sigma_M^2\}. \quad (10)$$

C. Inference Capability

In this work, the metric *discriminant gain* proposed in [22] is adopted as the inference accuracy measure for classification tasks. For arbitrary two classes, the discriminant gain represents the distance between their centroids in the Euclidean feature space under normalized covariance, as presented in Fig. 3. That says, a larger discriminant gain between two classes means that they are more likely to be differentiated, and thus implies a higher inference accuracy. In the sequel, the mathematical model of discriminant gain is introduced.

Discriminant gain is derived from the well-known KL divergence proposed in [44]. Consider an arbitrary class pair, say classes ℓ and ℓ' , and the feature space expanded by the feature vector

x. Based on the distribution of \mathbf{x} in (7) and according to [22], the pair-wise discriminant gain is defined as

$$\begin{aligned} G_{\ell,\ell'}(\mathbf{x}) &= D_{KL}[\mathcal{F}_\ell(\mathbf{x}) \parallel \mathcal{F}_{\ell'}(\mathbf{x})] + D_{KL}[\mathcal{F}_{\ell'}(\mathbf{x}) \parallel \mathcal{F}_\ell(\mathbf{x})], \\ &= \int_{\mathbf{x}} \mathcal{F}_\ell(\mathbf{x}) \log \left[\frac{\mathcal{F}_{\ell'}(\mathbf{x})}{\mathcal{F}_\ell(\mathbf{x})} \right] d\mathbf{x} + \int_{\mathbf{x}} \mathcal{F}_{\ell'}(\mathbf{x}) \log \left[\frac{\mathcal{F}_\ell(\mathbf{x})}{\mathcal{F}_{\ell'}(\mathbf{x})} \right] d\mathbf{x}, \\ &= (\boldsymbol{\mu}_\ell - \boldsymbol{\mu}_{\ell'})^T \boldsymbol{\Sigma}^{-1} (\boldsymbol{\mu}_\ell - \boldsymbol{\mu}_{\ell'}), \quad \forall(\ell, \ell'), \end{aligned} \quad (11)$$

where $D_{KL}(\cdot \parallel \cdot)$ is the KL divergence defined in [44]. As different feature elements are independent, it follows that

$$G_{\ell,\ell'}(\mathbf{x}) = \sum_{m=1}^M G_{\ell,\ell'}(x_m), \quad (12)$$

where x_m is the m -th element in \mathbf{x} and $G_{\ell,\ell'}(x_m)$ is given as

$$G_{\ell,\ell'}(x_m) = \frac{(\mu_{\ell,m} - \mu_{\ell',m})^2}{\sigma_m^2}, \quad 1 \leq m \leq M, \quad (13)$$

and the other notations follow that in (9) and (10). Then, the overall discriminant gain is defined as the average of all pair-wise discriminant gains in (11), given as

$$G(\mathbf{x}) = \frac{2}{L(L-1)} \sum_{\ell'=1}^L \sum_{\ell < \ell'} G_{\ell,\ell'}(\mathbf{x}) = \frac{2}{L(L-1)} \sum_{\ell'=1}^L \sum_{\ell < \ell'} \sum_{m=1}^M G_{\ell,\ell'}(x_m) = \sum_{m=1}^M G(x_m), \quad (14)$$

where $G(x_m)$ is the discriminant gain of the m -th feature elements, given as

$$G(x_m) = \frac{2}{L(L-1)} \sum_{\ell'=1}^L \sum_{\ell < \ell'} \frac{(\mu_{\ell,m} - \mu_{\ell',m})^2}{\sigma_m^2}, \quad 1 \leq m \leq M. \quad (15)$$

D. AirComp Model

The technique of AirComp is used to aggregate the local feature vectors $\{\mathbf{x}_k\}$ from all devices, as it can suppress the sensing noise and significantly enhance the communication efficiency. Specifically, each device transmits a complex scalar symbol via the single antenna in each time slot. The real part and the imaginary part of the complex scalar symbol contain one feature element, respectively. At the server, AirComp is performed to aggregate the two feature elements and estimate their ground-true versions. Thereby, the whole feature vector can be grouped into different element pairs, which can be sequentially transmitted in a time-division way over several time slots. Obviously, the design of AirComp in all time slots is the same. Without loss of generality, in the sequel, the transmission in an arbitrary time slot is considered.

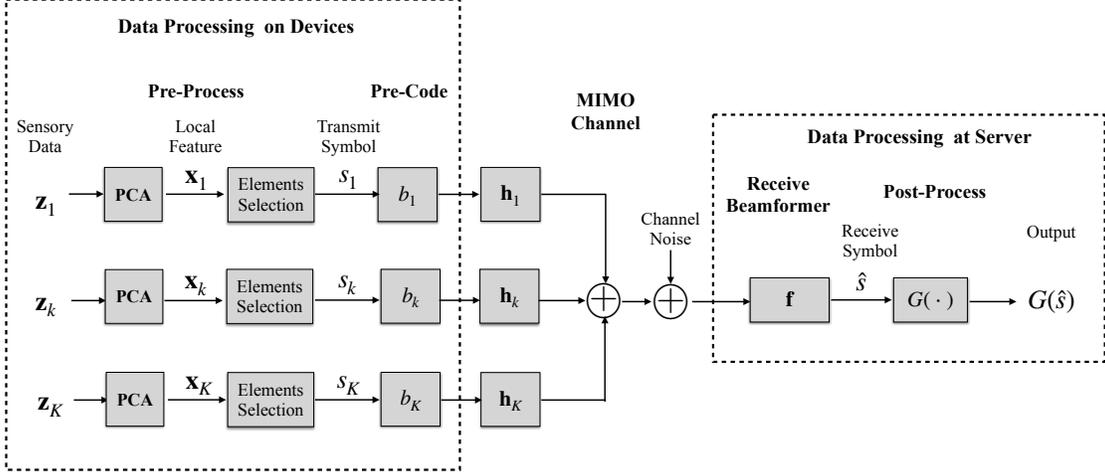


Fig. 4. Block diagram of AirComp for Feature Aggregation.

Consider the case where the server aggregates the m_1 -th and m_2 -th local elements from all devices, say $\{x_{k,m_1}, x_{k,m_2}, 1 \leq k \leq K\}$, to estimate the ground-true feature elements $\{x_{m_1}, x_{m_2}\}$. The procedure of AirComp is shown in Fig. 4 and is described as follows. For an arbitrary device k , its local sensory data is first pre-processed by PCA to extract the principal feature elements. Then, the m_1 -th and m_2 -th principal feature dimensions, say x_{k,m_1} and x_{k,m_2} , are combined in one symbol for transmission, as

$$s_k = x_{k,m_1} + jx_{k,m_2}, \quad 1 \leq k \leq K, \quad (16)$$

where $s_k \in \mathbb{C}$ is the transmitted symbol and j represents the imaginary unit. Next, s_k is further pre-coded with a scalar $b_k \in \mathbb{C}$ and transmitted over a MIMO channel. At the server, the receive signal is the aggregation of all transmit symbols, given as

$$\mathbf{y}_m = \sum_{k=1}^K \mathbf{h}_k b_k s_k + \mathbf{n}, \quad (17)$$

where \mathbf{n} is the additive white Gaussian noise with the following distribution:

$$\mathbf{n} \sim \mathcal{N}(\mathbf{0}, \delta_0^2 \mathbf{I}), \quad (18)$$

and δ_0^2 is the noise variance. Next, a receive beamforming vector $\mathbf{f} \in \mathbb{C}^N$ is used to aggregate all local symbols $\{s_k\}$ to generate the estimates of the ground-true feature elements x_{m_1} and x_{m_2} . Specifically, the received symbol after receive beamforming can be written as

$$\hat{s} = \mathbf{f}^H \mathbf{y}_m = \mathbf{f}^H \sum_{k=1}^K \mathbf{h}_k b_k s_k + \mathbf{f}^H \mathbf{n}. \quad (19)$$

It follows that the estimates are given by

$$\begin{cases} \hat{x}_{m_1} = \text{Re}(\hat{s}) = \text{Re}\left(\mathbf{f}^H \sum_{k=1}^K \mathbf{h}_k b_k s_k + \mathbf{f}^H \mathbf{n}\right), \\ \hat{x}_{m_2} = \text{Im}(\hat{s}) = \text{Im}\left(\mathbf{f}^H \sum_{k=1}^K \mathbf{h}_k b_k s_k + \mathbf{f}^H \mathbf{n}\right), \end{cases} \quad (20)$$

where \hat{x}_{m_1} and \hat{x}_{m_2} are the estimates of x_{m_1} and x_{m_2} respectively, $\text{Re}(\cdot)$ and $\text{Im}(\cdot)$ are the functions to extract real part and imaginary part of one complex number respectively, and other notations follow that in (17). Finally, \hat{x}_{m_1} and \hat{x}_{m_2} are post-processed to output the discriminant gain $G(\hat{x}_{m_1}) + G(\hat{x}_{m_2})$.

