# Federated Model Aggregation via Self-Supervised Priors for Highly Imbalanced Medical Image Classification

Marawan Elbatel[1,2], Hualiang Wang[1], Robert Martí[2], Huazhu Fu[3], and Xiaomeng Li[1]

[1] The Hong Kong University of Science and Technology
[2] Computer Vision and Robotics Institute, University of Girona
[3] Institute of High Performance Computing (IHPC), Agency for Science, Technology and Research (A*STAR), Singapore

**Abstract.** In the medical field, federated learning commonly deals with highly imbalanced datasets, including skin lesions and gastrointestinal images. Existing federated methods under highly imbalanced datasets primarily focus on optimizing a global model without incorporating the intra-class variations that can arise in medical imaging due to different populations, findings, and scanners. In this paper, we study the inter-client intra-class variations with publicly available self-supervised auxiliary networks. Specifically, we find that employing a shared auxiliary pre-trained model, like MoCo-V2, locally on every client yields consistent divergence measurements. Based on these findings, we derive a dynamic balanced model aggregation via self-supervised priors (MAS) to guide the global model optimization. Fed-MAS can be utilized with different local learning methods for effective model aggregation toward a highly robust and unbiased global model. Our code is available at https://github.com/xmed-lab/Fed-MAS.

## 1 Introduction

Federated learning (FL) has emerged as a way to train models with decentralized data while preserving privacy. Due to the inherent nature of data heterogeneity in medical imaging, training in a decentralized manner exhibits performance degradation compared to centralized training. With FedAvg [23] as the main baseline, multiple works proposed to improve the model's generic performance under data decentralization [19,20,24]. These methods have been successful in achieving positive results, assuming a balanced global data distribution. However, they struggle to address extreme data heterogeneity, especially in highly imbalanced medical datasets. There have been some methods proposed to address the imbalanced setting [25,21]. Nevertheless, these methods shared local features among clients, which may raise privacy concerns.
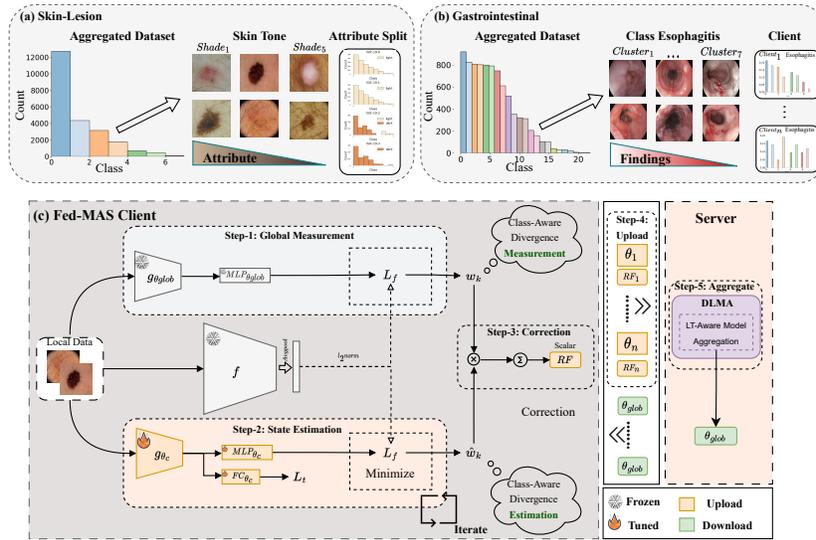
**Fig. 1.** (a) Skin lesion attribute imbalance, (b) Gastrointestinal findings imbalance (Ex: tracheleriazation, varices, leukoplakia). (c) Fed-MAS framework.

Label distribution skewness has been studied in the context of FL [33]. FedLC [33], inspired by LDAM [3], showed promising results by adjusting the local client class distribution. Additionally, multiple works proposed to tackle the issue of highly skewed label distribution (i.e. long-tailed) by decoupling the classifier and the feature extractor [29,32,5]. The rationale behind these methods is rooted in the understanding that the classifier is the bottleneck for majority class bias [18]. For instance, CReFF [29] retrained a balanced classifier on the server by leveraging federated features. A notable limitation of classifier re-training is its inability to address the intra-class attribute imbalance. Recently, [30] showed that training with imaging data with high attribute imbalance impedes representation learning by exacerbating the intra-class variations. In FL, the issue of intra-class imbalance is critical when dealing with highly imbalanced medical imaging datasets. As depicted in Fig. 1 (a), different skin tones can arise across different clients for the same class [1]. For gastrointestinal recognition depicted in Fig. 1 (b), different findings can arise in different clients for the same class [2]. Hence, the challenge of an unbiased robust global model that takes into account both the attribute and class imbalance still remains. More recently, FedCE [16] showed promising results by calculating a fair client contribution estimation in gradient and data space for medical image segmentation; Nevertheless, it relies on local validation samples, which may not adequately represent attribute imbalance and rare diseases in highly imbalanced medical image datasets.

