# Communication-Efficient Multimodal Federated Learning: Joint Modality and Client Selection

Liangqi Yuan, *Graduate Student Member, IEEE,* Dong-Jun Han, *Member, IEEE,* Su Wang, Devesh Upadhyay, *Senior Member, IEEE* and Christopher G. Brinton, *Senior Member, IEEE*

**Abstract**—Multimodal federated learning (FL) aims to enrich model training in FL settings where clients are collecting measurements across multiple modalities. However, key challenges to multimodal FL remain unaddressed, particularly in heterogeneous network settings where: (i) the set of modalities collected by each client will be diverse, and (ii) communication limitations prevent clients from uploading all their locally trained modality models to the server. In this paper, we propose multimodal Federated learning with joint Modality and Client selection (mmFedMC), a new FL methodology that can tackle the above-mentioned challenges in multimodal settings. The joint selection algorithm incorporates two main components: (a) A modality selection methodology for each client, which weighs (i) the impact of the modality, gauged by Shapley value analysis, (ii) the modality model size as a gauge of communication overhead, against (iii) the frequency of modality model updates, denoted recency, to enhance generalizability. (b) A client selection strategy for the server based on the local loss of modality model at each client. Experiments on five real-world datasets demonstrate the ability of mmFedMC to achieve comparable accuracy to several baselines while reducing the communication overhead by over 20x. A demo video of our methodology is available at https://liangqiy.com/mmfedmc/.

**Index Terms**—Multimodal federated learning, Data fusion, Internet of things, Edge computing, Communication efficiency.

✦

## 1 INTRODUCTION

Federated learning (FL) is a distributed machine learning (ML) approach in which users collaboratively train ML models through sharing model parameters rather than raw measurements [2], [3]. The FL approach establishes a federation of learners, each updating their model based on local data. Subsequently, these locally learned parameters are uploaded to a central server or shared with other clients. The local model within each client is then updated by integrating an aggregation of these parameters, ensuring that each model benefits from the collective learning experience [4]–[6]. Since Internet of Things (IoT) applications, such as smartphones, robots, unmanned aerial vehicles (UAVs), etc., are frequently equipped with multimodal sensors and rely on their performance and robustness, there has been increasing interest in multimodal federated learning (mmFL) frameworks [7]–[9]. For example, consider a set of connected and automated vehicles (CAVs) with various sensors such as cameras, LiDAR, and Radar [10]–[12]. These multimodal sensors enable control decisions in various driving scenarios, including varying weather conditions and fields of view. CAVs could benefit greatly be leveraging methods such as mmFL to collaboratively learn ML models across vehicles, such as on-board prognostics and diagnostics to estimate the vehicle's state of health [13].
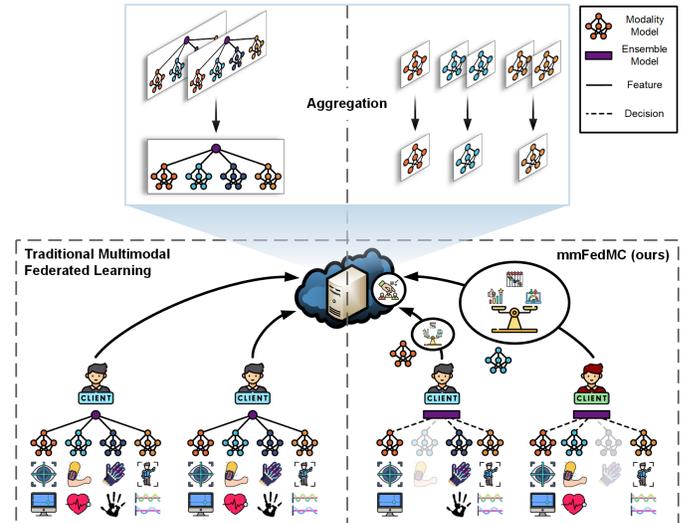


Fig. 1. Schematic representation of traditional multimodal federated learning vs. the proposed mmFedMC.

Multimodal fusion is a process of combining information from two or more modalities to make predictions, effectively enhancing the performance and robustness of ML models by leveraging cross modality interactions that provide richer representation [14]–[19]. Fusion methods are generally categorized as early, late, and intermediate, depending on the network layer where the modalities are merged. As the names suggest, early fusion implies fusion at the data-level, late fusion or decision-level fusion occurs at the output, and intermediate fusion leverages intermediate layers for the fusion of feature representations [20]–[23]. Generally, data-level fusion relies on the alignment of dimensions across different data modalities to concatenate

- *L. Yuan, D.-J. Han, and C. G. Brinton are with the School of Electrical and Computer Engineering, Purdue University, West Lafayette, 47907, USA. E-mail: {liangqiy, han762, cgb}@purdue.edu.*
- *S. Wang is with the School of Electrical and Computer Engineering, Princeton University, NJ, 08540, USA. E-mail: hw5731@princeton.edu.*
- *D. Upadhyay is with the Saab Inc., East Syracuse, 13057, USA. E-mail: deveshu@gmail.com.*

their channels; for example, an 3D RGB image can be combined with an 1D depth image to form an 4D RGB-D image for input into an ML model. Feature-level fusion combines features or outputs from intermediate layers of the various modalities. This allows for more nuanced interactions between modalities, but presents the challenge of a large search space within the network topologies. Decision-level fusion offers greater flexibility as it does not require concatenating feature vectors within the ML model. Instead, it aggregates decisions post-output from various models via some fusion decision logic to form a final prediction. However, more refined decision representations, such as those in decision-level fusion, may result in the loss of critical information during the fusion process, potentially limiting the model's ability to make holistic inferences. This trade-off between information retention and fusion sophistication is a key consideration in the design of multimodal learning systems.

## 1.1 Motivation

In the context of the IoT, where devices measure a diverse set of modalities, most existing mmFL frameworks call for massive parallel processing of extensive sensor data streams [24], [25]. They utilize multimodal fusion to boost model performance, especially in scenarios where heterogeneous clients lack one or more modalities. However, IoT devices often possess communication constraints due to bandwidth limitations, location, and varying capabilities. This has also been a key bottleneck in conventional (single modality) FL [26]–[29]. Thus, there is still a need to devise strategies to reduce communication overhead and improve learning efficiency within the mmFL paradigm. We summarize the existing gaps and challenges in the current mmFL frameworks as follows.

(i) **Client Heterogeneity** is commonly regarded as statistical heterogeneity, where each client's data distribution is non-independently and identically distributed (non-IID), with some clients potentially lacking samples from certain categories and thereby having differing densities. However, in real-world applications, especially in mmFL settings, client heterogeneity is more diverse, manifesting in various forms such as individual, group, and system differences [30]–[32]. Individual heterogeneity is seen in minor feature variations (e.g., walking frequencies), while group heterogeneity arises from physiological and background differences (e.g., left-hander vs. right-hander groups). System heterogeneity is influenced by external factors, such as the age of data collection devices. Moreover, in mmFL, some clients may lack certain modalities (e.g., some CAVs are equipped with LiDAR, while others rely solely on vision-based systems for autonomous driving). These heterogeneities are increasingly important considerations in mmFL scenarios.

(ii) **Communication Efficiency** is influenced not only by the balance between performance and communication overhead, but also by the varying data sizes across different modalities and data types, leading to differences in modality model sizes. Furthermore, the complexity of information patterns inherent in different data modalities affects the ease of learning. For instance, compared to high-resolution image data, simpler data modalities like time-series radar data, despite potentially lower accuracy, may require only simple ML models for effective training and recognition due to their smaller data size and straightforward pattern recognition.

(iii) **Impact of Modality and Client** in mmFL are critical yet under explored areas, particularly in the context of communication costs and client participation. It remains an open question as to which modality or combination of modalities should be prioritized in the predictions, considering their information richness. Moreover, the extent of client participation, influenced by factors such as data availability and quality, also plays a vital role in determining the effectiveness of different modalities within the mmFL framework.

Motivated by these challenges, we aim to answer the following key questions:

> *In resource-constrained and heterogeneous mmFL settings, (i) how should each client evaluate and select the best modalities, and (ii) how should the server determine the best clients to upload the models from so as to achieve the best trade-off between performance and communication?*

## 1.2 Overview and Contribution

We propose multimodal Federated learning with joint Modality and Client selection (mmFedMC), an mmFL methodology tailored for clients with heterogeneously absent modalities. An overview of our proposed method and comparison with traditional mmFL are given in Fig. 1. We propose a decision-level fusion approach, where predictions from *global modality models* are used as input to the individual *local ensemble model* in each client. This allows for the independent deployment of the modality models in various application scenarios, accommodating situations where a client may possess a different set of modalities. Furthermore, we consider classical ML models, such as Random Forest, for our ensemble instead of the complex neural networks typically used in traditional mmFL, resulting in reduced computational overhead and easily interpretable modality impacts. While modality models are uploaded to the server for aggregation to enhance generalization, each client retains the ensemble model locally to improve personalization and confidence. Retaining the ensemble model local to the client minimizes the risk of leakage of client sensitive information to the server either directly or indirectly via server side inferencing.

Building on this foundation, we introduce the concept of *joint modality and client selection*. First, *selective modality communication* acknowledges that clients may not always have the capacity or necessity to upload all modality models (e.g., resource-constrained IoT applications). This approach is based on factors such as modality performance or its impact on the final decision, communication overhead, and recency. To evaluate the impact of each modality, we propose employing Shapley values [33]–[35], measured on the ensemble model, to quantify their respective performances. Communication overheads are determined by quantifying the size of modality models, while recency is gauged by

tracking the upload history of modalities. Complementing this, *selective client uploading* is proposed to further reduce communication overheads and address the effects of client heterogeneity on the global model. This involves ranking each client based on the local loss of modality model, thereby optimizing both the efficiency and effectiveness of our FL framework.

In this paper, we present the following contributions:

- **Multimodal Federated Learning with Decision-level Fusion (Sec. 3.1).** We propose a framework that divides traditional holistic fusion approaches into separate modality and ensemble models. Modality models are uploaded to a server to enhance generalizability, while individual ensemble model is retained locally to further promote personalization and strengthen confidence. This bifurcated approach of training modality and ensemble models enables modular functionality of modalities and naturally accommodates client heterogeneity and missing modality scenarios.
- **Communication-Efficient Joint Modality and Client Selection (Sec. 3.4, Sec. 3.5).** We propose employing a Shapley component to quantify the impact of modality models, along with assessing their sizes to estimate communication overhead. Additionally, a recency term is introduced to track the frequency of uploads, facilitating selective modality communication. Furthermore, we implement selective client uploading, which leverages the local loss of modality model to quantify client heterogeneity. This joint modality and client selection strategy significantly reduces communication costs while maintaining model performance, thereby enhancing learning efficiency.
- **Five Real-World Experiments (Sec. 4.3).** Our proposed mmFedMC is evaluated for performance and communication overhead against four baseline methods across five heterogeneous multimodal datasets, including wearable sensors, healthcare, language, and satellite datasets. Experiment results highlight the superior performance of mmFedMC, achieving comparable accuracy while incurring less than 25% of the communication overhead compared to baseline methods.
- **Analytics on Modality Impact (Sec. 4.4).** We offer analyses on modality impact within the FL process. Utilizing Shapley values, we illustrate the interplay among modalities, revealing the dynamic impact of each modality on the final decision-making process within scenarios that take into account communication overhead and recency.

This paper is an extension of our previous paper [1]. Building upon the foundation laid by [1], this paper introduces the following key contributions: (i) For modality selection, we incorporate a new metric, recency, to prevent overemphasis on certain modalities and maintain generalization. (ii) We introduce a client selection strategy tailored to mmFL that synergistically optimizes communication overhead in conjunction with modality selection. (iii) We present a comprehensive set of experiments and analysis demonstrating that within the proposed mmFedMC framework, client selection based on the lower local loss outperforms the higher local loss approach utilized in [36]. (iv)

Beyond the previously experimented ActionSense dataset, we have added four new datasets, covering domains such as wearable sensors, healthcare, language, and satellite imagery, providing a holistic demonstration of the framework's applicability. (v) We extend our experimental suite with additional baselines, ablation studies, and extensive results presentation to thoroughly evaluate the performance of the mmFedMC framework.