III. PROBLEM FORMULATION AND SIMPLIFICATION

A. Problem Formulation

Different from the traditional AirComp design, which aims at minimizing the distortion between the estimated feature elements $\{\hat{x}_{m_1}, \hat{x}_{m_2}\}$ and the ground-true ones $\{x_{m_1}, x_{m_2}\}$ without taking into account the performance metric of the specific tasks, in this work, the design objective follows the task-oriented principle and maximizes the inference accuracy measured by the sum discriminant gains of \hat{x}_{m_1} and \hat{x}_{m_2} , given as

$$\max G = G(\hat{x}_{m_1}) + G(\hat{x}_{m_2}), \quad (21)$$

where \hat{x}_{m_1} and \hat{x}_{m_2} defined in (20) are the estimates of the ground-true feature elements, and $G(\hat{x}_{m_1})$ and $G(\hat{x}_{m_2})$ are the corresponding discriminant gains.

Besides, there is one constraint on the transmit power of each device, given by

$$b_k \mathbb{E}(s_k s_k^H) b_k^H \leq P_k, \quad 1 \leq k \leq K, \quad (22)$$

where b_k is the precoding scalar at device k , b_k^H is the hermitian of b_k , s_k is the transmit symbol, and P_k is the total transmit power of device k . The transmit symbol variance, say $\mathbb{E}(s_k s_k^H)$, can be estimated from the offline training data samples, and thus is known by the edge server as prior information. Therefore, the power constraint in (22) can be re-written as

$$b_k b_k^H \leq \hat{P}_k, \quad 1 \leq k \leq K, \quad (23)$$

where \hat{P}_k is the maximum transmit precoding power, given as

$$\hat{P}_k = \frac{P_k}{\mathbb{E}(s_k s_k^H)}, \quad 1 \leq k \leq K. \quad (24)$$

In summary, the discriminant gain maximization problem can be written as

$$\begin{aligned}
 \text{(P1)} \quad & \max_{\{b_k\}, \mathbf{f}} G = G(\hat{x}_{m_1}) + G(\hat{x}_{m_2}), \\
 & \text{s.t. } b_k b_k^H \leq \hat{P}_k, \quad 1 \leq k \leq K.
 \end{aligned} \tag{25}$$

The formulation of (P1) follows the task-oriented principle. To be specific, the inference accuracy measured by the discriminant gain is maximized instead of using the MMSE criterion. That's because the MMSE criterion ignores the fact that a same distortion level on different feature elements has different impacts on the inference accuracy, and thus leads to poor performance. The task-oriented formulation, however, causes new challenges. To begin with, the discriminant gain has a complicated non-convex sum-of-ratios form. Besides, the design of the receive beamforming \mathbf{f} and the precoding scalars $\{b_k\}$ are coupled [see (20)]. Moreover, the feature elements in the received symbol are cross coupled, i.e., each estimated feature element defined in (20) could be a linear combination of the ground-true elements x_{m_1} and x_{m_2} , due to channel rotation. This leads to a complicated distribution of \hat{x}_{m_1} and \hat{x}_{m_2} , and thus a complicated expression of the discriminant gains $G(\hat{x}_{m_1})$ and $G(\hat{x}_{m_2})$.

B. Discriminant Gains with Zero-Forcing Pre-coders

To address the challenges mentioned above, we simplify (P1) with two steps in this part. The well-known zero-forcing (ZF) precoders are first used to simplify the estimated feature elements $\{\hat{x}_{m_1}, \hat{x}_{m_2}\}$. Then, based on the ZF precoders, the discriminant gains, i.e., $G(\hat{x}_{m_1})$ and $G(\hat{x}_{m_2})$, are derived to simplify the objective function.

1) *ZF precoders:* First, the ZF design is given by

$$\mathbf{f}^H \mathbf{h}_k b_k = c_k, \quad 1 \leq k \leq K, \tag{26}$$

where \mathbf{f}^H is the receive beamforming vector, \mathbf{h}_k is the channel vector of device k , b_k is the precoder of device k , and $c_k \geq 0$ is a real number representing the receive signal strength, or called steering power, from device k . Then, the ZF precoders can be derived as

$$b_k = \frac{c_k \mathbf{h}_k^H \mathbf{f}}{\mathbf{h}_k^H \mathbf{f} \mathbf{f}^H \mathbf{h}_k}, \quad 1 \leq k \leq K. \tag{27}$$

It follows that the power constraint in (P1) can be re-written as

$$c_k^2 \leq \hat{P}_k \mathbf{h}_k^H \mathbf{f} \mathbf{f}^H \mathbf{h}_k, \quad 1 \leq k \leq K. \tag{28}$$

Besides, by substituting the precoders in (27) and \hat{s}_k in (16) into the estimates \hat{x}_{m_1} and \hat{x}_{m_2} in (20), we can obtain

$$\begin{cases} \hat{x}_{m_1} = \text{Re} \left(\sum_{k=1}^K c_k s_k + \mathbf{f}^H \mathbf{n} \right) = \sum_{k=1}^K c_k x_{k,m_1} + \text{Re}(\mathbf{f}^H \mathbf{n}), \\ \hat{x}_{m_2} = \text{Im} \left(\sum_{k=1}^K c_k s_k + \mathbf{f}^H \mathbf{n} \right) = \sum_{k=1}^K c_k x_{k,m_2} + \text{Im}(\mathbf{f}^H \mathbf{n}), \end{cases} \quad (29)$$

where the notations follow that in (16), (20), and (27).

2) *Discriminant Gains:* To achieve the discriminant gain G , in the sequel, the distributions of the local transmit feature elements $\{x_{k,m_1}, x_{k,m_2}\}$ are first derived. Then, based on the ZF precoders, the distribution of the received elements $\{\hat{x}_{m_1}, \hat{x}_{m_2}\}$ are derived. Next, the discriminant gains are obtained, followed by the derivation of a simplified problem of (P1).

First, recall the local elements x_{k,m_1} and x_{k,m_2} are given by

$$x_{k,m_i} = x_{m_i} + d_{k,m_i}, \quad i = 1, 2, 1 \leq k \leq K, \quad (30)$$

where the distribution of the ground-true element x_{m_i} is given by

$$x_{m_i} \sim \frac{1}{L} \sum_{\ell=1}^L \mathcal{N}(\mu_{\ell,m_i}, \sigma_{m_i}^2), \quad i = 1, 2, \quad (31)$$

according the distribution of \mathbf{x} in (7), (9), and (10), and the distribution of the noise d_{k,m_i} is given by

$$d_{k,m_i} \sim \mathcal{N}(0, \epsilon_k^2), \quad i = 1, 2, \quad (32)$$

according to the distribution of \mathbf{d}_k in (6). Subsequently, the following lemma in terms of x_{k,m_i} 's distribution can be obtained.

Lemma 1. *The distribution of the local elements $\{x_{k,m_i}\}$ can be derived as*

$$x_{k,m_i} \sim \frac{1}{L} \sum_{\ell=1}^L \mathcal{N}(\mu_{\ell,m_i}, \sigma_{m_i}^2 + \epsilon_k^2), \quad i = 1, 2, 1 \leq k \leq K, \quad (33)$$

Proof: Please see Appendix A.

Then, by substituting the distributions of $\{x_{k,m_1}, x_{k,m_2}\}$ in (33) and the distribution of the channel noise \mathbf{n} in (18) into the received feature elements $\{\hat{x}_{m_1}, \hat{x}_{m_2}\}$ in (29), their distributions can be derived as shown in Lemma 2.

Lemma 2. *The distribution of the estimated feature elements $\{\hat{x}_{k,m_i}\}$ are given by*

$$\hat{x}_{m_i} \sim \frac{1}{L} \mathcal{N}(\hat{\mu}_{\ell,m_i}, \hat{\sigma}_{m_i}^2), \quad i = 1, 2, \quad (34)$$

where the centroids $\{\hat{\mu}_{\ell,m_i}\}$ and the variance $\{\hat{\sigma}_{m_i}^2\}$ are

$$\begin{cases} \hat{\mu}_{\ell,m_i} = \sum_{k=1}^K c_k \mu_{\ell,m_i}, & i = 1, 2, \\ \hat{\sigma}_{m_i}^2 = \sigma_{m_i}^2 \left(\sum_{k=1}^K c_k \right)^2 + \sum_{k=1}^K c_k^2 \epsilon_k^2 + \frac{\delta_0^2}{2} (\mathbf{f}_1^T \mathbf{f}_1 + \mathbf{f}_2^T \mathbf{f}_2), & i = 1, 2, \end{cases} \quad (35)$$

δ_0^2 is the channel noise power, $\mathbf{f}_1 = \text{Re}(\mathbf{f})$ and $\mathbf{f}_2 = \text{Im}(\mathbf{f})$ are the real part and imaginary part of the receive beamforming \mathbf{f} respectively, and other notations follow that in (31) and (32).