Publicly available pre-trained models, such as MoCo-V2 [12] that were trained without any labels using a large set of naturals images, have been utilized with their batch statistics in calculating image priors [11] and have been utilized with

their generalizable representation to improve the performance in highly imbalanced medical imaging tasks [8]. In this paper, we leverage these pre-trained models locally to propose Fed-MAS as a novel approach to incorporate the client's local variations with consistent self-supervised priors, estimating client contributing ratios toward an unbiased robust global model.

## 2    Methodology

Figure 1 shows the overview of our Fed-MAS framework. Each local client is provided with a publicly self-supervised pre-trained model (e.g., MoCo-RN50 [12]) that is not involved in the training or communication process of the federated learning framework. Consequently, these pre-trained models do not increase communication costs while ensuring that each client can access the same consistent pre-trained model. With $n$ local clients and one global server, Fed-MAS performs the following steps in each round: (1) Each client receives the global model to measure its global class-aware divergence, $w_k$, and update its local model; (2) Each client trains its local model while estimating its class-aware divergence, $\hat{w}_k$; (3) Each client corrects $\hat{w}_k$ with $w_k$ to generate a rescue scalar, $RF$; (4) Client uploads the parameters of its local model and $RF$ to the server; (5) The server applies our proposed MAS to aggregate a new model from the parameters of the received client models, weighted by $RF$;

### 2.1    Class Aware Global Observation via Self-Supervised Priors

In highly imbalanced medical image datasets, both extreme class imbalance and inter-client intra-class variations can lead to client drift. Due to the decentralization of data, estimating the global intra-class attribute distribution in medical imaging within the FL framework is a challenge that is yet to be explored.

At the beginning of each round in the FL process, each client receives the model from the global server $\theta_{global}$. We study locally the distance between the distribution of the self-supervised pre-trained model, $f_\xi$, and $\theta_{glob}$ over each client's local data.

Given an input image $x$, we feed $x$ to the local feature encoder $g$ to generate a representation $z = g_\theta(x)$. This representation is then fed to an $MLP$ projector to generate a projection $y = MLP_\theta(z)$ in a space comparable with the self-supervised model. From the same discriminative pre-trained model in all clients, we can get a target representation $y' = f_\xi(x)$, where both $y$ and $y'$ are L2 normalized. We can measure the distribution difference using mean squared error as:

$$\mathcal{L}_f^\theta = 2 - 2 \cdot \langle y, y' \rangle \cdot \tag{1}$$

From Eq. (1), we can generate a class-aware distance for class $k$ with $M_k$ total samples as:

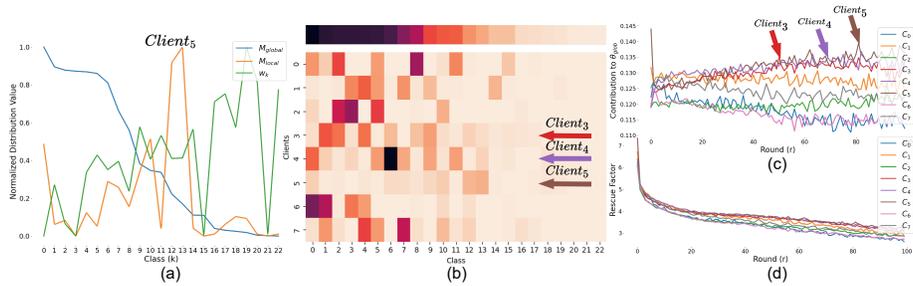$$\mathcal{L}_k^\theta = \frac{1}{M_k} \sum_{i=1}^{M_k} \mathcal{L}_f^\theta(x_{k,i}). \tag{2}$$

**Fig. 2.** Analysis of MAS on HyperKvasir: (a) The globally aggregated class counts, $M_{global}$, client local count, $M_{local}$, and $w_k$ in one round. (b) HyperKvasir non-IID setting, (c) Client's Contribution to $\theta_{glob}$ throughout rounds, (d) Rescue Factor (RF) on different clients throughout rounds

We define $w_k = \mathcal{L}_k^{\theta_{glob}}$. The factor $w_k$ can help to capture the distance in distribution between the global server and the self-supervised model on each client's local data. This divergence can provide insights into the sensitivity of the global model, $\theta_{glob}$, in effectively capturing the specific class attribute in each client's local data. A high $w_k$ indicates the failure of $\theta_{glob}$ in capturing a local class $k$. In Fig. 2 (a), we can see that $w_k$ is inversely proportional to the global class distribution, even if the local client distribution is not necessarily the same.