## 1.3 Organization

The rest of this paper is organized as follows. Section 2 reviews related works, and Section 3 details our mmFedMC algorithm. We present our experiments and their corresponding results in Section 4. Some discussion of the experimental results are given in Section 5 and we summarize our conclusions and future works in Section 6.

## 2 RELATED WORKS

### 2.1 Multimodal Federated Learning

The real-world nature of multimodal data and the complementary advantages of multimodality have gradually drawn attention to mmFL. Recently, a variety of algorithms have been proposed to improve the performance of mmFLs based on different fusion techniques. Qi *et al.* [37] proposed a data-level fusion FL system tailored for the combination of wearable sensor signals and images for user fall detection. Xiong *et al.* [38] implemented a feature-level fusion approach with attention modules. Feng *et al.* [39] presented two decision-layer fusion strategies using concatenation and attention modules, respectively, to address three types of client heterogeneities: modality absence, label absence, and label errors. Outside of fusion strategies, Salehi *et al.* [40] incorporated the mmFL framework into CAVs and conducted real-world experiments. In their FLASH framework, clients randomly select and upload one of the three modality models or an ensemble model for aggregation. Chen *et al.* [41] integrated mmFL into decentralized FL, aiming to facilitate collaborative training within client networks without server support.

However, most current mmFL implementations rely solely on end-to-end models for FL. For example, in data-level fusion, it is common to see input concatenation either merely expand channels in images or increase feature quantity through concatenation. These end-to-end models face limitations not only in scenarios where heterogeneous clients may miss certain modalities, but also inherit the drawbacks of multimodality, such as larger model sizes, increased computational costs, and reduced interpretability. The limitations of data-level fusion encompass challenges such as the high dimensionality of the input space, which hampers scalability to new modalities. In addition, the integration of modalities with varying granularities, such as audio and video, poses significant difficulties in forming a uniform input vector, further complicating the fusion process. Therefore, our paper considers a decision-level fusion algorithm, training each single-modality model separately to address the shortcomings of end-to-end models, especially in missing-modality scenarios. Moreover, by employing a traditional ML model as the ensemble model, we offer

insightful demonstrations of modality impacts throughout the FL process, which further facilitates subsequent modality selection.

## 2.2 Modality Selection

From feature selection and sample filtering to modality selection, such feature engineering techniques have been widely applied in the ML domain to improve accuracy, reduce computational overhead, and eliminate outliers. Recent studies employing Shapley values have shown that, in data fusion, different input features exert varying impacts on model output [42], [43]. For example, Yuan *et al.* [44] used Shapley values to demonstrate that sensor placement significantly affects the fusion process. Yuan *et al.* [45] reduced over 400 frequency-corresponding received signal strength indicators (RSSI) to just five using the Shapley value, thereby markedly improving the sampling rate of sensor configurations. Compared to well-established feature selection methods, modality selection remains a nascent field, despite numerous studies employing single-modality or various modality combinations as baselines to validate the effectiveness of multimodal approaches [46]. However, there is no universally accurate method to guide proper modality selection, especially in dynamic contexts like FL.

## 2.3 Client Selection within Federated Learning

Due to the extensive heterogeneity among clients, encompassing individual, group, and system heterogeneity, as well as varying richness of information, information asymmetry, and concerns of security and privacy, it becomes essential for FL frameworks to incorporate static or dynamic client selection for uploading and aggregation [47]–[50]. Wang *et al.* [51] utilized a Graph Convolutional Network (GCN) to intelligently select clients in decentralized FL, addressing network heterogeneity and overlapping data distributions, thereby enhancing training accuracy and minimizing resource consumption. Yuan *et al.* [52] improved the effectiveness of aggregation and cybersecurity by selectively excluding clients with significant heterogeneity. Cho *et al.* [36] posited that client selection biased towards those with higher local losses leads to faster error convergence and conducted a convergence analysis to support this claim.

However, to date, no work has addressed client selection in the context of mmFL, particularly in scenarios like the proposed mmFedMC, which involves separate global modality models. In such cases, varying biases in different global modality models can influence the selection outcome of various modalities within clients. Notably, our extensive empirical experiment demonstrates that, due to the profound heterogeneity among clients, those with lower local losses can achieve faster error convergence. A case in point is in FL frameworks involving left-hander and right-hander groups, where the global loss exhibits a bimodal loss function (i.e., two local minima). Client selection focusing on higher local losses would cause the global model to oscillate between the two groups. In contrast, selecting clients with lower local losses ensures that the current modality model at least tends towards one group.

## 3 FORMULATION AND METHODOLOGY

### 3.1 Federated Learning and Decision-level Fusion

We assume that there are $K$ clients participating in the FL framework. Each client, denoted $k$, possesses a dataset $\mathbb{D}^k$ and a label set $Y^k$, comprising multimodal data represented as

$$\mathbb{D}^k = \{\mathcal{D}_1^k, \mathcal{D}_2^k, \ldots, \mathcal{D}_{M_k}^k\}, \tag{1}$$

where $\mathcal{D}_m^k$ denotes datasets corresponding to modality $m$, such as images, LiDAR, RF, etc., with $m = 1, 2, \ldots, M_k$ indicating specific data modality. We note that the heterogeneous client $k$ can accommodate a different number of data modalities, represented by $M_k$. Each client has a *modality model*, $\theta_m^k$, for every modality dataset $\mathcal{D}_m^k$, which is designed to capture the relationship between the input data and its corresponding labels. Therefore, each client possesses a set of models, represented as

$$\Theta^k = \{\theta_1^k, \theta_2^k, \ldots, \theta_{M_k}^k\}. \tag{2}$$

Each model, when trained on the dataset, yields a predicted label, denoted as

$$\hat{y}_m^k = \theta_m^k(\mathcal{D}_m^k), \tag{3}$$
$$\widehat{\mathbb{Y}}^k = \{\hat{y}_1^k, \hat{y}_2^k, \ldots, \hat{y}_{M_k}^k\}, \tag{4}$$

where $\hat{y}_m^k$ represents the predicted label generated from the model $\theta_m^k$ for the dataset $\mathcal{D}_m^k$, and $\widehat{\mathbb{Y}}^k$ is the collection of all the predicted labels for client $k$. Our goal is to fuse the outputs of all modality models at the decision layer through a post-processing *ensemble model* $\omega^k$, represented as

$$\widehat{Y}^k = \omega^k \left( \theta_1^k(\mathcal{D}_1^k), \theta_2^k(\mathcal{D}_2^k), \ldots, \theta_{M_k}^k(\mathcal{D}_{M_k}^k) \right)$$
$$= \omega^k \left( \hat{y}_1^k, \hat{y}_2^k, \ldots, \hat{y}_{M_k}^k \right)$$
$$= \omega^k \left( \widehat{\mathbb{Y}}^k \right), \tag{5}$$

where $\widehat{Y}^k$ denotes the predicted label set corresponding to the dataset $\mathbb{D}^k$.

### 3.2 Overview of Proposed mmFedMC Algorithm

The proposed mmFedMC is described in Fig. 2 and Algorithm 1. The primary objective of mmFedMC is the collaborative learning of the modality model $\theta_m$ for each modality $m$. For each client, individual and system heterogeneities (e.g., modality missing, noise, device errors, device malfunctions, etc.) are addressed through a personalized ensemble model $\omega^k$. Furthermore, an intrinsic objective of mmFedMC is to compensate for the communication constraints inherent in IoT edge devices by minimizing communication overhead, thus ensuring efficient learning efficacy for clients within the FL framework. In the following, we provide a detailed description of our mmFedMC algorithm.

### 3.3 Client Learning

For each data modality for every client, the local objective is to minimize the difference between the predicted and
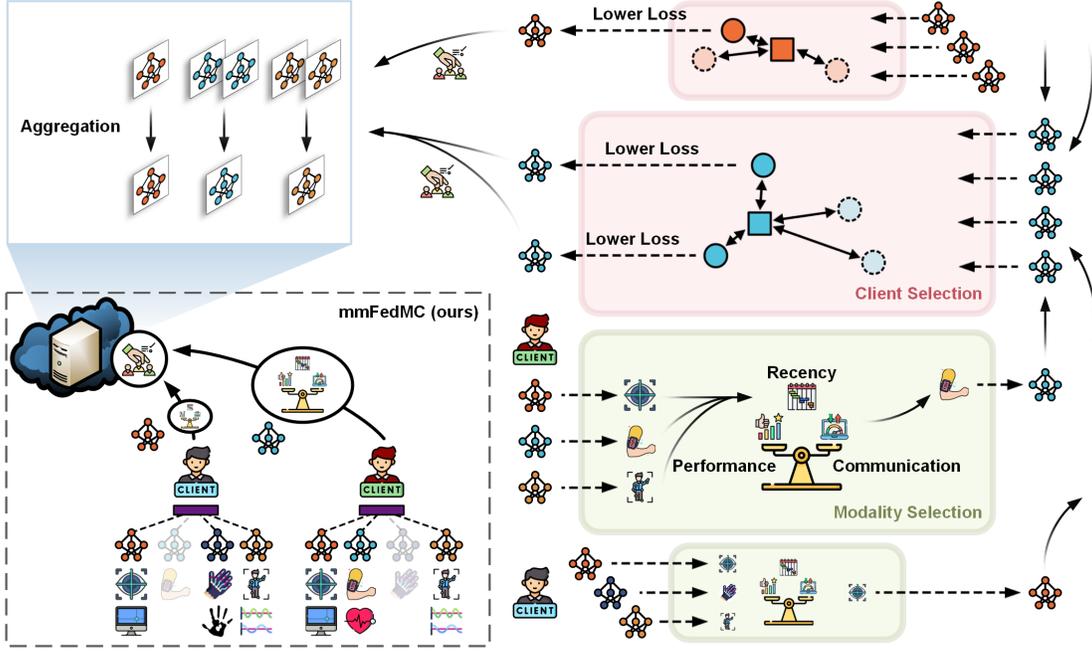
Fig. 2. System diagram of the proposed mmFedMC illustrating the process of **Modality Selection** and **Client Selection** as detailed in Algorithm 1.

the true labels. This objective can be achieved using various optimizers, such as stochastic gradient descent (SGD). Formally, we describe the optimization problem as

$$\min_{\theta_m^k} f(Y^k, \hat{y}_m^k), \tag{6}$$

where $f$ is a loss function that measures the discrepancy between the true label set and the predicted label set. For the ensemble model $\omega^k$, the learning objective is to minimize the discrepancy between the true label set and the predicted label sets across all modalities, as

$$\min_{\omega^k} f(Y^k, \widehat{Y}^k). \tag{7}$$

### 3.4 Modality Selection

Due to the resource constraints on edge devices serving as clients, they may not possess adequate storage capacity to house extensive multimodal data, computational capability to learn on multimodal datasets, or communication bandwidth to upload several models to the server. Thus, we introduce two metrics to assist clients to determine whether to upload their models to the server:

- **Shapley value** ($\varphi$) represents the impact of a modality model on the final prediction, where a higher value of $\varphi$ corresponds to a higher priority ↑.
- **Modality model size** ($|\theta|$) pertains to the communication overhead, where a lower value of $|\theta|$ corresponds to a higher priority ↓.
- **Recency** ($\mathcal{T}$) denotes the freshness or recentness of a client's modality model upload, where a higher value indicates that the modality has not been updated recently, thus it is given higher priority ↑.