Proof: Please see Appendix B.

Next, based on the distributions in Lemma 2 and the definition of discriminant gain in (15), the discriminant gains of $\{x_{k,m_1}, x_{k,m_2}\}$ can be derived as

$$G(\hat{x}_{m_i}) = \frac{2}{L(L-1)} \sum_{\ell'=1}^L \sum_{\ell < \ell'} \frac{(\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2}{\hat{\sigma}_{m_i}^2}, \quad i = 1, 2, \quad (36)$$

where $\{\hat{\mu}_{\ell,m_i}\}$ and $\{\hat{\sigma}_{m_i}^2\}$ are defined in (35).

Finally, by substituting the discriminant gains of $\{x_{k,m_1}, x_{k,m_2}\}$ in (36) and the power constraint in (28) into (P1), it can be equivalently derived as

$$\begin{aligned} \text{(P2)} \quad \max_{\{c_k\}, \mathbf{f}_1, \mathbf{f}_2} \quad G &= \frac{2}{L(L-1)} \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \frac{(\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2}{\hat{\sigma}_{m_i}^2}, \\ \text{s.t.} \quad c_k^2 &\leq \hat{P}_k \mathbf{h}_k^H (\mathbf{f}_1 \mathbf{f}_1^T + \mathbf{f}_2 \mathbf{f}_2^T) \mathbf{h}_k, \quad 1 \leq k \leq K, \end{aligned} \quad (37)$$

where the notations follow that in (35).

IV. JOINT POWER CONTROL AND RECEIVE BEAMFORMING FOR TASK-ORIENTED AIRCOMP

In this section, variables transformation is first applied to derive (P2) into an equivalent d.c. problem. Then, the method of SCA is adopted to address it and obtain the joint design of steering power control and receive beamforming. Finally, the struction of the obtained solution is investigated.

A. An Equivalent D.C. Problem

In this part, to simplify (P2), the following variables are first defined:

$$\alpha_{\ell,\ell',m_i} = \frac{(\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2}{\hat{\sigma}_{m_i}^2}, \quad \forall(\ell, \ell', m_i), \quad (38)$$

where α_{ℓ,ℓ',m_i} represents the per class pair discriminant gain of the m_i -th received element \hat{x}_{m_i} .

It follows that (P2) can be equivalently derived as

$$\begin{aligned} \max_{\{c_k\}, \mathbf{f}_1, \mathbf{f}_2, \{\alpha_{\ell,\ell',m_i}\}} \quad & G = \frac{2}{L(L-1)} \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \alpha_{\ell,\ell',m_i}, \\ \text{s.t.} \quad & c_k^2 \leq \hat{P}_k \mathbf{h}_k^H (\mathbf{f}_1 \mathbf{f}_1^T + \mathbf{f}_2 \mathbf{f}_2^T) \mathbf{h}_k, \quad 1 \leq k \leq K, \\ & (\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2 = \alpha_{\ell,\ell',m_i} \hat{\sigma}_{m_i}^2, \quad \forall(\ell, \ell', m_i), \end{aligned} \quad (39)$$

where

$$(\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2 = \left(\sum_{k=1}^K c_k \right)^2 (\mu_{\ell,m_i} - \mu_{\ell',m_i})^2, \quad \forall(\ell, \ell', m_i), \quad (40)$$

and

$$\hat{\sigma}_{m_i}^2 = \left[\sigma_{m_i}^2 \left(\sum_{k=1}^K c_k \right)^2 + \sum_{k=1}^K c_k^2 \epsilon_k^2 + \frac{\delta_0^2}{2} (\mathbf{f}_1^T \mathbf{f}_1 + \mathbf{f}_2^T \mathbf{f}_2) \right], \quad \forall(\ell, \ell', m_i). \quad (41)$$

Then, it can be shown that using symmetric real and imaginary receive beamformers can achieve the optimal solution of the problem in (39), as presented in the following lemma.

Lemma 3 (Symmetric Receive Beamformers). *Symmetric real and imaginary receive beamformers, as in (42), will not influence the optimality of the problem in (39).*

$$\mathbf{f}_1 = \mathbf{f}_2 = \hat{\mathbf{f}}. \quad (42)$$

Proof: Please see Appendix C.

Besides, it can be further proved that extending the feasible region of the second constraint of the problem in (39), i.e., the equality constraint, has no influence on its optimal solution, as equality should be achieved to obtain the optimum, as presented in Lemma 4.

Lemma 4 (Equivalent Extended Feasible Region). *A problem, which extends the feasible region of the second constraint of the problem in (39) as*

$$\left(\sum_{k=1}^K c_k \right)^2 \left[\frac{(\mu_{\ell,m_i} - \mu_{\ell',m_i})^2}{\alpha_{\ell,\ell',m_i}} - \sigma_{m_i}^2 \right] \geq \sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}}, \quad \forall(\ell, \ell', m_i), \quad (43)$$

and keeps the other constraints and the objective function, achieves the same optimal solution to the problem in (39).

Proof: Please see Appendix D.

Next, based on Lemmas 3 and 4, the problem in (39) can be equally derived as

$$\begin{aligned}
 \max_{\{c_k\}, \hat{\mathbf{f}}, \{\alpha_{\ell, \ell', m_i}\}} G &= \frac{2}{L(L-1)} \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \alpha_{\ell, \ell', m_i}, \\
 \text{(P3)} \quad \text{s.t.} \quad c_k^2 - R_k(\hat{\mathbf{f}}) &\leq 0, \quad 1 \leq k \leq K, \\
 \sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}} + \sigma_{m_i}^2 \left(\sum_{k=1}^K c_k \right)^2 - Q_{\ell, \ell', m_i}(\{c_k\}, \alpha_{\ell, \ell', m_i}) &\leq 0, \quad \forall(\ell, \ell', m_i),
 \end{aligned}$$

where $R_k(\hat{\mathbf{f}})$ and $Q_{\ell, \ell', m_i}(\{c_k\}, \alpha_{\ell, \ell', m_i})$ are the functions defined as

$$\begin{cases} R_k(\hat{\mathbf{f}}) = 2\hat{P}_k \mathbf{h}_k^H \hat{\mathbf{f}} \hat{\mathbf{f}}^T \mathbf{h}_k, & q \leq k \leq K, \\ Q_{\ell, \ell', m_i}(\{c_k\}, \alpha_{\ell, \ell', m_i}) = \left(\sum_{k=1}^K c_k \right)^2 \times \frac{(\mu_{\ell, m_i} - \mu_{\ell', m_i})^2}{\alpha_{\ell, \ell', m_i}}, & \forall(\ell, \ell', m_i). \end{cases} \quad (44)$$

Although (P3) is non-convex due to the two constraints therein, it is a d.c. problem as presented in the following lemma.

Lemma 5. (P3) is a d.c. problem, since c_k^2 , $R_k(\hat{\mathbf{f}})$, $\sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}} + \sigma_{m_i}^2 \left(\sum_{k=1}^K c_k \right)^2$, $Q_{\ell, \ell', m_i}(\{c_k\}, \alpha_{\ell, \ell', m_i})$, and the objective function are differentiable and convex.

Proof: See Appendix E.

In the sequel, the SCA method is used to get a sub-optimal solution based on Lemma 5.

B. SCA Based Solution Approach

In this part, the SCA approach is used to address (P3) for obtaining a sub-optimal solution based on Lemma 5 by iterating over the following two steps:

- *Convex relaxation:* Based on a feasible reference point, (P3) is relaxed into a convex problem, whose feasible region is a subset of that of (P3). Hence, the solution to the relaxed problem is guaranteed to be feasible for (P3).
- *Reference point updating:* The solution of the relaxed convex problem is used as the new reference point for the next iteration.

This process iterates till convergence and the final result can be guaranteed to satisfy the KKT conditions of (P3) [45]. In the sequel, the approach of convex relaxation is first presented, followed by the summary of the overall joint steering power control and receive beamforming algorithm.