### 2.2 State Estimation via Knowledge Distillation

While $w_k$ provides class-aware global divergence measurement with the same consistent local frozen self-supervised model, a client receives the global model, $\theta_{glob}$, and takes subsequent optimization steps for $E$ local epochs with uncertainty to generate $\theta_c'$. Hence, the client's drift from the global model is hideous after its uncertain optimization.

With a running average, a client can provide a class-aware divergence likelihood $\hat{w}_k$, where $\hat{w}_k = \sum_{e=1}^{E} \mathcal{L}_k^{\theta_c'}$. The factor $\hat{w}_k$ can help to capture how far the client drifted from $f_\xi$ since the global measurement, $w_k$, was taken. A client can then correct this estimation, $\hat{w}_k$, with the global observation, $w_k$, to generate a posterior rescue factor, $RF$, in every round.

$$RF = \sum_{k=1}^{K} w_k \hat{w}_k. \tag{3}$$

A higher $RF$ indicates that the client has information that the global model has not appropriately captured.

To train the projector $MLP_\theta(\cdot)$, we propose to minimize Eq. (1) along with the local balanced risk minimization [28] to minimize a total loss $\mathcal{L}_{total}$ concerning $\theta$ only as:

$$\mathcal{L}_{total} = \mathcal{L}_{sup} + \lambda_f \; \mathcal{L}_f, \tag{4}$$
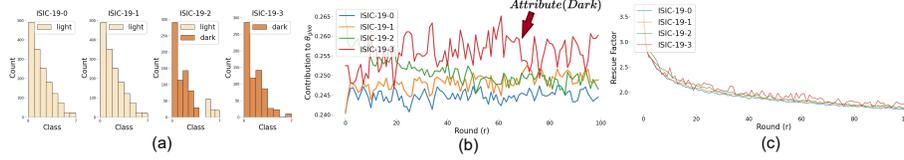
**Fig. 3.** (a) ISIC-FL Attribute Split, (b) Client's Contribution to $\theta_{glob}$ throughout rounds, (c) Rescue Factor (RF) on different clients throughout rounds

where $\mathcal{L}_{sup}$ refers to the original supervised loss and $\lambda_f$ as a weighting factor. This can be seen as restricting the client optimization direction. However, the self-supervised model ensures clients align with a common reference distribution and possess implicit regularization capabilities for minority classes through generalizable features [9].

### 2.3   Model Aggregation via Self-Supervised Posteriors

Inspired by the fact that client-specific models should contribute more to the global server to capture local variance, we propose a novel model aggregation via the corrected self-supervised posteriors (MAS) . We use our proposed $RF$ to indicate client-specific models that should contribute more to the global model than client-generic models to capture their attribute-class variations. While our proposed $RF$ can be used for biased client selection [15], we use it to aggregate a global model. Instead of aggregating based on the weighted samples as in FedAvg [23], we propose to weight the global model, $\theta_{glob}$, based on the $RF$ value as follows:

$$\bar{RF}_c = \frac{RF_c}{\sum_j RF_j} \text{ , and } \theta_{glob}^{r+1} = \sum_{c=1}^{C} \bar{RF}_c \theta_c'. \tag{5}$$

For instance, Client 3,4,5 in Fig. 2 (b) have mostly minority classes and contribute the most to $\theta_{glob}$ in Fig. 2 (c). Morever, in Fig. 3 (a) Client ISIC-3 have mostly underrepresented attribute and contributes the most in  Fig. 3 (b). Additionally, we show in Fig. 2 (d) and  Fig. 3 (c) that the rescue factor for all clients is decreasing throughout rounds. This highlights the ability of MAS to accommodate different clients. (See Algorithm 1 in Appendix).