**Shapley Value (Impact).** During each communication round, clients evaluate the impact of the models $\Theta^k$ on the outcomes utilizing interpretability techniques and choose to upload only a singular selected model $\theta^k$. We consider

using the Shapley value as an assessment to evaluate the relationship between input $\widehat{\mathbb{Y}}^k$ and output $\widehat{Y}^k$ of the ensemble model $\omega^k$:

$$\varphi_m^k = \sum_{\mathcal{Y} \subseteq \widehat{\mathbb{Y}}^k \setminus \{m\}} \frac{|\mathcal{Y}|!(|\widehat{\mathbb{Y}}^k| - |\mathcal{Y}| - 1)!}{|\widehat{\mathbb{Y}}^k|!} \left( \omega^k(\mathcal{Y} \cup \{m\}) - \omega^k(\mathcal{Y}) \right), \tag{8}$$

where $\varphi_m^k$ is the Shapley value of input modality $m$, $\mathcal{Y}$ is a subset of $\widehat{\mathbb{Y}}^k$ excluding modality $m$, and $\omega^k(\mathcal{Y})$ is the predicted value using only modalities in set $\mathcal{Y}$. For all modalities, we assess the magnitude of each Shapley value by taking its absolute value and construct the following set:

$$\Phi^k = \left\{ |\varphi_1^k|, |\varphi_2^k|, \ldots, |\varphi_{M_k}^k| \right\}. \tag{9}$$

**Modality Model Size (Communication Overhead).** Given modality models with parameters $\Theta^k = \{\theta_1^k, \theta_2^k, \ldots, \theta_{M_k}^k\}$, the communication overhead for each modality model is directly proportional to the model size given by

$$\bar{\Theta}^k = \left\{ |\theta_1^k|, |\theta_2^k|, \ldots, |\theta_{M_k}^k| \right\}. \tag{10}$$

**Recency.** Our primary aim with the recency metric is to encourage clients to update a specific modality model to the server. This is to ensure that the mmFedMC framework doesn't overly prioritize data modalities that are more easily obtainable, possess simpler data structures, or have conspicuous features, thus sidelining the diversity of other data modalities. To encapsulate the concept of recency in our model, we denote $\mathcal{T}_m^k$ as:

$$\mathcal{T}_m^k = t - t_m^k - 1, \tag{11}$$

where $t = 1, 2, \ldots, T$ represents the current communication round, and $t_m^k$ indicates the communication round during which the modality model from modality $m$ of client $k$ was last uploaded.

**Priority (Composite Score).** Considering the impact of the modality model, as quantified by the Shapley value, the

---

**Algorithm 1** mmFedMC: Multimodal Federated Learning with Joint **<u>Modality Selection</u>** and **<u>Client Selection</u>**

---

**Input:** Communication rounds ($T$), clients' dataset ($\mathbb{D}^k$), clients' label set ($Y^k$), local training epoch ($E$), initial models ($\theta_{m,0}^k$ & $\omega_0^k$), loss function ($f$), learning rate ($\eta$), modality model upload count ($\gamma$), performance and communication weights ($\alpha_s$ & $\alpha_c$)
**Output:** Generalized global modality models ($\theta_m$) and personalized local ensemble models for each client ($\omega^k$)

**# Global Iteration**
1: **for** $t = 0$ **to** $T - 1$ **do**

    **# Local Learning**
    1: **for** each client $k$ **in parallel do**
    2:     **for** each data modality $m$ **do**
    3:         Backpropagate the loss function and update the modality models $\theta_m^{k,t} \leftarrow \arg\min_{\theta_m^{k,t}} f(Y^{k,t}, \hat{y}_m^{k,t})$.
    4:         (Stage #1) Update the ensemble model $\omega^{k,t} \leftarrow \arg\min_{\omega^{k,t}} f(Y^{k,t}, \widehat{Y}^{k,t})$.
    5:     **end for**
    6: **end for**

    **# Modality Selection**
    1: Compute the impact of each modality $\Phi^{k,t}$ on the prediction using the Shapley values.         ▷ (8), (9)
    2: Compute the modality model size $\bar{\Theta}^k$.         ▷ (10)
    3: Compute the recency $\mathcal{T}_m^k$.         ▷ (11)
    4: **Clients select the modality models with top-$\gamma$ priority for uploading and update recency.**     ▷ (14), (15), (16)

    **# Client Selection & Server Aggregation**
    1: **Server selects top-$\delta$ clients with the lower local loss of modality model for uploading.**     ▷ (18), (19), (20)
    2: **for** each data modality $m$ **do**
    3:     Server calculates the weight $\beta_m^{k,t}$ for each client $k$ and aggregates modality model $\theta_m^t$.     ▷ (21), (22)
    4: **end for**

    **# Local Deploying**
    1: **for** each client $k$ **in parallel do**
    2:     **for** each data modality $m$ **do**
    3:         Deploy the downloaded global modality model $\theta_m^t$.
    4:         Re-calculate $\widehat{Y}^{k,t}$ based on the global modality model $\theta_m^t$.
    5:         (Stage #2) Update the ensemble model $\omega^{k,t} \leftarrow \arg\min_{\omega^{k,t}} f(Y^{k,t}, \widehat{Y}^{k,t})$.
    6:     **end for**
    7: **end for**

2:     Update the modality models for each modality $\theta_m \leftarrow \theta_m^T$.
3:     Update the ensemble models for each client $\omega^k \leftarrow \omega^{k,T}$.
4: **end for**

---

communication overhead as characterized by the modality model size, and the timeliness of model updates captured by the recency, we propose priority $P$ as a composite score. To derive the priority, we proceed with individual normalization for each criterion:

$$
\begin{cases}
\tilde{\varphi}_m^k = \dfrac{\varphi_m^k - \min(\Phi^k)}{\max(\Phi^k) - \min(\Phi^k)}, \\[2mm]
|\tilde{\theta}_m^k| = \dfrac{\theta_m^k - \min(\bar{\Theta}^k)}{\max(\bar{\Theta}^k) - \min(\bar{\Theta}^k)}, \text{ for } m = 1, 2, \ldots, M_k, \\[2mm]
\tilde{\mathcal{T}}_m^k = \dfrac{\mathcal{T}_m^k}{t},
\end{cases} \quad (12)
$$

where $\tilde{\varphi}_m^k$ represents the normalized Shapley Value, $|\tilde{\theta}_m^k|$ denotes the normalized communication overhead, and $\tilde{\mathcal{T}}_m^k$ indicates the normalized recency. With a focus on identifying modality models for server communication, we devise the priority $P_m^k$ for each modality and the corresponding set $\mathcal{P}^k$ to determine whether modality models should be sent to the server. They are formulated as:

$$
\begin{aligned}
P_m^k &= \alpha_s \times \tilde{\varphi}_m^k + \alpha_c \times (1 - |\tilde{\theta}_m^k|), \\
\mathcal{P}^k &= \left\{ P_1^k, P_2^k, \ldots, P_{M_k}^k \right\},
\end{aligned} \quad (13)
$$

where $\alpha_s$ and $\alpha_c$ are the predetermined metric weights, satisfying $\alpha_s + \alpha_c = 1$. Naturally, a modality with the maximal priority is considered optimal.

**Modality Selection.** To streamline our decision-making, we focus on modalities with scores among the top-$\gamma$ priority:

$$
\begin{aligned}
\mathcal{P}_\gamma^k &= \operatorname{top}\max_\gamma(\mathcal{P}^k) \\
&= \left\{ x : x \in \mathcal{P}^k \text{ and } |\mathcal{P}^k \cap \{y \mid y \geq x\}| \leq \gamma \right\}. \quad (14)
\end{aligned}
$$

Hence, the selected modality set of client $k$ becomes:

$$
\mathcal{M}_k = \left\{ m : P_m^k \in \mathcal{P}_\gamma^k, \text{ for } m = 1, 2, \ldots, M_k \right\}. \quad (15)
$$

With $\mathcal{M}_k$ determined, the set of modality models ready for communication to the server by the client $k$ can be defined as:

$$
\Theta_\gamma^k = \left\{ \theta_m^k : \theta_m^k \in \Theta^k \text{ and } m \in \mathcal{M}_k \right\}, \quad (16)
$$

where the set $\Theta_\gamma^k$ represents the modality models corresponding to the top-$\gamma$ priority from $\mathbb{S}^k$. Each client will upload a data packet with various details to the server for aggregation, including model parameters $\theta^k$, modality information $m$, the number of samples $|\mathcal{D}_m^k|$, among others. Likewise, upon downloading from the server, this information will also be retrieved. Note that only model $\theta^k$ will be uploaded/downloaded to/from the server. The ensemble model $\omega^k$ varies across clients, determined by the unique deployment scenarios of each client, such as geographical location, operational duration, external interference, etc.

## 3.5 Client Selection

Considering the heterogeneity of clients in the real world, the proposed mmFedMC is designed to further optimize communication overhead and enhance learning efficiency through a client selection strategy. During the training process of client models, we consider loss $\ell$ as a selection criterion, focusing on the modalities chosen to be communicated to the server. For each modality $m$, the set $\mathcal{L}_m$ includes the loss values $\ell_m^k$ from each client $k$ that has modality $m$ in their selected set $\mathcal{M}_k$:

$$\mathcal{L}_m = \left\{ \ell_m^k \mid k = 1, 2, \ldots, K \text{ and } m \in \mathcal{M}_k \right\}. \quad (17)$$

This collection $\mathcal{L}_m$ allows the server to make informed decisions about which clients to select to participate in the learning process, based on the reported loss values across the different modalities. The client selection based on top of modality selection can be expressed as

$$\begin{aligned} \mathcal{L}_m^\delta &= \text{top} \max_{\lceil \delta \times K \rceil} (\mathcal{L}_m) \\ &= \{x : x \in \mathcal{L}_m \text{ and } |\mathcal{L}_m \cap \{y \mid y \geq x\}| \leq \lceil \delta \times K \rceil \}, \end{aligned} \quad (18)$$

where the ratio parameter $\delta$ establishing the proportion of the higher loss values to be considered. This strategy utilizes the top-$\delta$ criterion to pick a subset of clients, concentrating on those with the higher loss values, which are indicative of their potential effect on the learning efficiency. By utilizing the rounding operation $\lceil \cdot \rceil$ to the product of $\delta$ and the total number of clients $K$, the exact number of clients to be included is determined. Hence, the set of selected clients of modality $m$ can be expressed as

$$\mathcal{K}_m = \left\{ k : \ell_m^k \in \mathcal{L}_\delta^k, \text{ for } k = 1, 2, \ldots, K \right\}. \quad (19)$$

In conclusion, we commence with modality selection to determine the modality set $\mathcal{M}_k$ and the model set $\Theta_\gamma^k$. Subsequently, building upon this foundation, we proceed to client selection, culminating in the identification of the client set $\mathcal{K}_m$ and the final model set to communicate to the server as

$$\Theta_\gamma^\delta = \left\{ \theta_m^k : \theta_m^k \in \Theta_\gamma^k \text{ and } k \in \mathcal{K}_m \right\}, \quad (20)$$

where the set $\Theta_\gamma^k$ represents the modality models corresponding to the top-$\gamma$ priority defined in (16). This client selection approach improves the learning process by optimizing communication overhead and directing resources towards clients that have the most significant impact on the model's performance.

## 3.6 Modality Model Aggregation

Upon receiving data packets from the clients, the server performs a weighted aggregation of the models based on the number of samples for each data modality. For a given data modality $m$, the server aggregates the model parameters from client $k$ with modality $m$. The update is given by

$$\theta_m \leftarrow \sum_{\theta_m^k \in \Theta_p^k} \beta_m^k \theta_m^k, \quad (21)$$

where $\beta_m^k$ represents the aggregation weight coefficient. Following the methodology adopted in FedAvg [2], these weights are determined based on the number of samples, and can be expressed as:

$$\beta_m^k = \frac{|\mathcal{D}_m^k|}{\sum_{k=1}^{K_m} |\mathcal{D}_m^k|}, \quad (22)$$

where $K_m$ denotes the number of client models received by the server for modality $m$.

## 3.7 Deployment

In the mmFedMC framework, only the modality models $\theta_m$ are aggregated, while the ensemble models $\omega_k$ remain separate. The purpose is to give the modality models a global perspective, ensuring a wider generalization. These modality models adhere to the classical FL iterative process and are deployed as global modality models. On the contrary, the personalized ensemble model undergoes a two-stage update. In Stage #1, the ensemble model serves as an intermediate state, primarily facilitating the calculation of the Shapley values. It is not deployed in any application (i.e., it is not tested in any test set). In Stage #2, after receiving the global modality models from the server, the clients subsequently update their ensemble model using global modality models in conjunction with local data to achieve the final state.

## 4 EXPERIMENT AND RESULTS

### 4.1 Dataset

Our experiments are conducted across five diverse multimodal real-world datasets, each incorporating varying numbers of clients, modalities, application contexts, and data types and structures, as depicted in Table 1. Specifically, we consider:

(i) ActionSense [53] is an extensive multimodal dataset that captures human daily activities. It integrates data from six different types of wearable sensors, documenting human interactions with objects and the environment in a kitchen context. Activities include tasks such as peeling a cucumber, slicing a potato, cleaning a plate with a sponge, and so forth. Fig. 3 illustrates the six modality data from the ActionSense dataset for Subject 00 while peeling a potato. The eye tracking data show that the subject's gaze is primarily focused on the potato within their field of vision for most of the time. The left hand is engaged in holding and controlling the rotation of the potato, while the right hand is occupied with gripping the peeler and moving it forward and back to peel the potato. The subject exhibits minimal movement in the torso, and most of the movement occurs in the arms. Consequently, the data reveal that the Tactile Right modality registers high pressure amplitudes due to the peeling action, and some of the Body Tracking signals exhibit regular variations, while other Body Tracking signals remain relatively unchanged. The proposed joint modality and client selection, along with demonstrations of modality impact and client upload frequency, can be found in a demo video available at https://liangqiy.com/mmfedmc/.