1) *Convex Relaxation of (P3)*: Consider an arbitrary SCA iteration $(t+1)$, the reference point is the solution of the relaxed problem in the previous iteration and is denoted as $\{\hat{\mathbf{f}}^{[t]}, c_k^{[t]}, \alpha_{\ell, \ell', m_i}^{[t]}\}$. According to Lemma 5, $R_k(\hat{\mathbf{f}})$ and $Q_{\ell, \ell', m_i}(\{c_k\}, \alpha_{\ell, \ell', m_i})$ are differentiable and convex, and hence are no less than their corresponding first-order Taylor expansions at the reference point:

$$\begin{cases} R_k(\hat{\mathbf{f}}) \geq \hat{R}_k^{[t]}(\hat{\mathbf{f}}), \\ Q_{\ell, \ell', m_i}(\{c_k\}, \alpha_{\ell, \ell', m_i}) \geq \hat{Q}_{\ell, \ell', m_i}^{[t]}(\{c_k\}, \alpha_{\ell, \ell', m_i}), \quad \forall(\ell, \ell', m_i). \end{cases} \quad (45)$$

In the equation above, $\hat{R}_k^{[t]}(\hat{\mathbf{f}})$ and $\hat{Q}_{\ell, \ell', m_i}^{[t]}(\{c_k\}, \alpha_{\ell, \ell', m_i})$ are the corresponding first-order linear expansion functions, given by

$$\hat{R}_k^{[t]}(\hat{\mathbf{f}}) = R(\hat{\mathbf{f}}^{[t]}) + 4\hat{P}_k(\hat{\mathbf{f}} - \hat{\mathbf{f}}^{[t]})^H (\mathbf{h}_k^H \hat{\mathbf{f}}^{[t]} \mathbf{h}_k), \quad 1 \leq k \leq K, \quad (46)$$

and

$$\begin{aligned} \hat{Q}_{\ell, \ell', m_i}^{[t]}(\{c_k\}, \alpha_{\ell, \ell', m_i}) &= Q(\{c_k^{[t]}\}, \alpha_{\ell, \ell', m_i}^{[t]}) + \sum_{k=1}^K A_k^{[t]}(c_k - c_k^{[t]}) \\ &\quad + B_{\ell, \ell', m_i}^{[t]}(\alpha_{\ell, \ell', m_i} - \alpha_{\ell, \ell', m_i}^{[t]}), \end{aligned} \quad (47)$$

where

$$\begin{cases} A_k^{[t]} = \frac{\partial Q}{\partial c_k} \Big|_{c_k=c_k^{[t]}} = \frac{2 \sum_{k=1}^K c_k^{[t]} (\mu_{\ell, m_i} - \mu_{\ell', m_i})^2}{\alpha_{\ell, \ell', m_i}^{[t]}}, \\ B_{\ell, \ell', m_i}^{[t]} = \frac{\partial Q}{\partial \alpha_{\ell, \ell', m_i}} \Big|_{\alpha_{\ell, \ell', m_i}=\alpha_{\ell, \ell', m_i}^{[t]}} = - \left[\frac{(\sum_{k=1}^K c_k^{[t]}) (\mu_{\ell, m_i} - \mu_{\ell', m_i})}{\alpha_{\ell, \ell', m_i}^{[t]}} \right]^2. \end{cases} \quad (48)$$

Next, by substituting the inequalities in (44) into (P3), a relaxed problem can be derived as

$$\begin{aligned} \max_{\{c_k\}, \hat{\mathbf{f}}, \{\alpha_{\ell, \ell', m_i}\}} \quad & G = \frac{2}{L(L-1)} \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \alpha_{\ell, \ell', m_i}, \\ \text{(P4)} \quad & \text{s.t. } c_k^2 \leq \hat{R}_k^{[t]}(\hat{\mathbf{f}}), \quad 1 \leq k \leq K, \\ & \hat{Q}_{\ell, \ell', m_i}^{[t]}(\{c_k\}, \alpha_{\ell, \ell', m_i}) - \left(\sum_{k=1}^K c_k \right)^2 \sigma_{m_i}^2 \geq \sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}}, \quad \forall(\ell, \ell', m_i), \end{aligned}$$

where $\hat{R}_k^{[t]}(\hat{\mathbf{f}})$ and $\hat{Q}_{\ell,\ell',m_i}^{[t]}(\{c_k\}, \alpha_{\ell,\ell',m_i})$ are defined in (46) and (47), respectively. (P4) is convex. The proof is straightforward and hence omitted. To address (P4), the well-known CVX toolbox can be used [47].

2) *Task-oriented AirComp Design*: Based on the convex relaxation approach above, (P3) can be addressed by using the SCA method, which iteratively solves the relaxed convex problem (P4), and updates the reference point using the obtained solution. The detailed procedure is summarized in Algorithm 1.

Algorithm 1: Joint Power Control and Receive Beamforming for Task-oriented AirComp

- 1: **Input:** Channel gains $\{\mathbf{h}_k\}$.
 - 2: **Initialize** $t = 0$ and $\{\hat{\mathbf{f}}^{[0]}, c_k^{[0]}, \alpha_{\ell,\ell',m_i}^{[0]}\}$, which is in the feasible region of (P3).
 - 3: **Loop**
 - 4: $t = t + 1$.
 - 5: Derive (P4), based on the reference point $\{\hat{\mathbf{f}}^{[t-1]}, c_k^{[t-1]}, \alpha_{\ell,\ell',m_i}^{[t-1]}\}$.
 - 5: Solve (P4) and obtain the optimum as $\{\hat{\mathbf{f}}^{[t]}, c_k^{[t]}, \alpha_{\ell,\ell',m_i}^{[t]}\}$.
 - 6: **Until Convergence**
 - 7: The solution is
- $$\hat{\mathbf{f}}^* = \hat{\mathbf{f}}^{[t]}, \quad \left\{ c_k^* = c_k^{[t]}, 1 \leq k \leq K \right\}, \quad \left\{ \alpha_{\ell,\ell',m_i}^* = \alpha_{\ell,\ell',m_i}^{[t]}, \forall (\ell, \ell', m_i) \right\}.$$
- 8: **Output:** $\hat{\mathbf{f}}^*$, $\{c_k^*\}$, and $\{\alpha_{\ell,\ell',m_i}^*\}$.
-

C. A Property of Transmit Power Control

As mentioned, the SCA based solution obtained by Algorithm 1 satisfies the KKT conditions [45]. In this part, some of these conditions are used to find the structure of the solved steering power of each device. To begin with, the Lagrangian function of (P3) is given by

$$\begin{aligned} \mathcal{L}_{\text{P3}} = & -\frac{2}{L(L-1)} \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \alpha_{\ell,\ell',m_i} + \sum_{k=1}^K \beta_k \left(c_k^2 - 2\hat{P}_k \mathbf{h}_k^H \hat{\mathbf{f}} \hat{\mathbf{f}}^T \mathbf{h}_k \right) \\ & + \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell,\ell',m_i} \left[\sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}} + \sigma_{m_i}^2 \left(\sum_{k=1}^K c_k \right)^2 - Q_{\ell,\ell',m_i}(\{c_k\}, \alpha_{\ell,\ell',m_i}) \right], \end{aligned} \quad (49)$$

where $\{\beta_k \geq 0\}$ and $\{\lambda_{\ell,\ell',m_i} \geq 0\}$ are Lagrange multipliers.

Some useful KKT conditions are given by

$$\begin{cases} \frac{\partial \mathcal{L}_{P3}}{\partial c_k} = 0, \quad 1 \leq k \leq K, \\ c_k^2 \leq 2\hat{P}_k \mathbf{h}_k^H \hat{\mathbf{f}} \hat{\mathbf{f}}^T \mathbf{h}_k, \quad 1 \leq k \leq K, \\ \beta_k \left(c_k^2 - 2\hat{P}_k \mathbf{h}_k^H \hat{\mathbf{f}} \hat{\mathbf{f}}^T \mathbf{h}_k \right) = 0, \quad 1 \leq k \leq K. \end{cases} \quad (50)$$

From the first condition, we can obtain

$$\beta_k c_k + \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell, \ell', m_i} \left[c_k \epsilon_k^2 + \left(\sum_{k=1}^K c_k \right) \left(\sigma_{m_i}^2 - \frac{(\mu_{\ell, m_i} - \mu_{\ell', m_i})^2}{\alpha_{\ell, \ell', m_i}} \right) \right] = 0, \quad (51)$$

for all $1 \leq k \leq K$. By using a normalized steering power as $c'_k = c_k / (\sum_{k=1}^K c_k)$ and substituting β_k in (51), the above condition can be further derived as

$$c'_k = \frac{\sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell, \ell', m_i} \left((\mu_{\ell, m_i} - \mu_{\ell', m_i})^2 / \alpha_{\ell, \ell', m_i} - \sigma_{m_i}^2 \right)}{\beta_k + \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell, \ell', m_i} \epsilon_k^2}, \quad 1 \leq k \leq K. \quad (52)$$

From the second condition in (50), we have

$$c_k \leq \sqrt{2\hat{P}_k \mathbf{h}_k^H \hat{\mathbf{f}} \hat{\mathbf{f}}^T \mathbf{h}_k}, \quad (53)$$

where the equality is achieved when $\beta_k \neq 0$ according to the third condition in (50). Then, together with (52), the normalized steering power of device k is given as

$$c'_k = \begin{cases} \frac{\sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell, \ell', m_i} \left((\mu_{\ell, m_i} - \mu_{\ell', m_i})^2 / \alpha_{\ell, \ell', m_i} - \sigma_{m_i}^2 \right)}{\sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell, \ell', m_i} \epsilon_k^2}, & \text{if } \beta_k = 0, \\ \frac{\sqrt{2\hat{P}_k \mathbf{h}_k^H \hat{\mathbf{f}} \hat{\mathbf{f}}^T \mathbf{h}_k}}{\sum_{k_1=1, k_1 \neq k}^K c_{k_1} + \sqrt{2\hat{P}_k \mathbf{h}_k^H \hat{\mathbf{f}} \hat{\mathbf{f}}^T \mathbf{h}_k}}, & \text{if } \beta_k \neq 0, \end{cases} \quad (54)$$

where c_{k_1} with $k_1 \neq k$ is irrelevant to the channel gain and sensing noise power of device k according to (52) and (53). Several observations can be made from (54). If the transmit power or the channel magnitude is large enough, i.e., $\beta_k = 0$ and the equality in (53) is not achieved, the normalized steering power of device k , say c'_k , is inversely proportional to its sensing data noise power ϵ_k^2 . Otherwise (i.e., $\beta_k \neq 0$), c'_k is an increasing function of its channel magnitude.