## 3   Experiments

**Dataset. HyperKvasir** [2] is a long-tailed (LT) dataset of 10,662 gastrointestinal tract images with 23 classes from different anatomical and pathological landmarks and findings. We divide the 23 classes into Head ($> 700$ images per class), Medium ($70 \sim 700$ images per class), and Tail ($< 70$ images per class) with respect to their class counts. Additionally, we partition the data across eight

**Table 1.** Comparison with other methods on HyperKvasir Dataset. All clients are initialized with ImageNet pre-trained weights; each result is averaged over five runs.

| Methods | IID | | | | | non-IID $Dir(\alpha = 0.5)$ | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | Head | Medium | Tail | All | B-acc | Head | Medium | Tail | All | B-acc |
| Federated Learning Methods (FL-Methods) | | | | | | | | | | |
| FedAvg [23] | $94.1 \pm 1.3$ | $72.9 \pm 1.3$ | $3.1 \pm 0.9$ | $56.69 \pm 0.6$ | $58.1 \pm 0.6$ | $86.2 \pm 2.7$ | $70.3 \pm 0.5$ | $8.0 \pm 1.2$ | $54.83 \pm 1.0$ | $56.17 \pm 0.9$ |
| FedProx [20] | $94.6 \pm 0.4$ | $72.1 \pm 0.2$ | $3.0 \pm 1.2$ | $56.58 \pm 0.4$ | $57.93 \pm 0.4$ | $88.1 \pm 2.2$ | $73.1 \pm 2.7$ | $3.6 \pm 2.5$ | $54.93 \pm 1.3$ | $56.51 \pm 1.3$ |
| MOON [19] | $94.7 \pm 0.7$ | $74.6 \pm 0.4$ | $4.0 \pm 1.8$ | $57.77 \pm 0.6$ | $59.23 \pm 0.6$ | $84.4 \pm 3.6$ | $73.1 \pm 1.6$ | $5.5 \pm 2.1$ | $54.3 \pm 1.2$ | $55.93 \pm 1.1$ |
| LT-integerated FL Methods | | | | | | | | | | |
| LDAM-FL [3] | $95.4 \pm 0.5$ | $72.2 \pm 1.1$ | $5.7 \pm 3.9$ | $57.77 \pm 1.4$ | $59.03 \pm 1.3$ | $86.9 \pm 2.8$ | $70.9 \pm 1.2$ | $4.7 \pm 4.6$ | $54.16 \pm 1.4$ | $55.61 \pm 1.4$ |
| BSM-FL [28] | $93.2 \pm 1.5$ | $74.6 \pm 2.6$ | $9.1 \pm 3.7$ | $58.92 \pm 0.6$ | $60.28 \pm 0.7$ | $89.6 \pm 3.9$ | $68.7 \pm 3.0$ | $16.4 \pm 5.4$ | $58.24 \pm 1.2$ | $59.15 \pm 1.3$ |
| Label-Skew FL Methods | | | | | | | | | | |
| CReFF [29] | $95.1 \pm 0.8$ | $72.0 \pm 1.5$ | $2.6 \pm 1.8$ | $56.53 \pm 1.4$ | $57.88 \pm 1.4$ | $89.3 \pm 0.7$ | $70.1 \pm 1.6$ | $9.0 \pm 4.5$ | $56.12 \pm 1.3$ | $57.34 \pm 1.2$ |
| FedLC [33] | $96.5 \pm 0.4$ | $75.3 \pm 2.5$ | $7.4 \pm 5.5$ | $59.73 \pm 1.8$ | $61.08 \pm 1.7$ | $95.8 \pm 0.6$ | $73.1 \pm 2.4$ | $6.6 \pm 4.1$ | $58.51 \pm 1.5$ | $59.78 \pm 1.5$ |
| **Fed-Mas (ours)** | $94.3 \pm 1.2$ | $72.9 \pm 1.0$ | $15.9 \pm 2.7$ | $\mathbf{61.05 \pm 0.3}$ | $\mathbf{62.08 \pm 0.2}$ | $93.0 \pm 0.9$ | $72.5 \pm 2.6$ | $16.2 \pm 1.3$ | $\mathbf{60.57 \pm 1.1}$ | $\mathbf{61.61 \pm 1.0}$ |

