Decentralized Over-the-Air Computation for Edge AI Inference with Integrated Sensing and Communication

Zeming Zhuang, Dingzhu Wen and Yuanming Shi

School of Information Science and Technology, Shanghai Tech University









Background

System Model

Problem Formulation

Joint Design Algorithm

Simulation

Background - Edge Al Inference



Network Edge Intelligent Services (e.g., auto-driving and XR)

□ Fast response to the requirements

□ Intelligent decisions depends on the deployment of well-trained AI models (Edge Inference)



Background - Inference Paradigms



On-device Inference

- Execute AI model on device
- High computation overhead
- On-server Inference
 - Communication bottleneck
 - Data privacy

□Split Inference (Device-Edge Co-inference)

- Feature extraction on device
 - Low-dimensional feature vector
 - Preserve data privacy
- Computation offloading



What if the server is not always available for all devices simaltaneously? Decentralized co-inference!

System Model – Network Model



□We consider a system containing K ISCC devices each equipped with a full-duplex transceiver of N_t transmit antennas and N_r receive antennas.

Each device has its own inference task requiring features from all other devices.



Sensing Model

上海科技大 ShanghaiTech Univers

Latency T



Feature Model



Latency T



*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.

Communication Model



 \Box The feature received through over-the-air computation by device $k \in \mathcal{K}$

$$\mathbf{y}_k(m) = \sum_{j \in \mathcal{K}, j \neq k} \mathbf{H}_{j,k} \mathbf{p}_{j,m} x_j(m) + \mathbf{H}_{k,k} \mathbf{p}_{k,m} x_k(m) + \mathbf{w}_k(m)$$

 $\Box \mathbf{p}_k$ is the multicast beamformer, when k' = k, $\mathbf{H}_{k,k}$ is the Self-Interference channel gain

lacksquare By designing \mathbf{p}_k , SI channel is exploited to aggregate \mathbf{x}_k with other \mathbf{x}'_k

Then a receive beamformer $f_{k,m}$ is applied at device k

$$\hat{x}_k(m) = \mathbf{f}_{k,m}^H \cdot \left(\sum_{j \in \mathcal{K}} \mathbf{H}_{j,k} \mathbf{p}_{j,m} x_j(m) + \mathbf{w}_k(m) \right)$$

Then the aggregated features are used for inference task later.



Superposition of

Over-the-air Computation





Distribution of Extracted Feature Elements

 \succ For the *m*-th element

$$x_k(m) \sim \frac{1}{L} \sum_{\ell=1}^{L} \mathcal{N}\left(\mu_{\ell,m}, \sigma_{k,m}^2 + \sigma_{s,k}^2 + \frac{\sigma_r^2}{P_{s,k}}\right)$$

Aggregated Feature Elements

$$\hat{\mu}_{\ell,k,m} = \mu_{\ell,m} \cdot \left(\sum_{j \in \mathcal{K}} \mathbf{f}_{k,m}^H \mathbf{H}_{j,k} \mathbf{p}_{j,m} \right)$$
$$\hat{\sigma}_{k,m}^2 = \sigma_{k,m}^2 \cdot \left(\sum_{j \in \mathcal{K}} \mathbf{f}_{k,m}^H \mathbf{H}_{j,k} \mathbf{p}_{j,m} \right)^2 + \sum_{j \in \mathcal{K}} \left(\mathbf{f}_{k,m}^H \mathbf{H}_{j,k} \mathbf{p}_{j,m} \right)^2 \left(\sigma_{s,j}^2 + \frac{\sigma_r^2}{P_{s,j}} \right) + N_0 \mathbf{f}_{k,m}^H \mathbf{f}_{k,m}$$

Designing Metric

>MMSE ignores the heterogeneous distortion sensitivities of different elements

 \geq Inference accuracy is hard to calculate instantaneously

上海科技大学 ShanghaiTech University

Task-oriented Metric

Discriminant gain is defined based on the KL divergence between two Gaussian dist. classes $G_{\ell,\ell'}(\hat{x}_k(m)) = D_{KL} [f_\ell(\hat{x}_k(m)) || f_{\ell'}(\hat{x}_k(m))]$ $+ D_{KL} [f_{\ell'}(\hat{x}_k(m)) || f_\ell(\hat{x}_k(m))]$

Normally discriminant gain of every pair of classes are summed up

$$G(\mathbf{x}) = \sum_{i=1}^{-} \sum_{j < i} G_{ij}(\mathbf{x})$$

Shortage of avg. DG: unbalance

Class 2 & 3 are too close and far away from class I

Solution: Minimum Discriminant Gain $G_{\min}(\hat{\mathbf{x}}_k) = \min_{1 \le \ell \ne \ell' \le L} \sum_{m=1}^M G_{\ell,\ell'}(\hat{x}_k(m))$



Problem Formulation



Decision variable: Sensing power $P_{s,k}$

Multicast beamformer $\mathbf{p}_{k,m}$

Receive beamformer $f_{k,m}$

Long-term Power constraint: $P_{s,k}T_s + E_p + T_c \sum_{m=1}^M \|\mathbf{p}_{k,m}\|^2 X_k(m) \le E_k$