V. PERFORMANCE EVALUATION

A. Experiment Setup

1) *Communication model*: In this experiment, a multi-user single-input multiple-output network is considered, where K single-antenna devices are distributed randomly within a circle with a radius of 50 meters. The multi-antenna AP (edge server) is located at the circle center. For the k -th device, the channel gain \mathbf{h}_k is modeled as $\mathbf{h}_k = |\varphi_k \boldsymbol{\rho}_k|^2$, where φ_k and $\boldsymbol{\rho}_k$ stand for the large-scale and small-scale fading propagation coefficients, respectively. The large-scale propagation coefficient (in dB) is modeled as $[\varphi_k]_{\text{dB}} = -[\text{PL}_k]_{\text{dB}} + [\zeta_k]_{\text{dB}}$, where $[\text{PL}_k]_{\text{dB}} = 128.1 + 37.6 \log_{10} \text{dist}_k$ (dist_k is the distance in kilometer) is the path loss in dB, and $[\zeta_k]_{\text{dB}}$ accounts for the shadowing in dB. In the simulation, $[\zeta_k]_{\text{dB}}$ is a Gauss-distributed random variable with mean zero and variance σ_ζ^2 . Besides, Rayleigh small-scale fading is assumed, i.e., $\boldsymbol{\rho}_k \sim \mathcal{CN}(0, \mathbf{I})$.

2) *Inference task*: A concrete classification task of human motion recognition to identify four distinct human motions, i.e., child walking, child pacing, adult walking, and adult pacing, is considered. The wireless sensing simulator proposed in [46] is adopted to generate the datasets for this task. Using similar settings as [48], the heights of adults and children are assumed to be uniformly distributed in the intervals [1.6m, 1.9m] and [0.9m, 1.2m], respectively. The speeds of standing, walking, and pacing are set as 0 m/s, $0.5H$ m/s, and $0.25H$ m/s, respectively, where H is the height value. The heading of the moving human is set to be uniformly distributed in $[-180^\circ, 180^\circ]$.

3) *Inference model*: Two AI models based on SVM and MLP neural networks are used for the inference task. The neural network model consists of two hidden layers, each with 80 and 40 neurons. The total number of training data samples is 6400, which are assumed to have no noise corruption during the training of both AI models. The testing dataset includes 1600 noise-corrupted data samples, where the noise power is determined by the three schemes.

Unless specified otherwise, other simulation parameters are stated in Table I. All experiments are implemented using Python 3.8 on a Linux server with one NVIDIA[®] GeForce[®] RTX 3090 GPU 24GB and one Intel[®] Xeon[®] Gold 5218 CPU.

B. Inference Algorithms

For comparison, we consider three schemes as follows.

TABLE I
SIMULATION PARAMETERS

Parameter	Value	Parameter	Value
Number of ISAC devices, K	3	Channel noise variance, δ_0^2	1
feature noise variance, ϵ_k^2	0.4	Number of receive antennas, N_r	8
Number of dimension after PCA, N_K	12	Number of classes, L	4
Training data sizes, B	6400	Transmit power, P_k	12 mdB
Variance of shadow fading, σ_ζ^2	8 dB	Communication channel noise power, δ_c^2	10^{-11} W

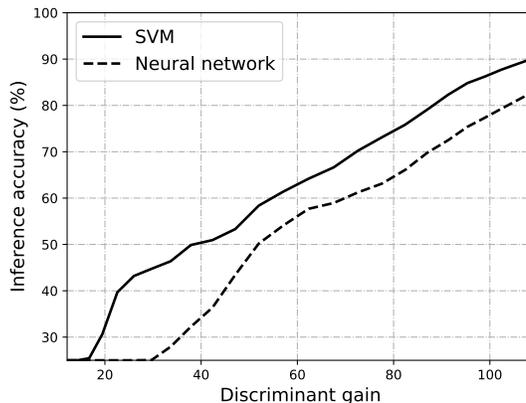
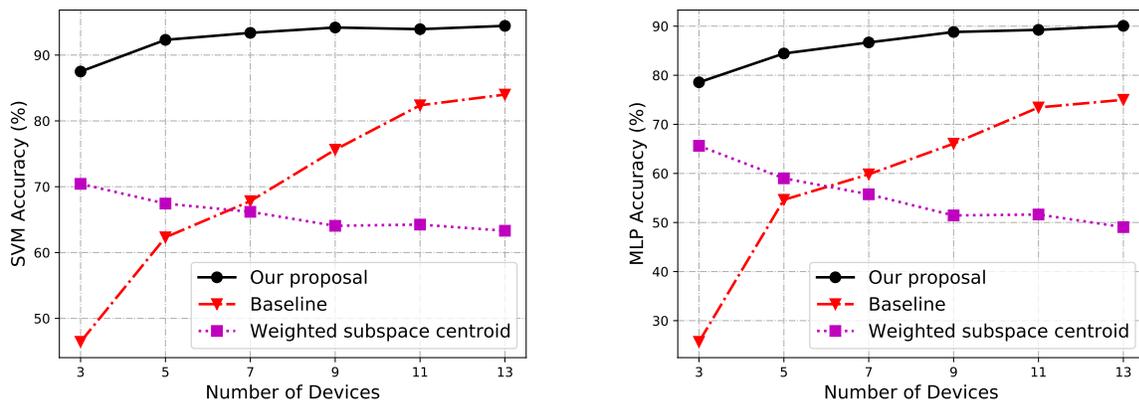


Fig. 5. Inference accuracy versus discriminant gain.

- *Baseline*: In this scheme, random receive beamformer is first used, and then the transmit precoders are selected to satisfy the constraint in (P1).
- *Weighted subspace centroid*: All the parameters are allocated following the AirComp scheme in [31], where the design criterion is MMSE and channel equalization of all devices is performed.
- *Joint design of transmit precoding and receive beamforming (our proposal)*: All parameters are set to follow the proposed scheme Algorithm 1.

C. Experimental Results

This part starts from presenting the relation between the discriminant gain and the corresponding inference accuracy of the two models. Then, the three schemes are compared for both models in terms of the changing number of devices and transmit power. Finally, the influence of feature elements' number on the inference accuracy is shown.



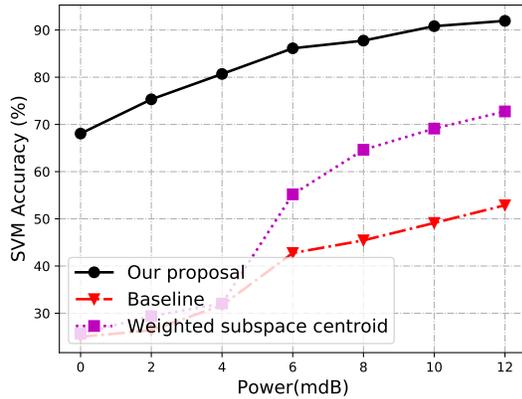
(a) Inference accuracy with SVM versus number of devices (b) Inference accuracy with MLP versus number of devices

Fig. 6. Inference accuracy comparison among different models under different number of devices.

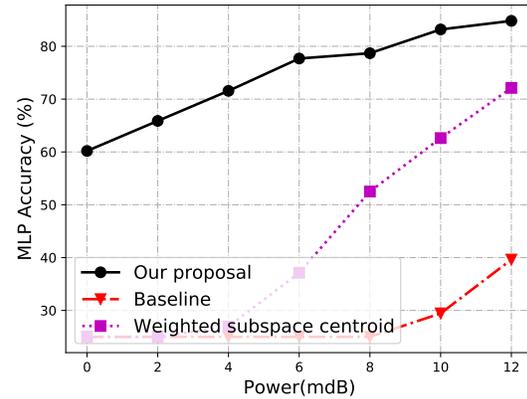
1) *Inference accuracy v.s. discriminant gain*: In Fig. 5, the relation between inference accuracy and discriminant gain for the SVM model and the MLP neural network is presented. To investigate the relation, different values of discriminant gain are obtained by using different transmit power on devices. From the figure, for both models, it is seen that the inference accuracy increases as the discriminant gain grows. Additionally, the SVM beats the neural network, as the training of the latter is overfitting, which has a complex model compared to the simple dataset.