(ii) UCI-HAR [54], similar to ActionSense, employs wearable sensors to monitor people during routine activities such as Walking, Walk Downstairs, Sitting, etc.

TABLE 1
Description of Datasets

| Dataset | Client | Task | Modality | Feature |
|---|---|---|---|---|
| ActionSense [1] | 9 Subjects | 20 Kitchen Activities | Eye Tracking | 2 |
| | | | EMG (×2) (Left and Right Arm) | 8(×2) |
| | | | Tactile (×2) (Left and Right Hand) | $32 \times 32$(×2) |
| | | | Body Tracking | $22 \times 3$ |
| UCI-HAR | 30 Subjects | 6 Daily Activities | Accelerometer | $128 \times 3$ |
| | | | Gyroscope | $128 \times 3$ |
| PTB-XL | 39 Hospitals | 5 Diagnosis | Limb Lead ECG [2] | $1000 \times 6$ |
| | | | Precordial Lead ECG [3] | $1000 \times 6$ |
| MELD | 42 Speakers | 4 Emotions | Audio | Client max length×80 [4] |
| | | | Text | 100 |
| DFC2023 | 10 Cities of GF2 Satellite + 17 Cities of SV Satellite | 12 Roof Types | SAR | $32 \times 32 \times 1$ |
| | | | Optical | $32 \times 32 \times 3$ |

[1] Subjects 06 through 09 miss tactile data.
[2] Limb Lead ECG refers to I, II, III, aVL, aVR, aVF leads.
[3] Precordial Lead ECG refers to V1–V6 leads.
[4] Client max length refers to the max length of audio utterances for each client.
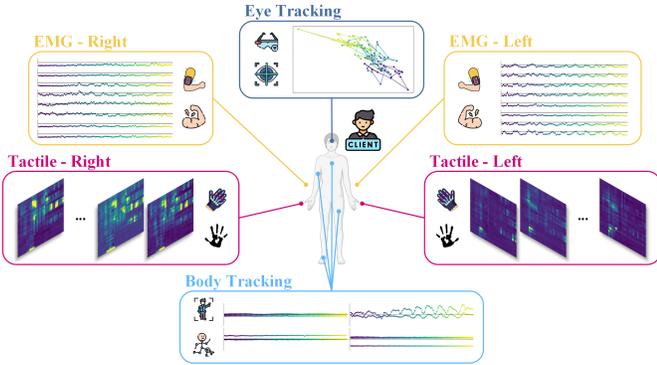


Fig. 3. Data visualization of ActionSense dataset for subject 00 engaged in peeling a potato.

Compared to ActionSense, UCI-HAR features a larger number of subjects but fewer modalities.

(iii) PTB-XL [55], a substantial electrocardiography (ECG) dataset, comprises diverse clinical patient data, featuring 12-lead ECGs from various hospitals. Following established protocols in [56], we treat the Limb Lead ECG and the Precordial Lead ECG as separate modalities. Each data entry includes five labels: Normal, Myocardial Infarction, ST/T Change, Conduction Disturbance, and Hypertrophy.

(iv) MELD [57] is a natural language processing (NLP) dataset derived from dialogues in the TV-series Friends. In each dialogue scene, each speaker is represented by audio and text modalities, reflecting different emotions such as neutral, sadness, joy, and anger.

(v) The 2023 IEEE GRSS Data Fusion Contest (DFC23) [58] is a large-scale dataset featuring rooftop satellite imagery of buildings. It encompasses images from the Gaofen-2 (GF2) and SuperView-1 (SV) satellites, covering rooftops in seventeen cities across six continents. The dataset includes synthetic aperture radar (SAR) and optical images, categorizing twelve types of rooftops.

For each dataset, we conduct experiments in both independent and identically distributed (IID) settings and natural distribution settings. We adhere to a natural client division that reflects real-world scenarios, categorizing entities such as subjects, speakers, hospitals, and satellites as distinct clients.

- In the IID setting, we shuffle all the data from clients and then redistribute these shuffled samples to all clients. In this scenario, individual, group, and system heterogeneities are significantly diminished due to the random and uniform distribution of samples. For samples that are missing certain modalities, the data corresponding to the absent modalities are filled in to facilitate training.
- In the natural distribution setting, each client possesses only its local data and has no knowledge of the other clients' data. This approach not only contends with individual, group, and system heterogeneities but also grapples with heterogeneous data distributions that can greatly impact performance. Challenges in this setting also include wide variations in the number of samples per client, an uneven distribution of classes, and a large number of clients with small sample sizes.

## 4.2 Experiment Setup

**Training Setup.** To ensure a fair comparison across datasets (i – iv), we initially reshape all data modalities into a two-dimensional format, i.e., time × features. For all modalities, we employ a consistent long short-term memory (LSTM) network structure, consisting of a single LSTM layer with 128 hidden units, followed by a fully connected layer. The learning rate $\eta$ for these LSTM models is set to 0.1. For the dataset (v), DFC23, which is composed entirely of images, we use a convolutional neural network (CNN) with one 5x5 convolution layer containing 32 channels, followed by a rectified linear unit (ReLU) activation, a $2 \times 2$ max pooling, and a fully connected layer. The learning rate $\eta$ for these CNN models is set at 0.01. We adopt cross-entropy loss with stochastic gradient descent (SGD) as the optimizer, a batch size of 32, and a local training epoch $E = 5$.

In FL frameworks, the average communication overhead of clients accumulates until it reaches our predefined threshold (5 MB or 25 MB). Taking the ActionSense dataset

as an example, the six modalities – Eye Tracking, EMG Left, EMG Right, Tactile Left, Tactile Right, and Body Tracking – have sizes of 0.27 MB, 0.28 MB, 0.28 MB, 2.26 MB, 2.26 MB, and 0.39 MB, respectively. We simply use the size of the uploaded modality models to simulate communication overhead. However, in real-world scenarios, these overheads would be larger, including client information (e.g., local loss for client selection), metadata of modality models, error correction codes, protocol overheads, encryption, and more.

**Baselines.** At the system level, we use two types of baselines for comprehensive comparisons: the traditional multimodal fusion technique and the ablation study of our method. We consider four state-of-the-art (SOTA) multimodal FL frameworks, encompassing data-level (e.g., [37]), feature-level (e.g., [38]), decision-level (e.g., [39]) fusion, as well as random submodel uploading (e.g., [40] with uniform probabilities). Furthermore, to validate the effectiveness of the proposed mmFedMC, we also conduct experiments where the modality selection strategy, client selection strategy, or both strategies are substituted with random selections. To ensure a fair systematic comparison within the FL context, we do not incorporate various specialized techniques present in these baselines, such as co-attention mechanisms. All these baselines employ a uniform network architecture, specifically, an LSTM/convolution layer followed by a fully connected layer, utilizing a concatenate strategy for fusion.

**Proposed mmFedMC.** For the proposed mmFedMC framework, in addition to the fundamental configurations mentioned above, these modality models do not output logarithmic probabilities but rather provide definitive predicted categories ($\widehat{\mathbb{Y}}$) for the ensemble model ($\omega$). The ensemble model can adopt various choices depending on the specific use case and the resources available on the client, such as voting methods, linear models, $k$-nearest neighbors ($k$-NN), etc. Here, we use the Random Forest (RF) as our ensemble model because of its robust interpretability. We perform a subsampling on the dataset, selecting 50 samples to compute the Shapley values to reduce computational complexity. Note that the ensemble model, Shapley values, the modality model sizes, as well as the recency are kept private by the client and used for modality selection. They are neither uploaded to the server nor does the server possess knowledge of the client's computing methodology.

### 4.3 Results and Overall Comparison

The results of the proposed mmFedMC, in comparison with seven baselines on the ActionSense dataset, are presented in Table 2 and Fig. 4.

**Comparison with SOTA Baselines.** The proposed mmFedMC, along with its three variants employing random selection strategies, demonstrates substantial performance improvements over four other SOTA approaches, while significantly reducing communication costs. The experiment results show that even random selection strategies, despite leading to more aggregation iterations due to reduced communication overhead, contribute to substantial performance gains. Moreover, this improvement in accuracy is strongly attributed to the proposed decision-level fusion approach, where a client receives the global
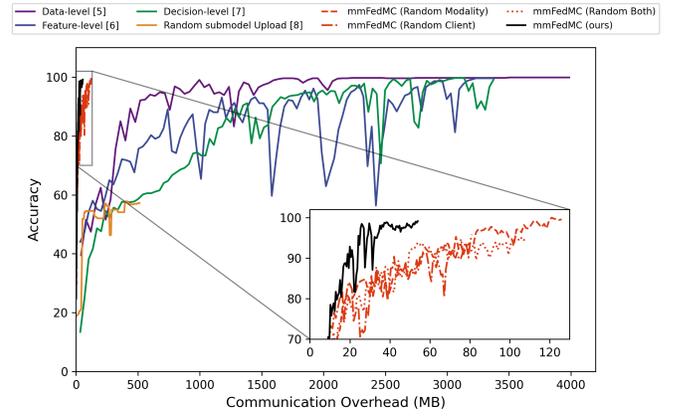


Fig. 4. Comparison of accuracy between proposed mmFedMC and seven baselines on the communication overhead scale.

modality models from the server and undergoes a local ensemble model update. The modality models, through aggregation, not only become more generalized by acquiring information from other clients, but also allow the ensemble model to further fine-tune and personalize based on the new global modality models. The proposed mmFedMC not only exceeds in terms of accuracy but also manages to reduce communication overhead by nearly an order of magnitude, thanks to its selective communication strategy. Insights obtained from mmFedMC's selection approach are:

(i) Different data modalities contribute distinctively to recognition accuracy.

(ii) Aggregation of all modality models is not always necessary.

(iii) Utilizing all modalities for fusion is not imperative.

Although such a [40] random submodel upload strategy effectively reduces communication costs to approximately $\frac{1}{M_k+1}$, it lacks a performance and communication overhead-based selection mechanism, leading to suboptimal results.

**Comparison with Ours Baselines.** Compared to the three random selection methods, the proposed mmFedMC generally achieves a similar or even lower communication overhead, while outperforming in most scenarios. This is attributed to the selective modality and client communication strategy, where not only does the lower communication cost promote more frequent aggregation, but also the Shapley value and client selection based on local loss can identify the clients with a greater impact in each modality. Due to the inherent randomness, random modality and client selections in some scenarios might marginally outperform mmFedMC. This is because randomness, to some extent, ensures that all modalities for all clients have the same probability of being selected. However, this is not the optimal solution as excessive randomness and instability are not preferable.