Unit receive beamforming constraint: $\|\mathbf{f}_{k,m}\|^2 = 1$

Joint design of Sensing Power and Beamforming



\Box Introduce α_k as a slack variable and extend the feasible region of constraints

$$\max_{\substack{\{P_{s,k}\},\{\mathbf{f}_{k,m}\},\{\mathbf{p}_{k',m}\},\\\{\alpha_k\},\{\beta_{\ell,\ell',k,m}\}}} \sum_{k=1}^{n} \eta_k \alpha_k,$$
(1)

s.t.
$$P_{s,k}T_{s,k} + E_{p,k} + T_c \sum_{m=1}^{M} \|\mathbf{p}_{k,m}\|^2 \le E_k, \quad \|\mathbf{f}_{k,m}\|^2 \le 1,$$
(2)

$$\alpha_k - \sum_{m=1}^M \beta_{\ell,\ell',k,m} \le 0,\tag{3}$$

$$\frac{\left(\mu_{\ell,k,m}-\mu_{\ell',k,m}\right)^2}{\beta_{\ell,\ell',k,m}} \left(\sum_{k'=1}^K \mathbf{f}_{k,m}^H \mathbf{H}_{k'k,m} \mathbf{p}_{k',m}\right)^2 \ge Z_{\ell,\ell',k,m}\left(\{P_{s,k'}\},\{\mathbf{p}_{k',m}\},\mathbf{f}_{k,m}\right) \tag{4}$$

Approximate the last constraint with Taylor expansion

$$\begin{array}{ll} \mathbf{P4} & \max_{\substack{P_{s,k}, \mathbf{f}_{k,m}, \mathbf{p}_{j,m}, \\ \alpha_{k}, \beta_{\ell,\ell',k,m}}} & \sum_{k=1}^{\kappa} \eta_{k} \alpha_{k}, \\ & \\ & \text{s.t.} & (2), (3), \\ & \\ & \\ \mathbf{Fernativelv} & Z_{\ell,\ell',k,m} \left(\{P_{s,j}\}, \{\mathbf{p}_{j,m}\}, \mathbf{f}_{k,m}\right) - Q_{\ell,\ell',k,m}^{[t]} \left(\{\mathbf{p}_{j,m}\}, \mathbf{f}_{k,m}, \beta_{\ell,\ell',k,m}\right) \leq 0 \end{array}$$

Then alternatively

P5: Fix
$$\mathbf{p}_{j,m}$$
 , $P_{s,k}$, update $\mathbf{f}_{k,m}$

 $ightarrow \mathbf{P6}$: Fix $\mathbf{f}_{k,m}$, update $\mathbf{p}_{j,m}$, $P_{s,k}$



Algorithm 1 Joint Sensing Power Allocation, Multicast and Receive Beamforming Design

Input: Channel gain $\{\mathbf{H}_{j,k}\}$, Device energy E_k **Output:** $\{P_{s,k}^*\}, \{\mathbf{p}_{j,m}^*\}, \{\mathbf{f}_{k,m}^*\}$

- 1: Initialize t = 0, $\{P_{s,k}^{[0]}\}$, $\{\mathbf{p}_{j,m}^{[0]}\}$, $\{\mathbf{f}_{k,m}^{[0]}\}$, $\alpha_k^{[0]}$, $\beta_{\ell,\ell',k,m}^{[0]}$ in feasible region of **P3**;
- 2: Initialize function $Q_{\ell,\ell',k,m}^{[0]}\left(\{\mathbf{p}_{j,m}^{[0]}\},\mathbf{f}_{k,m}^{[0]},\beta_{\ell,\ell',k,m}^{[0]}\right);$

3: repeat

- 4: Derive the approximated problem P4;
- 5: repeat
- 6: Update $\mathbf{f}_{k,m}$ by solving **P5**;
- 7: Update $\mathbf{p}_{j,m}$, $P_{s,k}$ by solving **P6**;
- 8: **until** convergence
- 9: Update function $Q_{\ell,\ell',k,m}^{[t+1]}\left(\{\mathbf{p}_{j,m}^{[t+1]}\},\mathbf{f}_{k,m}^{[t+1]},\beta_{\ell,\ell',k,m}^{[t+1]}\right);$
- 10: $t \leftarrow t + 1;$
- 11: **until** convergence
- 12: Optimal solution $P_{s,k}^* \leftarrow P_{s,k}^{[t]}, c_{k,m}^* \leftarrow c_{k,m}^{[t]}, \mathbf{f}_m^* \leftarrow \mathbf{f}_m^{[t]};$



Multi-task Radar Sensing of human activities*

□ 3 edge devices each equipped with 8 receive antennas and 8 transmit antennas, retrieving wide-view sensing data of same target

Tasks

- > Activity recognition and Person identification (Gender and Age)
- > 4 activities: sit down, stand up, pick up an object, drink water
- I 500 samples of 54 volunteers

Schemes

- Joint designing ISCC scheme (proposed)
- > Fixed beamforming vectors and sensing powers







Simulation

Simulation – Accuracy vs. Device Energy



- It shows that a higher upper limit of energy leads to better inference accuracy, because a higher device energy threshold indicates that the devices can allocate more power to suppress the distortion and noise.
- Moreover, our proposal outperforms the baseline since the sensing power and beamforming vectors are not adjusted to different feature elements in the baseline scheme.



Simulation – Accuracy vs. Feature Dimension



Accuracies of all three tasks increase as the feature dimension increases. This is because different feature dimensions are orthogonal and independent from the algorithm of PCA.

It follows that more feature dimensions can keep more data about the sensing target, which improves the performance of inference tasks.







We proposed a decentralized ISCC framework with the exploitation of the self-interference channel in full-duplex communication.

We jointly designed the sensing power, multicast and receive beamformer under the criterion of maximizing the weighted sum of minimum discriminant gain in different tasks.

Simulation results demonstrate that the proposed scheme can achieve higher inference accuracy than the benchmarks under various limitations of energy and feature dimensions.