2) *Inference accuracy v.s. number of devices*: The inference accuracy of both models is shown in Fig. 6 in terms of a changing number of devices. It is observed that our proposed scheme has the best performance. Besides, the performance of the weighted subspace centroid scheme decreases with the number of devices. The reason is as follows. Channel equalization is performed among all devices under the target of MMSE in this scheme. As a result, with a growing number of devices, the possibility of deep fading channels increases, which leads to a higher distortion level. Better inference accuracy is obtained in the baseline scheme and our proposed scheme, as the number of devices increases. This is because under the task-oriented principle, different steering powers are permitted for different devices, and thus the data diversity provided by more devices can be fully exploited.

3) *Inference accuracy v.s. transmit power*: The inference accuracy of both models under various transmit powers is shown in Fig. 7. In both cases, improved inference accuracy is acquired as the transmit power rises, since larger transmit powers can more effectively suppress the channel noise. As well, our proposed scheme outperforms the other two schemes.

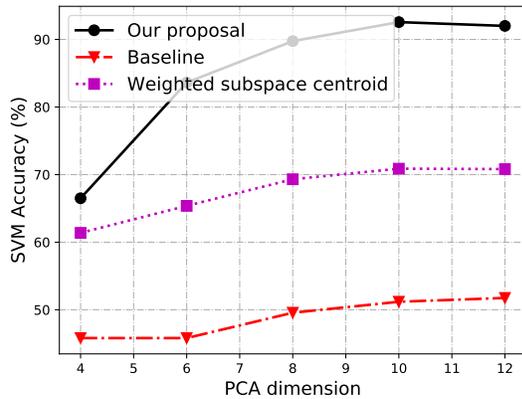


(a) Inference accuracy with SVM versus transmit power

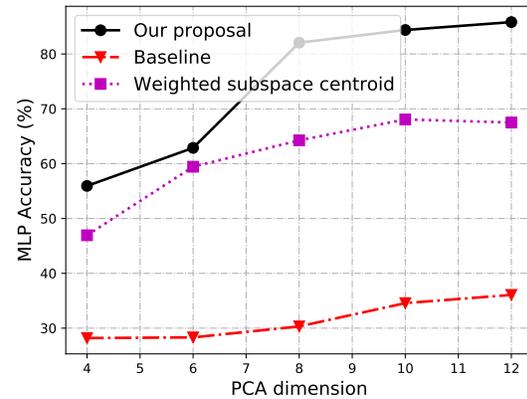


(b) Inference accuracy with MLP versus transmit power

Fig. 7. Inference accuracy comparison among different models under different transmit power.



(a) Inference accuracy with SVM versus PCA dimension



(b) Inference accuracy with MLP versus transmit power

Fig. 8. Inference accuracy comparison among different models under different PCA dimension.

4) *Inference accuracy v.s. number of feature elements*: The inference accuracy of both models in terms of the number of used feature elements in the task is shown in Fig 8. Specifically, the number of used feature elements is sequentially increased following the order of the PCA dimensions from the largest to the least. From the figure, the inference accuracy increases as the number of used feature elements in the inference task. That's because more feature elements increase the dimensions of feature space so that different classes can be better differentiated and a better discriminant gain can be achieved. In addition, the accuracy turns to be saturate at a large number of feature elements, since the added least important feature elements have less contribution to the discriminant gain and the inference accuracy.

The extensive experimental results presented above demonstrate the best performance of the proposed optimal scheme and verify our theoretical analysis.

VI. CONCLUSION

To enhance the performance of multi-device edge inference systems, this paper proposed a task-oriented AirComp scheme. To alleviate the influence of sensing noise on the inference accuracy, it aggregated the noise-corrupted local feature vectors to generate a global one at the server for completing the task. Instead of using the conventional design criterion MMSE, the task-oriented AirComp scheme aimed at directly maximizing the inference accuracy measured by an approximate but tractable metric, called discriminant gain, which represents the averaged centroids distances of different class pairs normalized by their covariance in the Euclidean feature space. A larger discriminant gain means that the classes are better separated, and thus indicates a higher inference accuracy. This task-oriented problem, however, was non-convex due to the complicated form of discriminant gain as well as the couple of transmit precoding and receive beamforming. To tackle this problem, variables transformation was first applied to derive an equivalent d.c. problem. Then, a joint scheme of transmit precoding and receive beamforming was proposed to address the d.c. problem based on the SCA approach. The performance of the proposed scheme was verified using extensive numerical results of a concrete classification task of human motion recognition.

It is noteworthy that the theoretical analysis presented in this work holds for all types of the AI models, from general linear models like SVM to deep neural networks (DNNs), but highly depends on the assumption that the feature vector follows a mixture of Gaussians distribution. We remark that, in some AI tasks, this assumption may not strictly hold, since either the raw data generated by the source or the feature maps generated by the intermediate layers of a DNN may not follow Gaussian distribution. To tackle this issue, a practical approach is to first fit the data to a mixture of Gaussians distribution approximating the ground-truth. Then, the proposed scheme can be extended to these AI tasks, and its performance is verified via extensive experiments based on a high-fidelity human-motion recognition dataset (Please refer to Section V).

This work opens several interesting directions. One is the device selection for further accuracy enhancement by excluding the devices with a weak channel or high sensing noise. Another is to extend the current design to the case where devices are equipped with multi antennas.

APPENDIX

A. Proof of Lemma 1

First, according to (30), x_{m_i} can be decomposed into the average of L independent Gaussian variables, as $x_{m_i} = \frac{1}{L} \sum_{\ell=1}^L x_{\ell, m_i}$, $i = 1, 2$, where the distribution of $x_{m_i, \ell}$ is

$$x_{\ell, m_i} \sim \mathcal{N}(\mu_{\ell, m_i}, \sigma_{m_i}^2), \quad 1 \leq \ell \leq L, \quad i = 1, 2. \quad (55)$$

Then, by substituting the above equation into the local elements defined in (31), we have

$$x_{k, m_i} = \frac{1}{L} \sum_{\ell=1}^L x_{\ell, m_i} + d_{k, m_i}, \quad 1 \leq k \leq K, \quad i = 1, 2. \quad (56)$$

It follows that

$$x_{k, m_i} = \frac{1}{L} \sum_{\ell=1}^L x_{\ell, k, m_i}, \quad 1 \leq k \leq K, \quad i = 1, 2. \quad (57)$$

where $x_{\ell, k, m_i} = x_{\ell, m_i} + d_{k, m_i}$. Next, by substituting the distributions of x_{ℓ, m_i} in (55) and the distribution of d_{k, m_i} in (32), the distribution of x_{ℓ, k, m_i} can be derived as

$$x_{\ell, k, m_i} \sim \mathcal{N}(\mu_{\ell, m_i}, \sigma_{m_i}^2 + \epsilon_k^2), \quad 1 \leq k \leq K, \quad \& \leq \ell \leq L, \quad \& i = 1, 2. \quad (58)$$

It follows that the distribution of x_{k, m_i} can be derived as

$$x_{k, m_i} \sim \frac{1}{L} \sum_{\ell=1}^L \mathcal{N}(\mu_{\ell, m_i}, \sigma_{m_i}^2 + \epsilon_k^2), \quad i = 1, 2, \quad \& 1 \leq k \leq K. \quad (59)$$

B. Proof of Lemma 2

First, for the received symbol \hat{s} in (19), the received noise can be derived as

$$\mathbf{f}^H \mathbf{n} = (\mathbf{f}_1 + j\mathbf{f}_2)^H (\mathbf{n}_1 + j\mathbf{n}_2) = \mathbf{f}_1^T \mathbf{n}_1 + \mathbf{f}_2^T \mathbf{n}_2 + j(\mathbf{f}_1^T \mathbf{n}_2 - \mathbf{f}_2^T \mathbf{n}_1), \quad (60)$$

where \mathbf{f}_1 and \mathbf{f}_2 are the real part and imaginary part of \mathbf{f} respectively, and \mathbf{n}_1 and \mathbf{n}_2 are the real part and imaginary part of the Gaussian noise \mathbf{n} respectively. Specifically, we have

$$\mathbf{n}_1 \sim \mathcal{N}\left(\mathbf{0}, \frac{\delta_0^2}{2} \mathbf{I}\right), \quad \mathbf{n}_2 \sim \mathcal{N}\left(\mathbf{0}, \frac{\delta_0^2}{2} \mathbf{I}\right) \quad (61)$$

where δ_0^2 is the noise variance. Then, for the real part of the received noise, its expectation and co-variance can be derived as

$$\mathbb{E}[\text{Re}(\mathbf{f}^H \mathbf{n})] = \mathbb{E}[\mathbf{f}_1^T \mathbf{n}_1 + \mathbf{f}_2^T \mathbf{n}_2] = \mathbf{0}, \quad (62)$$

and

$$\mathbb{C} = [\text{Re}(\mathbf{f}^H \mathbf{n})] = \mathbb{E}\left[(\mathbf{f}_1^T \mathbf{n}_1 + \mathbf{f}_2^T \mathbf{n}_2)(\mathbf{f}_1^T \mathbf{n}_1 + \mathbf{f}_2^T \mathbf{n}_2)^T\right] = \frac{\delta_0^2}{2} (\mathbf{f}_1^T \mathbf{f}_1 + \mathbf{f}_2^T \mathbf{f}_2). \quad (63)$$