**Comparison between IID and Natural Distribution.** From the results, it is not always the case that the IID setting provides a clear advantage in model generalization during training. In fact, both IID and natural distribution settings present unique challenges due to client heterogeneity. Under the natural distribution scenario, we typically presume

TABLE 2
Overall Comparison with Baselines at Cumulative Communication Consumption of 5 MB per Client

### IID Setting

| Method | ActionSense Acc. [1] (%) ↑ | ActionSense Comm. [2] (MB) ↓ | UCI-HAR Acc. (%) ↑ | UCI-HAR Comm. (MB) ↓ | PTB-XL Acc. (%) ↑ | PTB-XL Comm. (MB) ↓ | MELD Acc. (%) ↑ | MELD Comm. (MB) ↓ | DFC23 Acc. (%) ↑ | DFC23 Comm. (MB) ↓ |
|---|---|---|---|---|---|---|---|---|---|---|
| Data-level (e.g., [37]) | 13.32 | 4.43 | 33.27 | 0.27 | 19.80 | 0.28 | 55.01 | 3.62 | 59.64 | 1.51 |
| Feature-level (e.g., [38]) | 14.00 | 3.74 | 23.50 | 0.53 | 4.76 | 0.54 | 55.05 | 3.88 | 60.03 | 3.01 |
| Decision-level (e.g., [39]) | 9.03 | 3.75 | 20.13 | 0.53 | 4.22 | 0.54 | 54.26 | 3.88 | 59.52 | 3.01 |
| Random Submodal Upload (e.g., [40]) | 8.29 | 0.57 | 24.28 | 0.18 | 4.29 | 0.18 | 55.39 | 1.29 | 59.80 | 1.00 |
| mmFedMC (Random Modality) | 68.73 | 0.20 | 76.92 | **0.05** | 86.50 | **0.05** | 54.89 | 0.40 | 61.16 | **0.28** |
| mmFedMC (Random Client) | 81.79 | 0.10 | 70.90 | **0.05** | 85.66 | **0.05** | 55.58 | **0.13** | 63.81 | **0.28** |
| mmFedMC (Random Both) [3] | 68.00 | 0.21 | 71.63 | **0.05** | 86.90 | **0.05** | 55.44 | 0.52 | **64.36** | **0.28** |
| **mmFedMC (Ours)** [4] | **92.28** [5] | 0.10 | **78.10** | 0.05 | **87.04** | 0.05 | **55.82** | 0.16 | 62.79 | **0.28** |

### Natural Distribution

| Method | ActionSense Acc. (%) ↑ | ActionSense Comm. (MB) ↓ | UCI-HAR Acc. (%) ↑ | UCI-HAR Comm. (MB) ↓ | PTB-XL Acc. (%) ↑ | PTB-XL Comm. (MB) ↓ | MELD Acc. (%) ↑ | MELD Comm. (MB) ↓ | DFC23 Acc. (%) ↑ | DFC23 Comm. (MB) ↓ |
|---|---|---|---|---|---|---|---|---|---|---|
| Data-level (e.g., [37]) | 51.57 | 4.43 | 35.89 | 0.27 | 22.28 | 0.28 | 51.73 | 3.62 | 62.49 | 0.27 |
| Feature-level (e.g., [38]) | 47.24 | 3.74 | 21.34 | 0.53 | 13.97 | 0.54 | 51.85 | 3.88 | 61.81 | 0.53 |
| Decision-level (e.g., [39]) | 24.81 | 3.75 | 17.70 | 0.53 | 16.40 | 0.54 | 49.39 | 3.88 | 62.40 | 0.53 |
| Random Submodal Upload (e.g., [40]) | 34.77 | 0.57 | 22.88 | 0.18 | 17.44 | 0.18 | 49.10 | 1.29 | 62.13 | 0.18 |
| mmFedMC (Random Modality) | 88.85 | 0.14 | 71.12 | **0.05** | 49.80 | **0.05** | 48.35 | 0.41 | 66.47 | **0.05** |
| mmFedMC (Random Client) | 93.29 | 0.09 | 60.54 | **0.05** | 51.44 | **0.05** | 53.67 | **0.16** | 67.41 | **0.05** |
| mmFedMC (Random Both) | 89.30 | 0.12 | 60.22 | **0.05** | 54.18 | **0.05** | **57.43** | 0.49 | 66.69 | **0.05** |
| **mmFedMC (Ours)** | **98.19** | 0.06 | **71.28** | 0.05 | **55.09** | 0.05 | 53.31 | 0.23 | **67.61** | 0.05 |

[1] Average accuracy between clients, ↑ refers to the higher (preferred).
[2] Communication overhead per iteration, ↓ refers to the lower (preferred).
[3] Random Both refers to Random Modality and Client Selection.
[4] **mmFedMC (Ours)** utilize a consistent configuration with $\delta = 0.2, \gamma = 1, \alpha_s = 1/3, \alpha_c = 1/3, \alpha_r = 1/3$ across all datasets.
[5] **Bold** refers to the highest accuracy or the lowest communication overhead.

that each client's training and testing data distributions are similar yet biased. For example, a hospital in Asia is unlikely to receive a large influx of patients from North America when deploying a model. In such scenarios, the ensemble model adapts to Asian patient data by fine-tuning and personalizing based on the global modality models. Therefore, in the IID setting, the ensemble model may not effectively leverage its potential. Moreover, in the natural distribution scenario, the number of samples per client might more likely follow a normal distribution, whereas in the IID setting, it would exhibit a uniform distribution. Take the PTB-XL dataset as an example, where three hospital sites hold 93.54% of the data in the FL framework, with 34 out of 39 hospital sites having fewer than 100 samples each. Similarly, in the MELD dataset, six speakers possess 92.68% of the data in the FL framework, and 36 out of 42 speakers have fewer than 100 samples each. Consequently, many clients with few samples could lead to overfitting, resulting in small local losses and impacting our client selection strategy. We will further discuss the implications of client sample size in the natural distribution scenario in Section 5.2.

## 4.4 Ablation Study for Modality Selection

### 4.4.1 Trade-Off Analysis

To better investigate the impact of configuring the four parameters, $\gamma$, $\alpha_s$, $\alpha_c$, and $\alpha_r$, in modality selection on

performance and communication overhead, Tables 3 and 4 present the results on the ActionSense and UCI-HAR datasets when client selection is not performed (i.e., 100%, client $\delta = 1$), respectively. For this ablation study, we set the cumulative communication overhead threshold at 25 MB per client to observe changes over more extended iterations.

Considering the communication constraints, our results underscore the need to find a compromise between modality impact $\alpha_s$ and communication overhead $\alpha_c$ to optimize performance when $\gamma$ is constant. Increasing $\gamma$ does not always lead to better results, as it can exacerbate communication overhead. In some instances, communicating only a select few informative modalities yields enhanced performance. In this context, the proposed mmFedMC framework facilitates a flexible determination of the number of modality models to upload, balancing model performance, communication overhead, and learning efficiency.

Building on this, the introduction of the recency term $\alpha_r$ also provides generalizability throughout the framework to avoid falling into a *single modality optimization trap*. As the sizes of the modality models are constant, if a modality with a smaller feature dimension contains comparable information, the FL framework might initially choose it due to its larger Shapley value. However, due to server aggregation, this modality model will have an even greater Shapley value in the next communication round, leading the FL framework to continuously focus only on this modality,

TABLE 3
ActionSense - Comparison of Accuracy and Communication Overhead
at Cumulative Communication Consumption of 25 MB per client

| Method | $\gamma$ | $\alpha_s$ | $\alpha_c$ | $\alpha_r$ | Acc. (%) ↑ | Comm. (MB) ↓ | Comm. Round |
|---|---|---|---|---|---|---|---|
| Data-level (e.g., [37]) | | | | | 62.43 | 4.43 | 5 |
| Feature-level (e.g., [38]) | | | | | 59.28 | 3.74 | 6 |
| Decision-level (e.g., [39]) | | | | | 54.69 | 3.75 | 6 |
| Random Submodal Upload (e.g., [40]) | | | | | 56.93 | 0.57 | 44 |
| mmFedMC ($\delta = 1$) (Random Modality) | 1 | - | - | - | 95.97 | 0.67 | 37 |
| | | 1.0 | 0.0 | 0.0 | 85.14 | 0.92 | 27 |
| | | 0.0 | 1.0 | 0.0 | 90.91 | 0.27 | 93 |
| | | 0.0 | 0.0 | 1.0 | 86.34 | 0.63 | 38 |
| | 1 | 0.0 | 0.5 | 0.5 | 88.95 | 0.29 | 84 |
| | | 0.5 | 0.0 | 0.5 | 97.92 | 0.54 | 46 |
| | | 0.5 | 0.5 | 0.0 | 98.05 | 0.34 | 73 |
| | | **1/3** | **1/3** | **1/3** | **99.58** | **0.31** | **80** |
| | | 1.0 | 0.0 | 0.0 | 88.06 | 1.56 | 16 |
| | | 0.0 | 1.0 | 0.0 | 97.36 | 0.55 | 45 |
| | | 0.0 | 0.0 | 1.0 | 87.76 | 1.27 | 19 |
| mmFedMC ($\delta = 1$) | 2 | 0.0 | 0.5 | 0.5 | 86.45 | 0.59 | 42 |
| | | 0.5 | 0.0 | 0.5 | 91.94 | 1.13 | 22 |
| | | 0.5 | 0.5 | 0.0 | 97.92 | 0.59 | 42 |
| | | 1/3 | 1/3 | 1/3 | 97.08 | 0.67 | 37 |
| | | 1.0 | 0.0 | 0.0 | 93.45 | 1.91 | 13 |
| | | 0.0 | 1.0 | 0.0 | 88.97 | 0.83 | 30 |
| | | 0.0 | 0.0 | 1.0 | 87.91 | 1.87 | 12 |
| | 3 | 0.0 | 0.5 | 0.5 | 87.28 | 0.87 | 28 |
| | | 0.5 | 0.0 | 0.5 | 91.66 | 2.30 | 10 |
| | | 0.5 | 0.5 | 0.0 | 94.44 | 1.07 | 23 |
| | | 1/3 | 1/3 | 1/3 | 94.71 | 1.32 | 18 |

TABLE 4
UCI-HAR - Comparison of Accuracy and Communication Overhead at
Cumulative Communication Consumption of 25 MB per client

| Method | $\gamma$ | $\alpha_s$ | $\alpha_c$ | $\alpha_r$ | Acc. (%) ↑ | Comm. (MB) ↓ | Comm. Round |
|---|---|---|---|---|---|---|---|
| Data-level (e.g., [37]) | | | | | 71.90 | 0.27 | 93 |
| Feature-level (e.g., [38]) | | | | | 43.85 | 0.53 | 47 |
| Decision-level (e.g., [39]) | | | | | 37.36 | 0.53 | 47 |
| Random Submodel Upload (e.g., [40]) | | | | | 46.85 | 0.18 | 142 |
| mmFedMC ($\delta = 1$) (Random Modality) | 1 | - | - | - | 75.53 | 0.26 | 95 |
| | | 1.0 | 0.0 | 0.0 | 74.51 | 0.26 | 95 |
| | | 0.0 | 1.0 | 0.0 | 37.61 | 0.26 | 95 |
| | | 0.0 | 0.0 | 1.0 | 56.97 | 0.26 | 95 |
| | 1 | 0.0 | 0.5 | 0.5 | 57.49 | 0.26 | 95 |
| | | 0.5 | 0.0 | 0.5 | 74.22 | 0.26 | 95 |
| | | **0.5** | **0.5** | **0.0** | **75.58** | **0.26** | **95** |
| mmFedMC ($\delta = 1$) | | 1/3 | 1/3 | 1/3 | 75.57 | 0.26 | 95 |
| | | 1.0 | 0.0 | 0.0 | 53.81 | 0.53 | 47 |
| | | 0.0 | 1.0 | 0.0 | 54.53 | 0.53 | 47 |
| | | 0.0 | 0.0 | 1.0 | 54.36 | 0.53 | 47 |
| | 2 | 0.0 | 0.5 | 0.5 | 54.34 | 0.53 | 47 |
| | | 0.5 | 0.0 | 0.5 | 56.93 | 0.53 | 47 |
| | | 0.5 | 0.5 | 0.0 | 50.55 | 0.53 | 47 |
| | | 1/3 | 1/3 | 1/3 | 43.29 | 0.53 | 47 |



(a) ActionSense Dataset with Six Modalities



(b) UCI-HAR with Two Modalities

Fig. 5. The mean Shapley value (i.e., impact) of modality models throughout the mmFedMC iteration.
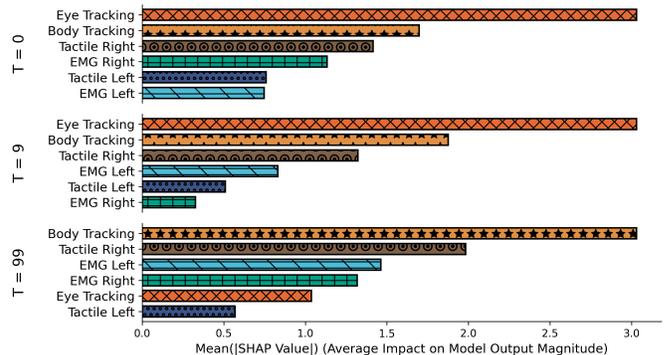
a situation we refer to as the single modality optimization trap. In such cases, a modality might be sufficiently trained within just 10 communication rounds, but the FL framework should shift its focus to other modalities to fully leverage the advantages of multimodality. Thus, the introduction of the recency term is essentially a means of injecting uncertainty, allowing FL to focus on other modalities. It is important to note that our goal is not to train all modality models perfectly, especially in edge clients with limited resources. Instead, we aim for FL to focus on certain modalities that not only possess rich information but also have a sufficiently small feature dimension.