clients with IID (similar label distributions) and non-IID (heterogeneous partition with Dirichlet distribution). **ISIC** [7] is a highly imbalanced dataset of skin lesion images with 8 classes that exhibits skin-tone attribute imbalance [1]. For instance, melanoma incidence is lower in quantity and higher in mortality rates in black patients than in others [6]. We partition the dataset on four clients based on two attributes, light and dark skin tones, with [1] labeling. Additionally, we split the data between the four clients for training, validation, and testing with 70%, 15%, and 15%, respectively. We also benchmark Fed-MAS over Flamby-ISIC split [31] with six different hospitals with stratified 5-fold cross-validation. **Implementation Details.** For both datasets, we use resnet-18 [13] as the local target model. For the long-tailed HyperKvasir dataset, we employ an SGD optimizer and a cosine annealing scheduler [22] with a maximum learning rate of 0.1. For ISIC, we employ Adam optimizer with the 3e-4 learning rate. Additionally, we employ balanced risk minimization [28] and train methods for 200 communication rounds with 10 local epochs. We set $\lambda_f$ to 3 and provide an ablation in Tab. 4 in Appendix.

**Evaluation Metrics.** We evaluate the model performance of the global model in this paper. To assess the unequal treatment of each class in HyperKvasir, we report the top-1 accuracy on shot-based division (head, medium, tail) and their average results denoted as "All" as existing works [17]. Following prior work [10,14,27], we also report the Balanced Accuracy "B-Acc", which calculates the average per-class accuracy and is resistant to class imbalance. As the test set of HyperKvasir contains only 12 classes, we follow previous work [10] to assess the model performance with a stratified 5-fold cross-validation. To evaluate the performance of attributes in ISIC-FL, we report the "B-Acc" separately for each attribute ("Light", "Dark") and the average of these scores "Avg". Additionally, we report the overall "B-Acc" across all attributes and distributions.

### 3.1    Performance on the HyperKvasir

We compare our methods with FL-methods [23,20,19], LT-integrated FL methods [3,28], and label-skew FL methods [33,29]

**FL-Methods** [23,20,19]. One simple solution for federated learning with highly imbalanced medical data is to apply existing FL methods to our setting directly.

To this end, we compare our methods with state-of-the-art FL methods, including FedAvg [23], FedProx [20], and MOON [19], under the same setting. As shown in Table 1, we find that our method outperforms the best existing FL method MOON by 2.85% and 5.68% on "B-acc" in both IID and non-IID settings, respectively. Notably, our Fed-MAS achieves similar results with MOON [19] on the "Head" while reaching large improvements on the "Tail" (11.9% on iid and 10.71% on non-iid), showing that our Fed-MAS can tackle LT distribution under FL more effectively. The limited results could be attributed to the use of local empirical risk minimization in MOON [19]. However, even when we applied a balanced risk minimization [28] in MOON, our method still outperformed it (60.69% vs. 62.08% on "B-acc" for IID); see results in Table 6 in Appendix.

**LT integrated FL methods** [3,28]. To design FL methods for local clients with long-tailed distribution, a straightforward idea is to directly use LT methods in each local client and then use an FL framework such as FedAvg to obtain the final results. In this regard, we implement LDAM-DRW [3] and BSM [28] into the FedAvg framework and rename them as LDAM-FL and BSM-FL respectively. From Table 1, we can notice the LT methods utilizing an FL framework have produced limited results on "Tail" primarily due to the extreme client drifting phenomenon. Please note that Fed-MAS does not focus on designing any specific long-tailed training for each local client. Instead, MAS enables the global server to effectively aggregate the model parameters from long-tailed distributed local clients. As a result, our Fed-MAS can successfully capture the "Tail" with a 6.84% accuracy gain on IID with lower variance than the best-performing LT method BSM-FL [28]. Notably, our method consistently outperforms the best-performing LT method on the "B-acc" with a lower variance (improvement of 1.8% on IID and 2.46% on non-IID).

**Label-Skew FL** We compare our method with the state-of-the-art label-skew FL method, FedLC [33], and the highly labeled skew (i.e. long-tailed) FL method, CReFF [29]. CReFF, as proposed by [29], involves a method of re-training the classifier by utilizing learnable features on the server at each communication round, holding an equal treatment of all clients' models. However, this technique fails to accommodate inter-client intra-class variations which could arise. From Table 1, we can notice that FedAvg with local LT such as BSM-FL [28] can outperform CReFF [29] on the HyperKvasir dataset in both IID and non-IDD by 2.4% and 1.8% on "B-acc", respectively. Our comparative analysis illustrates that Fed-MAS consistently outperforms CReFF in both IID and non-IID by 4.2% and 4.27% on "B-acc", respectively, by incorporating the client's local variations with MAS. FedLC [33] proposes a loss function to address label distribution skewness by locally calibrating logits and reducing local bias in learning. Their modification yields compelling performance. Nevertheless, our method surpasses them in both IID and non-IID, achieving improvements of 1.0% and 1.83% on "B-Acc", respectively. Remarkably, our method effectively captures the tail classes with reduced variance in both IID and non-IID, exhibiting improvements of 8.5% and 9.6%, respectively, while experiencing only a minor drop in performance for the head classes (96.5% vs 94.3% for IID and 95.8% vs 93.0% for non-IID).