That's to say,

$$\text{Re}(\mathbf{f}^H \mathbf{n}) \sim \mathcal{N} \left(\mathbf{0}, \frac{\delta_0^2}{2} (\mathbf{f}_1^T \mathbf{f}_1 + \mathbf{f}_2^T \mathbf{f}_2) \right). \quad (64)$$

Similarly, it can be derived that the imaginary part of the received noise has the same distribution:

$$\text{Im}(\mathbf{f}^H \mathbf{n}) \sim \mathcal{N} \left(\mathbf{0}, \frac{\delta_0^2}{2} (\mathbf{f}_1^T \mathbf{f}_1 + \mathbf{f}_2^T \mathbf{f}_2) \right). \quad (65)$$

By substituting the noise distributions in (64) and (65) into the global estimates in (20) and using the similar method in Appendix A, i.e., decomposing the local estimate x_{k,m_i} into the average of L independent Gaussia variables, the distributions of the global estimates can be derived as in (34).

C. Proof of Lemma 3

The Lagrange function of the problem in (39) can be written as

$$\begin{aligned} \mathcal{L} = & -\frac{2}{L(L-1)} \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \alpha_{\ell,\ell',m_i} + \sum_{k=1}^K \beta_k \left[c_k^2 - \hat{P}_k \mathbf{h}_k^H (\mathbf{f}_1 \mathbf{f}_1^T + \mathbf{f}_2 \mathbf{f}_2^T) \mathbf{h}_k \right] \\ & + \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell,\ell',m_i} \left[\alpha_{\ell,\ell',m_i} \hat{\sigma}_{m_i}^2 - (\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2 \right], \end{aligned} \quad (66)$$

where $(\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2$ is defined in (40) and $\hat{\sigma}_{m_i}^2$ is defined in (41). KKT conditions are necessary to achieve the optimal solution. Some useful KKT conditions are given below.

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{f}_1} &= -2 \sum_{k=1}^K \beta_k \hat{P}_k \mathbf{h}_k^H \mathbf{h}_k \mathbf{f}_1 + \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell,\ell',m_i} \alpha_{\ell,\ell',m_i} \delta_0^2 \mathbf{f}_1 = 0, \\ \frac{\partial \mathcal{L}}{\partial \mathbf{f}_2} &= -2 \sum_{k=1}^K \beta_k P_k \mathbf{h}_k^H \mathbf{h}_k \mathbf{f}_2 + \sum_{i=1}^2 \sum_{\ell'=1}^L \sum_{\ell < \ell'} \lambda_{\ell,\ell',m_i} \alpha_{\ell,\ell',m_i} \delta_0^2 \mathbf{f}_2 = 0. \end{aligned} \quad (67)$$

It is observed from the above equations that $\mathbf{f}_1 = \mathbf{f}_2$ won't influence the optimality of the problem.

D. Proof of Lemma 4

First, the second constraint in (39) can be equally written as

$$(\hat{\mu}_{\ell,m_i} - \hat{\mu}_{\ell',m_i})^2 \geq \alpha_{\ell,\ell',m_i} \hat{\sigma}_{m_i}^2, \quad \forall (\ell, \ell', m_i), \quad (68)$$

The reason is as follows. In (68), if the equality is not achieved, the value of α_{ℓ,ℓ',m_i} can be increased to make the objective function in (39) larger. In other words, it's necessary to achieve

equality for obtaining the optimal solution. Then, by substituting $(\hat{\mu}_{\ell, m_i} - \hat{\mu}_{\ell', m_i})^2$ in (40) and $\hat{\sigma}_{m_i}^2$ in (41) into (68), it can be derived as

$$\left(\sum_{k=1}^K c_k\right)^2 (\mu_{\ell, m_i} - \mu_{\ell', m_i})^2 \geq \alpha_{\ell, \ell', m_i} \left[\sigma_{m_i}^2 \left(\sum_{k=1}^K c_k\right)^2 + \sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}} \right], \quad \forall(\ell, \ell', m_i),$$

It follows that

$$\left(\sum_{k=1}^K c_k\right)^2 \left[\frac{(\mu_{\ell, m_i} - \mu_{\ell', m_i})^2}{\alpha_{\ell, \ell', m_i}} - \sigma_{m_i}^2 \right] \geq \sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}}, \quad \forall(\ell, \ell', m_i). \quad (69)$$

E. Proof of Lemma 5

It is straightforward that the objective function, c_k^2 , $R_k(\hat{\mathbf{f}})$, and $\sum_{k=1}^K c_k^2 \epsilon_k^2 + \delta_0^2 \hat{\mathbf{f}}^T \hat{\mathbf{f}} + \sigma_{m_i}^2 \left(\sum_{k=1}^K c_k\right)^2$ are convex and differentiable, since they are either linear or combination of quadratic functions. In the sequel, we show that $Q_{\ell, \ell', m_i}(\{c_k\}, \alpha_{\ell, \ell', m_i})$, are convex and differentiable. To begin with, $Q(\{c_k\}, \alpha_{\ell, \ell', m_i})$ can be linearly transformed from the convex function, say $f(x, y) = \frac{x^2}{y}$, $x > 0$, $y > 0$. As linear transformation preserves convexity, $Q(\{c_k\}, \alpha_{\ell, \ell', m_i})$ is convex and differentiable.

REFERENCES

- [1] K. B. Letaief, W. Chen, Y. Shi, J. Zhang, and Y.-J. A. Zhang, "The roadmap to 6G: AI empowered wireless networks," *IEEE Commun. Mag.*, vol. 57, no. 8, pp. 84–90, Aug. 2019.
- [2] K. B. Letaief, Y. Shi, J. Lu, and J. Lu, "Edge artificial intelligence for 6G: Vision, enabling technologies, and applications," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 1, pp. 5–36, Jan. 2022.
- [3] G. Zhu, D. Liu, Y. Du, C. You, J. Zhang, and K. Huang, "Toward an intelligent edge: Wireless communication meets machine learning," *IEEE Commun. Mag.*, vol. 58, no. 1, pp. 19–25, Jan. 2020.
- [4] Y. Shi, K. Yang, T. Jiang, J. Zhang, and K. B. Letaief, "Communication-efficient edge AI: Algorithms and systems," *IEEE Commun. Surv. & Tut.*, vol. 22, no. 4, pp. 2167–2191, Fourth quarter 2020.
- [5] D. Wen, K.-J. Jeon, and K. Huang, "Federated dropout—A simple approach for enabling federated learning on resource constrained devices," *IEEE Wireless Commun. Lett.*, vol. 11, no. 5, pp. 923–927, May 2022.
- [6] Y. Shi, K. Yang, Z. Yang, and Y. Zhou, *Mobile Edge Artificial Intelligence*. Academic Press, 2022. [Online]. Available: <https://www.sciencedirect.com/science/article/pii/B9780128238172000048>
- [7] D. Wen, M. Bennis, and K. Huang, "Joint parameter-and-bandwidth allocation for improving the efficiency of partitioned edge learning," *IEEE Trans. Wireless Commun.*, vol. 68, pp. 2128–2142, Dec. 2020.
- [8] K. Yang, Y. Shi, W. Yu, and Z. Ding, "Energy-efficient processing and robust wireless cooperative transmission for edge inference," *IEEE Internet Things J.*, vol. 7, no. 10, pp. 9456–9470, Oct. 2020.
- [9] P. Liu, G. Zhu, W. Jiang, W. Luo, J. Xu, and S. Cui, "Vertical federated edge learning with distributed integrated sensing and communication," *IEEE Commun. Lett.*, early access, Jun. 2022.