Shapley value, communication overhead, and recency each exhibit superior performance in specific scenarios, while potentially becoming less effective in others. For example, in the UCI-HAR dataset, the two modalities have identical data dimensions (i.e., 128×3), leading to the same modality model size (i.e., 0.26 MB). Consequently, using communication overhead for modality selection is ineffective in this context. Moreover, although we intend to use recency to encourage the FL framework to focus on different modalities, in cases with only two data modalities, recency can cause the framework to cyclically select these two modalities. This significantly diminishes the effectiveness of the Shapley value. Therefore, our proposed mmFedMC is expected to perform better in datasets like ActionSense, which feature multimodal data with varying feature dimensions. In such settings, the diverse dimensions enable more nuanced and effective modality selection, leveraging the strengths of Shapley value, communication overhead, and recency in a more balanced and impactful manner.
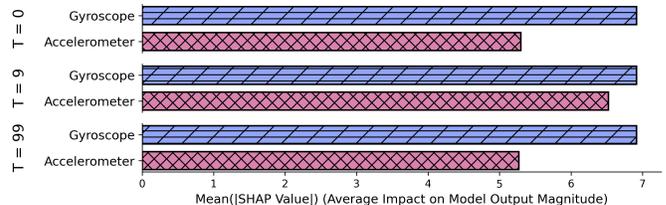
### 4.4.2 Analytics on Modality Impact

In addition to helping clients select the modality models, the Shapley value also offers an interpretative approach to quantify the modality models. During the FL process, we can see the clients' favor and the efficacy of each modality model within the mmFedMC framework. Fig. 5 illustrates this dynamics, showing the impact of each data modality on the final prediction of the ensemble model across different communication rounds $T$. In the initial stages of FL, data with simpler features exhibit higher impact because the modality model will more easily capture its information. As communication rounds progress, modalities with larger

TABLE 5
ActionSense - Client Selection at Cumulative Consumption of 5 MB per client

| $\delta$ | $\gamma$ | $\alpha_s$ | $\alpha_c$ | $\alpha_r$ | Higher Loss | | | Lower Loss | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Acc. (%) ↑ | Comm. (MB) ↓ | Comm. Round | Acc. (%) ↑ | Comm. (MB) ↓ | Comm. Round |
| 0.2 | 1 | 1.0 | 0.0 | 0.0 | 81.97(+29.16%) | 0.05(-80.77%) | 99(+80) | 98.30(+54.88%) | 0.09(-65.48%) | 57(+38) |
| | | 0.0 | 1.0 | 0.0 | 74.04(-10.00%) | 0.06(0.00%) | 84(0) | 98.87(+20.20%) | 0.06(0.00%) | 84(0) |
| | | 0.0 | 0.0 | 1.0 | 88.17(+3.95%) | 0.13(0.00%) | 38(0) | 84.68(-0.17%) | 0.13(0.00%) | 38(0) |
| | | 0.0 | 0.5 | 0.5 | 85.30(-2.90%) | 0.07(0.00%) | 76(0) | 88.59(+0.84%) | 0.07(0.00%) | 76(0) |
| | | 0.5 | 0.0 | 0.5 | 87.59(-6.27%) | 0.05(+1.37%) | 92(0) | 98.87(+5.80%) | 0.06(+12.25%) | 83(-9) |
| | | 0.5 | 0.5 | 0.0 | 90.62(+5.70%) | 0.08(+9.40%) | 61(-5) | 92.88(+34.04%) | 0.14(-64.29%) | 35(+23) |
| | | 1/3 | 1/3 | 1/3 | 85.06(-8.82%) | 0.07(-9.48%) | 76(+7) | 98.19(+14.52%) | 0.06(-15.66%) | 79(+13) |
| | 2 | 1.0 | 0.0 | 0.0 | 82.24(+18.69%) | 0.09(-75.98%) | 53(+41) | 98.89(+6.00%) | 0.06(-9.98%) | 77(+8) |
| | | 0.0 | 1.0 | 0.0 | 82.50(-7.70%) | 0.12(0.00%) | 41(0) | 94.55(+5.77%) | 0.12(0.00%) | 41(0) |
| | | 0.0 | 0.0 | 1.0 | 86.39(-0.95%) | 0.26(0.00%) | 19(0) | 86.11(-1.27%) | 0.26(0.00%) | 19(0) |
| | | 0.0 | 0.5 | 0.5 | 82.55(-5.27%) | 0.13(0.00%) | 37(0) | 89.66(+2.89%) | 0.13(0.00%) | 37(0) |
| | | 0.5 | 0.0 | 0.5 | 87.87(-2.03%) | 0.13(-19.20%) | 38(+7) | 91.48(+2.01%) | 0.16(+1.38%) | 30(-1) |
| | | 0.5 | 0.5 | 0.0 | 86.17(+0.16%) | 0.14(+0.68%) | 34(-1) | 98.60(+14.61%) | 0.14(+0.59%) | 35(0) |
| | | 1/3 | 1/3 | 1/3 | 86.47(-2.82%) | 0.13(-1.71%) | 38(+1) | 97.05(+9.07%) | 0.14(+5.67%) | 35(-2) |
| | 3 | 1.0 | 0.0 | 0.0 | 86.75(+18.81%) | 0.17(-68.80%) | 30(+21) | 89.25(+22.24%) | 0.21(-61.33%) | 24(+15) |
| | | 0.0 | 1.0 | 0.0 | 79.59(-7.19%) | 0.18(0.00%) | 27(0) | 93.30(+8.80%) | 0.18(0.00%) | 27(0) |
| | | 0.0 | 0.0 | 1.0 | 86.77(+5.03%) | 0.39(0.00%) | 12(0) | 86.67(+4.90%) | 0.39(0.00%) | 12(0) |
| | | 0.0 | 0.5 | 0.5 | 83.65(-7.01%) | 0.19(0.00%) | 25(0) | 91.21(+1.41%) | 0.19(0.00%) | 25(0) |
| | | 0.5 | 0.0 | 0.5 | 89.55(-2.39%) | 0.21(-2.11%) | 24(+1) | 90.20(-1.68%) | 0.26(+20.91%) | 19(-4) |
| | | 0.5 | 0.5 | 0.0 | 81.83(-4.73%) | 0.18(-9.07%) | 27(+3) | 93.00(+8.28%) | 0.21(+2.48%) | 24(0) |
| | | 1/3 | 1/3 | 1/3 | 81.97(-7.74%) | 0.19(-2.31%) | 26(0) | 93.44(+5.16%) | 0.21(+8.09%) | 24(-2) |

Improvement over Baseline (Random Client Selection)
| |
|---|
| ≥ 10% |
| 0% to 10% |
| −10% to 0% |
| ≤ -10% |

TABLE 6
UCI-HAR - Client Selection at Cumulative Consumption of 5 MB per cliente 5 MB

| $\delta$ | $\gamma$ | $\alpha_s$ | $\alpha_c$ | $\alpha_r$ | Higher Loss | | | lower Loss | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | Acc. (%) ↑ | Comm. (MB) ↓ | Comm. Round | Acc. (%) ↑ | Comm. (MB) ↓ | Comm. Round |
| 0.2 | 1 | 1.0 | 0.0 | 0.0 | 47.78(-22.79%) | 0.05(0.00%) | 95(0) | 68.77(+11.13%) | 0.05(0.00%) | 95(0) |
| | | 0.0 | 1.0 | 0.0 | 33.80(-39.88%) | 0.05(0.00%) | 95(0) | 55.77(-0.80%) | 0.05(0.00%) | 95(0) |
| | | 0.0 | 0.0 | 1.0 | 40.31(-18.33%) | 0.05(0.00%) | 95(0) | 68.27(+38.30%) | 0.05(0.00%) | 95(0) |
| | | 0.0 | 0.5 | 0.5 | 37.29(-23.51%) | 0.05(0.00%) | 95(0) | 67.60(+38.67%) | 0.05(0.00%) | 95(0) |
| | | 0.5 | 0.0 | 0.5 | 42.39(-34.28%) | 0.05(0.00%) | 95(0) | 70.62(+9.50%) | 0.05(0.00%) | 95(0) |
| | | 0.5 | 0.5 | 0.0 | 47.70(-14.81%) | 0.05(0.00%) | 95(0) | 65.16(+16.36%) | 0.05(0.00%) | 95(0) |
| | | 1/3 | 1/3 | 1/3 | 39.74(-24.53%) | 0.05(0.00%) | 95(0) | 71.28(+35.37%) | 0.05(0.00%) | 95(0) |
| | 2 | 1.0 | 0.0 | 0.0 | 36.74(-22.25%) | 0.11(0.00%) | 47(0) | 65.77(+39.19%) | 0.11(0.00%) | 47(0) |
| | | 0.0 | 1.0 | 0.0 | 35.43(-29.30%) | 0.11(0.00%) | 47(0) | 67.37(+34.43%) | 0.11(0.00%) | 47(0) |
| | | 0.0 | 0.0 | 1.0 | 41.61(-15.11%) | 0.11(0.00%) | 47(0) | 68.22(+39.17%) | 0.11(0.00%) | 47(0) |
| | | 0.0 | 0.5 | 0.5 | 37.24(-27.29%) | 0.11(0.00%) | 47(0) | 65.35(+27.58%) | 0.11(0.00%) | 47(0) |
| | | 0.5 | 0.0 | 0.5 | 39.52(-23.88%) | 0.11(0.00%) | 47(0) | 66.99(+29.02%) | 0.11(0.00%) | 47(0) |
| | | 0.5 | 0.5 | 0.0 | 37.60(-20.59%) | 0.11(0.00%) | 47(0) | 69.74(+47.28%) | 0.11(0.00%) | 47(0) |
| | | 1/3 | 1/3 | 1/3 | 36.02(-26.67%) | 0.11(0.00%) | 47(0) | 67.83(+38.10%) | 0.11(0.00%) | 47(0) |

Improvement over Baseline (Random Client Selection)
| |
|---|
| ≥ 10% |
| 0% to 10% |
| −10% to 0% |
| ≤ -10% |

feature sets and complex models, due to their higher communication overhead, take on a subordinate role in their selection within mmFedMC. As FL advances, more straightforward modalities that still convey ample information, such as Body Tracking, emerge as primary contributors.

For the ActionSense dataset, Fig. 5a illustrates a game of dramatic dynamic changes in modality selection throughout the FL process. In the early stages of FL, the Eye Tracking modality, due to its simple structure, facilitated easy extraction of information from features and had a smaller data size, leading to a lower communication overhead. This resulted in a higher upload frequency, placing it advantageously in the balance between information gain and communication cost. However, as communication rounds progressed, the selection frequency of the Eye Tracking modality tended to decrease, while that of the Body Tracking modality increased. This shift in modality preference is driven by the trade-off between information value and communication cost. Furthermore, as the FL process progresses and modality models become more well trained, the impact of simpler but less information-rich modalities like Eye Tracking gradually diminishes. Conversely, the Tactile Right modality, despite being the most complex and having the highest communication cost, maintained a preferential secondary position due to its rich information content. On the other hand, the Tactile Left modality, while having a broader feature dimension, remained at a disadvantage possibly due to less diverse actions from the subject's left hand. Over time, this increased communication cost impacts the overall efficiency of the FL process, leading to a preference for modalities like Body Tracking, which offer valuable information with less data transmission and hence are more cost-effective.

In contrast, for the UCI-HAR dataset, Fig. 5b presents only subtle changes in modality selection during the FL process. Given that this dataset comprises only two modalities with identical feature dimensions, the competition between them hinges solely on the Shapley value and recency. As discussed in Section 4.4.1, recency tends to diminish the effectiveness of Shapley values. Thus, in sit-

uations where feature dimensions, modality model sizes, and information richness are similar, the Accelerometer and Gyroscope display comparable the Shapley value. The slightly higher Shapley value of the Gyroscope compared to the Accelerometer can be attributed as Gyroscopes are particularly adept at capturing rotational movements and angular velocities, which are crucial in differentiating these activities. For example, the distinct patterns of rotation and angular momentum while running or climbing stairs are more pronounced and detectable by a gyroscope than by an accelerometer.