- [10] D. Wen, P. Liu, G. Zhu, Y. Shi, S. Cui, and Y. C. Eldar, "Task-oriented sensing, computation, and communication integration for multi-device edge AI," [Online]. Available: <https://arxiv.org/pdf/2207.00969.pdf>, Jul. 2022.
- [11] E. C. Strinati and S. Barbarossa, "6G networks: Beyond Shannon towards semantic and goal-oriented communications," *Comput. Netw.*, vol. 190, May 2021.
- [12] D. Ma, N. Shlezinger, T. Huang, Y. Liu, and Y. C. Eldar, "FRaC: FMCW-based joint radar-communications system via index modulation," *IEEE J. Sel. Topics Signal Process.*, vol. 15, no. 6, pp. 1348–1364, Nov. 2021.
- [13] Q. Lan, D. Wen, Z. Zhang, Q. Zeng, X. Chen, P. Popovski, and K. Huang, "What is semantic communication? A view on conveying meaning in the era of machine intelligence," *J. Commun. Inf. Netw.*, vol. 6, no. 4, pp. 336–371, 2021.
- [14] H. Xie, Z. Qin, X. Tao, and K. B. Letaief, "Task-oriented multi-user semantic communications," [Online]. Available: <https://arxiv.org/pdf/2112.10255.pdf>, 2021.
- [15] Z. Zhou, X. Chen, E. Li, L. Zeng, K. Luo, and J. Zhang, "Edge intelligence: Paving the last mile of artificial intelligence with edge computing," *Proceedings of the IEEE*, vol. 107, no. 8, pp. 1738–1762, 2019.
- [16] J. Soifer, J. Li, M. Li, J. Zhu, Y. Li, Y. He, E. Zheng, A. Oltean, M. Mosyak, C. Barnes, T. Liu, and J. Wang, "Deep learning inference service at Microsoft," in *USENIX Conf. Oper. Mach. Learn. (OpML)*, Santa Clara, CA, May 2019, pp. 15–17.
- [17] S. Jang, B. Kostadinov, and D. Lee, "Microservice-based edge device architecture for video analytics," in *2021 IEEE/ACM Symp. Edge Comput. (SEC)*, Los Alamitos, CA, Dec. 2021, pp. 165–177.
- [18] J. Shao and J. Zhang, "Communication-computation trade-off in resource-constrained edge inference," *IEEE Commun. Mag.*, vol. 58, no. 12, pp. 20–26, Jan. 2020.
- [19] X. Huang and S. Zhou, "Dynamic compression ratio selection for edge inference systems with hard deadlines," *IEEE Internet Things J.*, vol. 7, no. 9, pp. 8800–8810, Sep. 2020.
- [20] J. Shao, Y. Mao, and J. Zhang, "Learning task-oriented communication for edge inference: An information bottleneck approach," *IEEE J. Sel. Areas Commun.*, vol. 40, no. 1, pp. 197–211, Jan. 2022.
- [21] W. Shi, S. Zhou, and Z. Niu, "Device scheduling with fast convergence for wireless federated learning," in *IEEE Int. Conf. Commun. (ICC)*, Jun. 2020, pp. 1–6.
- [22] Q. Lan, Q. Zeng, P. Popovski, D. Gündüz, and K. Huang, "Progressive feature transmission for split inference at the wireless edge," [Online]. Available: <https://arxiv.org/abs/2112.07244>, 2021.
- [23] E. Li, L. Zeng, Z. Zhou, and X. Chen, "Edge AI: On-demand accelerating deep neural network inference via edge computing," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 447–457, Jan. 2019.
- [24] Z. Liu, Q. Lan, and K. Huang, "Resource allocation for multiuser edge inference with batching and early exiting (extended version)," [Online]. Available: <https://arxiv.org/abs/2204.05223>, Apr. 2022.
- [25] M. Jankowski, D. Gündüz, and K. Mikolajczyk, "Joint device-edge inference over wireless links with pruning," in *IEEE Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*. IEEE, May 2020, pp. 1–5.
- [26] Z. Yi, R. Zhang, B. Xu, Y. Chen, L. Zhu, F. Li, G. Yang, and Y. Luo, "A wide-angle beam scanning antenna in e-plane for k-band radar sensor," *IEEE Access*, vol. 7, pp. 171 684–171 690, Nov. 2019.
- [27] Z.-Q. Luo, "Universal decentralized estimation in a bandwidth constrained sensor network," *IEEE Trans. Inf. Theory*, vol. 51, no. 6, pp. 2210–2219, Jun. 2005.
- [28] J.-J. Xiao, S. Cui, Z.-Q. Luo, and A. J. Goldsmith, "Power scheduling of universal decentralized estimation in sensor networks," *IEEE Transactions on Signal Processing*, vol. 54, no. 2, pp. 413–422, 2006.
- [29] G. Zhu, J. Xu, K. Huang, and S. Cui, "Over-the-air computing for wireless data aggregation in massive IoT," *IEEE Wireless Commun.*, vol. 28, no. 4, pp. 57–65, 2021.

- [30] L. Chen, N. Zhao, Y. Chen, F. R. Yu, and G. Wei, "Over-the-air computation for iot networks: Computing multiple functions with antenna arrays," *IEEE Internet Things J.*, vol. 5, no. 6, pp. 5296–5306, 2018.
- [31] G. Zhu and K. Huang, "MIMO over-the-air computation for high-mobility multimodal sensing," *IEEE Internet Things J.*, vol. 6, no. 4, pp. 6089–6103, Sep. 2019.
- [32] X. Li, G. Zhu, Y. Gong, and K. Huang, "Wirelessly powered data aggregation for IoT via over-the-air function computation: Beamforming and power control," *IEEE Trans. Wireless Commun.*, vol. 18, no. 7, pp. 3437–3452, 2019.
- [33] D. Wen, G. Zhu, and K. Huang, "Reduced-dimension design of mimo over-the-air computing for data aggregation in clustered iot networks," *IEEE Trans. Wireless Commun.*, vol. 18, no. 11, pp. 5255–5268, 2019.
- [34] X. Zhai, X. Chen, J. Xu, and D. W. Kwan Ng, "Hybrid beamforming for massive mimo over-the-air computation," *IEEE Trans. Commun.*, vol. 69, no. 4, pp. 2737–2751, 2021.
- [35] X. Cao, G. Zhu, J. Xu, and K. Huang, "Optimized power control for over-the-air computation in fading channels," *IEEE Trans. on Wireless Commun.*, vol. 19, no. 11, pp. 7498–7513, 2020.
- [36] W. Liu, X. Zang, Y. Li, and B. Vucetic, "Over-the-air computation systems: Optimization, analysis and scaling laws," *IEEE Trans. Wireless Commun.*, vol. 19, no. 8, pp. 5488–5502, 2020.
- [37] M. Fu, Y. Zhou, Y. Shi, T. Wang, and W. Chen, "UAV-assisted over-the-air computation," in *IEEE Int. Conf. Commun.*, early access, Jun. 2021, pp. 1–6.
- [38] G. Zhu, Y. Wang, and K. Huang, "Broadband analog aggregation for low-latency federated edge learning," *IEEE Trans. Wireless Commun.*, vol. 19, no. 1, pp. 491–506, 2020.
- [39] K. Yang, T. Jiang, Y. Shi, and Z. Ding, "Federated learning via over-the-air computation," *IEEE Trans. Wireless Commun.*, vol. 19, no. 3, pp. 2022–2035, 2020.
- [40] M. M. Amiri and D. Gündüz, "Federated learning over wireless fading channels," *IEEE Trans. Wireless Commun.*, vol. 19, no. 5, pp. 3546–3557, 2020.
- [41] Y. Sun, S. Zhou, Z. Niu, and D. Gündüz, "Dynamic scheduling for over-the-air federated edge learning with energy constraints," *IEEE J. Sel. Areas in Commun.*, vol. 40, no. 1, pp. 227–242, 2022.
- [42] D. Liu and O. Simeone, "Privacy for free: Wireless federated learning via uncoded transmission with adaptive power control," *IEEE J. Sel. Areas Commun.*, vol. 39, no. 1, pp. 170–185, 2021.
- [43] S. F. Yilmaz, B. Hasircioglu, and D. Gunduz, "Over-the-air ensemble inference with model privacy," [Online]. Available: <https://arxiv.org/pdf/2202.03129.pdf>, May 2022.
- [44] S. Kullback and R. A. Leibler, "On information and sufficiency," *The annals of mathematical statistics*, vol. 22, no. 1, pp. 79–86, 1951.
- [45] M. Razaviyayn, "Successive convex approximation: Analysis and applications," Ph.D. dissertation, University of Minnesota, 2014.
- [46] G. Li, S. Wang, J. Li, R. Wang, X. Peng, and T. X. Han, "Wireless sensing with deep spectrogram network and primitive based autoregressive hybrid channel model," in *IEEE Int. Workshop Signal Process. Adv. Wireless Commun. (SPAWC)*, Sep. 2021, pp. 481–485.
- [47] M. Grant, S. Boyd, and Y. Ye, "CVX users' guide," [Online]. Available: <http://www.stanford.edu/boyd/software.html>, 2009.
- [48] MathWorks, "Pedestrian and bicyclist classification using deep learning," [Online]. Available: <https://ww2.mathworks.cn/help/radar/ug/pedestrian-and-bicyclist-classification-using-deep-learning.html>, 2022.