### 4.5 Ablation Study for Client Selection

#### 4.5.1 Higher vs. Lower Local Loss Selection

Tables 5 and 6 showcase the performance of client selection strategies based on higher and lower local losses on the ActionSense and UCI-HAR datasets, respectively, and compare these two contrasting strategies against a baseline (i.e., random client selection) across various configurations. In our proposed multimodal decision-level fusion FL framework, where each modality has a global modality model, our preference leans towards quickly converging these global modality models to a local minima before proceeding to fusion. It is evident that in most configurations, the strategy of selecting clients with lower local losses outperforms the random client selection approach, which in turn is superior to the strategy of selecting clients with higher local losses. Although [36] demonstrated, under certain strong assumptions, such as $L$–smooth and $\mu$–strongly convex, that selecting clients with the higher local loss could accelerate convergence, this approach may not be as effective in real-world scenarios. Particularly in cases where group heterogeneity results in a bimodal or even multimodal loss function (i.e., two or more local minima), choosing higher loss clients could cause the global model to oscillate between these local minima.

#### 4.5.2 Comparison of Client Selection Frequency

Fig. 6 shows the frequency with which clients are selected under strategies based on higher and lower local losses, with the aim of demonstrating the preference for certain clients by these two different strategies. It is observable that the strategy that favors higher local loss results in a more uniform distribution of client selection (i.e., each client is chosen with a more frequency). In contrast, the strategy that prioritizes lower local loss demonstrates a preference for a subset of clients. The rationale behind these distinct distributions of client selection frequency lies in the fact that local loss depends not only on the richness of information in client data, but also heavily on client heterogeneity and model optimization, among other factors. A client significantly deviating from most others due to heterogeneity and also high data noise, erroneous data, or even malicious attacks like data poisoning, would exhibit higher local loss. Therefore, while the strategy based on lower local loss may risk overfitting due to data homogeneity, the higher local loss strategy could be greatly perturbed by client heterogeneity. Furthermore, thanks to multimodality and decision-level fusion, the risk of overfitting can be further mitigated.



(a) ActionSense - Higher Loss    (b) ActionSense - Lower Loss

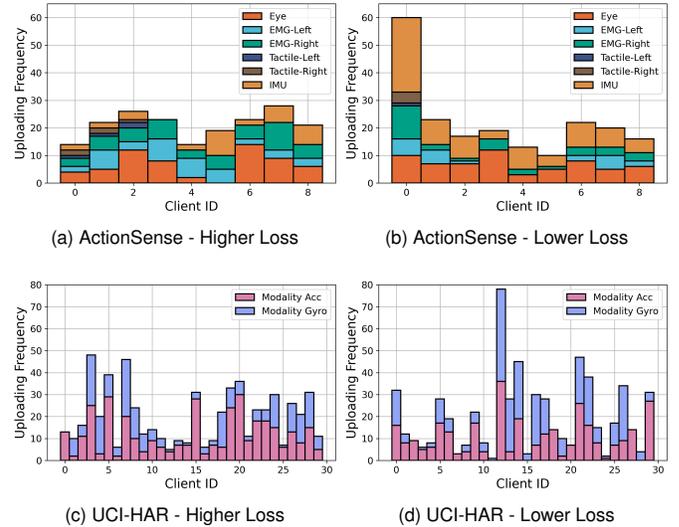(c) UCI-HAR - Higher Loss    (d) UCI-HAR - Lower Loss

Fig. 6. Histograms of the client selection frequency based on the criteria of the higher and lower loss.

## 5 Discussion

### 5.1 Efficiency of Modality Selection

This paper provides a modality selection trade-off based on Shapley value, communication overhead, and recency, with flexible parameters that play different roles in various scenarios. On the one hand, changing the needs of applications can lead to adjustments in parameter weights; for example, communication resources may be more abundant at night, allowing for a lower setting of the weight of communication overhead $\alpha_c$. On the other hand, in certain contexts, these parameters may not be effective; for example, if two modalities have identical data dimensions, as is the case with the UCI-HAR dataset, the trade-off related to communication overhead becomes irrelevant. Therefore, in such situations, we might consider alternative parameters to optimize modality selection, such as sensor sampling rates, energy consumption, sensor cost, and so on [59]. These considerations could provide additional dimensions to the selection strategy, enabling a more comprehensive approach tailored to the specific constraints and objectives of the FL system in question.

### 5.2 Influence of Overfitting on Client Selection

The proposed client selection strategy is predicated on local loss; however, in instances of extremely unbalanced data distribution, the ostensibly low local losses, which may arise from overfitting, could potentially compromise the efficacy of client selection. This susceptibility to overfitting arises for three factors: low statistical class divergence (e.g., clients possessing only one or two classes), a small number of client samples (e.g., clients with only tens of samples), and varied modality properties (e.g., different sensors have different sampling rates). Such scenarios are prevalent, particularly in healthcare, where many hospitals may have a small patient population, and the majority of these patients may present as normal. In the datasets experimented in this paper, for example, in the PTB-XL dataset, 34 out of 39

clients have fewer than 100 samples; in the MELD dataset, 36 out of 42 clients have fewer than 100 samples; and in the DFC23 dataset, 6 out of 27 clients exhibit significant class imbalance (i.e., one class represents more than 80% of the total sample size). In the ActionSense dataset, the six different sensors have varying sampling rates, with Tactile sensors at 15 Hz and Body Tracking sensors at 60 Hz, leading to fewer samples for Tactile modalities, and thus a heightened risk of overfitting. Additionally, certain sensors may not be well-suited for some applications, which can lead to the model learning from noise and overfitting, such as the Tactile Left sensor displaying minimal signal variation across different activities.

Disparities in data distribution can amplify overfitting risks, affecting the local loss metric, and potentially resulting in biased client selection. Clients with such extreme imbalances might be considered outliers and could be pre-emptively excluded from the FL process, akin to a static client pre-selection strategy. Post-FL, these outlier clients could benefit from using the trained global models as a starting point for further fine-tuning, thus addressing the challenges posed by their unique data characteristics.

## 6 CONCLUSION

In this paper, we introduce the mmFedMC framework that minimizes communication overhead and enhances learning efficiency through joint modality and client selection strategies. By optimizing modality selection with Shapley values, modality model sizes, and recency, as well as favoring clients with lower local loss for client selection, we achieve considerable recognition accuracy and up to 20x increase in communication efficiency. The proposed mmFedMC is highly flexible, suitable for heterogeneous clients, offers modular modality models that can be detached, and provides impact for data modalities. Extensive experiments across a diverse range of real-world datasets, including wearable sensors, healthcare, NLP, and remote sensing satellite datasets, demonstrate the outstanding performance of the mmFedMC framework across various data modalities.

Our future work will focus on enhancing the adaptability of mmFedMC through dynamic configuration. For modality selection, we could adjust the weight of communication overhead based on the dynamically available communication resources in the real world (such as higher bandwidth at night), allowing for the uploading of more modality models. For client selection, we may consider a dynamic strategy based on local loss, employing selection of higher local losses at the initial stages of FL to speed up convergence, and transitioning to the proposed lower local loss selection towards the end to optimize local minima. Additionally, we contemplate that Shapley values could further refine sensor deployment, potentially halting the training of underperforming modalities to reduce computational overhead, energy consumption, and sensor costs.

## REFERENCES

[1] L. Yuan, D.-J. Han, V. P. Chellapandi, S. H. Żak, and C. G. Brinton, "Fedmfs: Federated multimodal fusion learning with selective modality communication," *arXiv preprint arXiv:2310.07048*, 2023.

[2] B. McMahan, E. Moore, D. Ramage, S. Hampson, and B. A. y Arcas, "Communication-efficient learning of deep networks from decentralized data," in *Artificial intelligence and statistics*. PMLR, 2017, pp. 1273–1282.

[3] L. Yuan, L. Sun, P. S. Yu, and Z. Wang, "Decentralized federated learning: A survey and perspective," *arXiv preprint arXiv:2306.01603*, 2023.

[4] W. Fang, Z. Yu, Y. Jiang, Y. Shi, C. N. Jones, and Y. Zhou, "Communication-efficient stochastic zeroth-order optimization for federated learning," *IEEE Transactions on Signal Processing*, vol. 70, pp. 5058–5073, 2022.

[5] V. P. Chellapandi, A. Upadhyay, A. Hashemi, and S. H. Żak, "On the Convergence of Decentralized Federated Learning Under Imperfect Information Sharing," *IEEE Control Systems Letters*, 2023.

[6] L. Yuan, Z. Wang, and C. G. Brinton, "Digital ethics in federated learning," *arXiv preprint arXiv:2310.03178*, 2023.

[7] H. Yu, Z. Chen, X. Zhang, X. Chen, F. Zhuang, H. Xiong, and X. Cheng, "Fedhar: Semi-supervised online learning for personalized federated human activity recognition," *IEEE Transactions on Mobile Computing*, 2021.

[8] Q. Wu, X. Chen, Z. Zhou, and J. Zhang, "Fedhome: Cloud-edge based personalized federated learning for in-home health monitoring," *IEEE Transactions on Mobile Computing*, vol. 21, no. 8, pp. 2818–2832, 2020.

[9] S. Baghersalimi, T. Teijeiro, A. Aminifar, and D. Atienza, "Decentralized federated learning for epileptic seizures detection in low-power wearable systems," *IEEE Transactions on Mobile Computing*, 2023.

[10] L. Yuan, Y. Ma, L. Su, and Z. Wang, "Peer-to-peer federated continual learning for naturalistic driving action recognition," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2023, pp. 5249–5258.

[11] V. P. Chellapandi, L. Yuan, S. H. Zak, and Z. Wang, "A Survey of Federated Learning for Connected and Automated Vehicles," *2023 IEEE 26th International Conference on Intelligent Transportation Systems (ITSC)*, 2023.

[12] V. P. Chellapandi, L. Yuan, C. G. Brinton, S. H. Zak, and Z. Wang, "Federated Learning for Connected and Automated Vehicles: A Survey of Existing Approaches and Challenges," *IEEE Transactions on Intelligent Vehicles*, 2023.

[13] D. Upadhyay, H. Shui, and D. Filev, "Methods and systems for estimating a remaining useful life of an asset," Mar. 2 2023, uS Patent App. 17/446,708.

[14] D. Lahat, T. Adali, and C. Jutten, "Multimodal data fusion: an overview of methods, challenges, and prospects," *Proceedings of the IEEE*, vol. 103, no. 9, pp. 1449–1477, 2015.

[15] T. Baltrušaitis, C. Ahuja, and L.-P. Morency, "Multimodal machine learning: A survey and taxonomy," *IEEE transactions on pattern analysis and machine intelligence*, vol. 41, no. 2, pp. 423–443, 2018.

[16] G. Lee, S.-H. Jung, and D. Han, "An adaptive sensor fusion framework for pedestrian indoor navigation in dynamic environments," *IEEE Transactions on Mobile Computing*, vol. 20, no. 2, pp. 320–336, 2019.

[17] Y. Li, L. Liu, H. Qin, S. Deng, M. A. El-Yacoubi, and G. Zhou, "Adaptive deep feature fusion for continuous authentication with data augmentation," *IEEE Transactions on Mobile Computing*, 2022.

[18] J. Zhang, Q. Wang, Q. Wang, and Z. Zheng, "Multimodal fusion framework based on statistical attention and contrastive attention for sign language recognition," *IEEE Transactions on Mobile Computing*, 2023.

[19] P. Bharti, D. De, S. Chellappan, and S. K. Das, "Human: Complex activity recognition with multi-modal multi-positional body sensing," *IEEE Transactions on Mobile Computing*, vol. 18, no. 4, pp. 857–870, 2018.

[20] P. K. Atrey, M. A. Hossain, A. El Saddik, and M. S. Kankanhalli, "Multimodal fusion for multimedia analysis: a survey," *Multimedia systems*, vol. 16, pp. 345–379, 2010.

[21] S. K. D'mello and J. Kory, "A review and meta-analysis of multimodal affect detection systems," *ACM computing surveys (CSUR)*, vol. 47, no. 3, pp. 1–36, 2015.

[22] D. Ramachandram and G. W. Taylor, "Deep multimodal learning: A survey on recent advances and trends," *IEEE signal processing magazine*, vol. 34, no. 6, pp. 96–108, 2017.

[23] A. Kline, H. Wang, Y. Li, S. Dennis, M. Hutch, Z. Xu, F. Wang, F. Cheng, and Y. Luo, "Multimodal machine learning in precision health: A scoping review," *npj Digital Medicine*, vol. 5, no. 1, p. 171, 2022.

[24] P. Kairouz, H. B. McMahan, B. Avent, A. Bellet, M. Bennis, A. N. Bhagoji, K. Bonawitz, Z. Charles, G. Cormode, R. Cummings *et al.*, "Advances and open problems in federated learning," *Foundations and Trends® in Machine Learning*, vol. 14, no. 1–2, pp. 1–210, 2021.

[25] Y. Zhao, P. Barnaghi, and H. Haddadi, "Multimodal federated learning on iot data," in *2022 IEEE/ACM Seventh International Conference on Internet-of-Things Design and Implementation (IoTDI)*. IEEE, 2022, pp. 43–54.

[26] F. Sattler, S. Wiedemann, K.-R. Müller, and W. Samek, "Robust and communication-efficient federated learning from non-iid data," *IEEE transactions on neural networks and learning systems*, vol. 31, no. 9, pp. 3400–3413, 2019.

[27] W. Y. B. Lim, N. C. Luong, D. T. Hoang, Y. Jiao, Y.-C. Liang, Q. Yang, D. Niyato, and C. Miao, "Federated learning in mobile edge networks: A comprehensive survey," *IEEE Communications Surveys & Tutorials*, vol. 22, no. 3, pp. 2031–2063, 2020.

[28] A. Imteaj, U. Thakker, S. Wang, J. Li, and M. H. Amini, "A survey on federated learning for resource-constrained iot devices," *IEEE Internet of Things Journal*, vol. 9, no. 1, pp. 1–24, 2021.

[29] W. Fang, D.-J. Han, and C. G. Brinton, "Submodel partitioning in hierarchical federated learning: Algorithm design and convergence analysis," *arXiv preprint arXiv:2310.17890*, 2023.

[30] J. Zhang, X. Cheng, C. Wang, Y. Wang, Z. Shi, J. Jin, A. Song, W. Zhao, L. Wen, and T. Zhang, "Fedada: Fast-convergent adaptive federated learning in heterogeneous mobile edge computing environment," *World Wide Web*, vol. 25, no. 5, pp. 1971–1998, 2022.

[31] Y.-W. Chu, S. Hosseinalipour, E. Tenorio, L. Cruz, K. Douglas, A. Lan, and C. Brinton, "Mitigating biases in student performance prediction via attention-based personalized federated learning," in *Proceedings of the 31st ACM International Conference on Information & Knowledge Management*, 2022, pp. 3033–3042.

[32] A. M. Abdelmoniem, C.-Y. Ho, P. Papageorgiou, and M. Canini, "Empirical analysis of federated learning in heterogeneous environments," in *Proceedings of the 2nd European Workshop on Machine Learning and Systems*, 2022, pp. 1–9.

[33] M. Sundararajan and A. Najmi, "The many shapley values for model explanation," in *International conference on machine learning*. PMLR, 2020, pp. 9269–9278.

[34] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," *Advances in neural information processing systems*, vol. 30, 2017.

[35] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee, "From local explanations to global understanding with explainable ai for trees," *Nature Machine Intelligence*, vol. 2, no. 1, pp. 2522–5839, 2020.

[36] Y. J. Cho, J. Wang, and G. Joshi, "Client selection in federated learning: Convergence analysis and power-of-choice selection strategies," *arXiv preprint arXiv:2010.01243*, 2020.

[37] P. Qi, D. Chiaro, and F. Piccialli, "FL-FD: Federated learning-based fall detection with multimodal data fusion," *Information Fusion*, p. 101890, 2023.

[38] B. Xiong, X. Yang, F. Qi, and C. Xu, "A unified framework for multi-modal federated learning," *Neurocomputing*, vol. 480, pp. 110–118, 2022.

[39] T. Feng, D. Bose, T. Zhang, R. Hebbar, A. Ramakrishna, R. Gupta, M. Zhang, S. Avestimehr, and S. Narayanan, "Fedmultimodal: A benchmark for multimodal federated learning," *arXiv preprint arXiv:2306.09486*, 2023.

[40] B. Salehi, J. Gu, D. Roy, and K. Chowdhury, "Flash: Federated learning for automated selection of high-band mmwave sectors," in *IEEE INFOCOM 2022-IEEE Conference on Computer Communications*. IEEE, 2022, pp. 1719–1728.

[41] J. Chen and A. Zhang, "Fedmsplit: Correlation-adaptive federated multi-task learning across multimodal split networks," in *Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining*, 2022, pp. 87–96.

[42] L. Yuan, H. Chen, R. Ewing, and J. Li, "Passive radio frequency-based 3d indoor positioning system via ensemble learning," *arXiv preprint arXiv:2304.06513*, 2023.

[43] S. Yang, L. Yuan, and J. Li, "Extraction and denoising of human signature on radio frequency spectrums," in *2023 IEEE International Conference on Consumer Electronics (ICCE)*. IEEE, 2023, pp. 1–6.

[44] L. Yuan, J. Andrews, H. Mu, A. Vakil, R. Ewing, E. Blasch, and J. Li, "Interpretable passive multi-modal sensor fusion for human identification and activity recognition," *Sensors*, vol. 22, no. 15, p. 5787, 2022.

[45] L. Yuan, H. Chen, R. Ewing, E. Blasch, and J. Li, "Three dimensional indoor positioning based on passive radio frequency signal strength distribution," *IEEE Internet of Things Journal*, vol. 10, no. 15, pp. 13 933–13 944, 2023.

[46] Z. Wang, Z. Wan, and X. Wan, "Transmodality: An end2end fusion method with transformer for multimodal sentiment analysis," in *Proceedings of The Web Conference 2020*, 2020, pp. 2514–2520.

[47] Q. Pan, H. Cao, Y. Zhu, J. Liu, and B. Li, "Contextual client selection for efficient federated learning over edge devices," *IEEE Transactions on Mobile Computing*, 2023.

[48] Y. Xu, Z. Jiang, H. Xu, Z. Wang, C. Qian, and C. Qiao, "Federated learning with client selection and gradient compression in heterogeneous edge systems," *IEEE Transactions on Mobile Computing*, 2023.

[49] B. Luo, W. Xiao, S. Wang, J. Huang, and L. Tassiulas, "Tackling system and statistical heterogeneity for federated learning with adaptive client sampling," in *IEEE INFOCOM 2022-IEEE conference on computer communications*. IEEE, 2022, pp. 1739–1748.

[50] L. Fu, H. Zhang, G. Gao, M. Zhang, and X. Liu, "Client selection in federated learning: Principles, challenges, and opportunities," *IEEE Internet of Things Journal*, 2023.

[51] S. Wang, R. Morabito, S. Hosseinalipour, M. Chiang, and C. G. Brinton, "Device sampling and resource optimization for federated learning in cooperative edge networks," *arXiv preprint arXiv:2311.04350*, 2023.

[52] L. Yuan, L. Su, and Z. Wang, "Federated transfer-ordered-personalized learning for driver monitoring application," *IEEE Internet of Things Journal*, vol. 10, no. 20, pp. 18 292–18 301, May 2023.

[53] J. DelPreto, C. Liu, Y. Luo, M. Foshey, Y. Li, A. Torralba, W. Matusik, and D. Rus, "Actionsense: A multimodal dataset and recording framework for human activities using wearable sensors in a kitchen environment," *Advances in Neural Information Processing Systems*, vol. 35, pp. 13 800–13 813, 2022.

[54] D. Anguita, A. Ghio, L. Oneto, X. Parra, J. L. Reyes-Ortiz *et al.*, "A public domain dataset for human activity recognition using smartphones." in *Esann*, vol. 3, 2013, p. 3.

[55] P. Wagner, N. Strodthoff, R.-D. Bousseljot, D. Kreiseler, F. I. Lunze, W. Samek, and T. Schaeffter, "Ptb-xl, a large publicly available electrocardiography dataset," *Scientific data*, vol. 7, no. 1, p. 154, 2020.

[56] N. Strodthoff, P. Wagner, T. Schaeffter, and W. Samek, "Deep learning for ecg analysis: Benchmarks and insights from ptb-xl," *IEEE Journal of Biomedical and Health Informatics*, vol. 25, no. 5, pp. 1519–1528, 2020.

[57] S. Poria, D. Hazarika, N. Majumder, G. Naik, E. Cambria, and R. Mihalcea, "MELD: A multimodal multi-party dataset for emotion recognition in conversations," in *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*. Florence, Italy: Association for Computational Linguistics, Jul. 2019, pp. 527–536. [Online]. Available: https://aclanthology.org/P19-1050

[58] C. Persello, R. Hänsch, G. Vivone, K. Chen, Z. Yan, D. Tang, H. Huang, M. Schmitt, and X. Sun, "2023 ieee grss data fusion contest: Large-scale fine-grained building classification for semantic urban reconstruction," 2022. [Online]. Available: https://dx.doi.org/10.21227/mrnt-8w27

[59] R. Chen, L. Li, K. Xue, C. Zhang, M. Pan, and Y. Fang, "Energy efficient federated learning over heterogeneous mobile devices via joint design of weight quantization and wireless transmission," *IEEE Transactions on Mobile Computing*, 2022.

**Liangqi Yuan** (Student Member, IEEE) received the B.E. degree from the Beijing Information Science and Technology University, Beijing, China, in 2020, and the M.S. degree from the Oakland University, Rochester, MI, USA, in 2022. He is currently pursuing the Ph.D. degree with the School of Electrical and Computer Engineering, Purdue University, West Lafayette, IN, USA. His research interests are in the areas of sensors, the Internet of Things, signal processing, and machine learning.

**Dong-Jun Han** (Member, IEEE) received the B.S. degrees in mathematics and electrical engineering, and the M.S. and Ph.D. degrees in electrical engineering from Korea Advanced Institute of Science and Technology (KAIST), South Korea, in 2016, 2018, and 2022, respectively. He received the Best Ph.D. Dissertation Award from the School of Electrical Engineering at KAIST in 2022. He is currently a postdoctoral researcher in the School of Electrical and Computer Engineering at Purdue University. His research interest is at the intersection of communications, networking, and machine learning, specifically in distributed/federated machine learning and network optimization.

**Su Wang** received his B.S. and Ph.D. in Electrical Engineering from Purdue University, West Lafayette, IN, USA, in 2018 and 2023, respectively. He is currently a postdoctoral researcher in the School of Electrical and Computer Engineering at Princeton University.

**Devesh Upadhyay** (Senior Member, IEEE) received the M.S. and Ph.D. degree in mechanical engineering from The Ohio State University, Columbus. He is currently the Technical Director for AI/ML and Autonomy at Saab Inc. Before joining Saab Devesh was a Senior Technical Leader at Ford Research where he led the Core AI/ML Quantum Computing team.

**Christopher G. Brinton** (Senior Member, IEEE) is the Elmore Rising Star Assistant Professor of Electrical and Computer Engineering (ECE) at Purdue University. His research interest is at the intersection of networking, communications, and machine learning, specifically in fog/edge network intelligence, distributed machine learning, and data-driven wireless network optimization. Dr. Brinton is a recipient of the NSF CAREER Award, ONR Young Investigator Program (YIP) Award, DARPA Young Faculty Award (YFA), Intel Rising Star Faculty Award, and roughly $15M in sponsored research projects as a PI or co-PI. He has also been awarded Purdue College of Engineering Faculty Excellence Awards in Early Career Research, Early Career Teaching, and Online Learning. He currently serves as an Associate Editor for IEEE/ACM Transactions on Networking. Prior to joining Purdue, Dr. Brinton was the Associate Director of the EDGE Lab and a Lecturer of Electrical Engineering at Princeton University. He also co-founded Zoomi Inc., a big data startup company that has provided learning optimization to more than one million users worldwide and holds US Patents in machine learning for education. His book The Power of Networks: 6 Principles That Connect our Lives and associated Massive Open Online Courses (MOOCs) have reached over 400,000 students to date. Dr. Brinton received the PhD (with honors) and MS Degrees from Princeton in 2016 and 2013, respectively, both in Electrical Engineering.