**Table 2.** Ablation of minimizing Eq. (1) (KD) and MAS on HyperKvasir non-IID

| | KD | MAS | Metrics | | |
|---|---|---|---|---|---|
| | | | All (%) | B-acc (%) | p-value |
| BSM-FL [28] (Baseline) | × | × | 58.24 ± 1.2 | 59.15 ± 1.3 | — |
| [28] w/ KD | ✓ | × | 59.26 ± 1.2 | 60.19 ± 1.1 | <0.001 |
| **Fed-MAS** | ✓ | ✓ | **60.57 ± 1.1** | **61.61 ± 1.0** | <0.001 |

**Table 3.** Experimental Results on ISIC-FL. Results are averaged over 5 folds.

| Method | Attribute Setting (ours) | | | | Flamby-ISIC [31] |
|---|---|---|---|---|---|
| | Light | Dark | Avg | B-Acc | B-Acc |
| | With ImageNet Weight Initialization | | | | |
| FedLC [33] | 71.11 ± 1.8 | 73.64 ± 6.6 | 72.38 ± 2.9 | 71.63 ± 1.6 | 76.54 ± 2.6 |
| BSM-FL [28] (Baseline) | 73.88 ± 1.4 | 74.78 ± 5.4 | 74.33 ± 2.5 | 74.49 ± 1.4 | 78.19 ± 1.8 |
| [28] w/ KD | 73.87 ± 1.6 | 72.44 ± 5.9 | 73.16 ± 3.0 | 74.09 ± 1.5 | 79.17 ± 2.1 |
| **Fed-Mas (ours)** | 73.43 ± 1.6 | 77.0 ± 6.6 | **75.21 ± 2.9** | **74.61 ± 1.4** | **80.87 ± 2.2** |

**Effectiveness of KD and MAS** As shown in Table 2, minimizing Eq. (1) (KD) can enhance the "All" and "B-Acc" via 1.02% and 1.04% due to the implicit regularization of MoCo-V2 on the tail classes for extreme imbalance datasets. With both KD and MAS, the performance is further improved to the best via 2.33% and 2.46% on "All" and "B-Acc", respectively. MAS utilizes unbiased frozen generalizable representations to incorporate the inter-client intra-class characteristics in FL and combine them with the drifting belief. This combination helps in capturing client-specific models in the aggregation step.

### 3.2   Performance on ISIC-FL

We evaluate the best-performing and competitive methods with the ISIC-FL dataset to shorten the benchmark. While previous studies neglect weight initialization to provide better convergence analysis as pre-trained weights are architecture dependent. Recently, [26] and [4] studied the impact of pre-training initialization on reducing the data and system heterogeneity in FL. We present in Table 3 the results of the most competitive methods with weight initialization on the ISIC-FL attribute setting. FedLC [33] demonstrates compelling performance to address label skewness in Hyperkvasir-FL. Nevertheless, it falls short in accommodating attribute heterogeneity in ISIC-FL due to its local learning focus. Our method consistently outperforms FedLC [33] with a notable improvement of 2.8% and 3.0% in terms of the averaged balanced accuracies "Avg" and balanced accuracy "B-acc" respectively. Compared to the baseline [28], Fed-MAS notably captured the underrepresented attribute with 2.2% on the "B-acc" of the "Dark Attribute" with a minimal drop of 0.5% on the "B-acc" of the "Light Attribute", balancing the intra-class attribute characteristics in FL. On the highly heterogeneous Flamby-ISIC split resembling six hospitals, Fed-MAS outperform FedLC and the baseline on the "B-acc" with 4.33% and 2.68%, respectively.

### 3.3   Privacy Concerns

Similarly to traditional FL methods [23,19,20], Fed-MAS shares the model weights with an additional scalar, $RF$, which protects data privacy by not revealing in-

put data or label distribution. The *scalar*, $RF$, is calculated in the output feature space, safeguarding the input data distribution. Moreover, $RF$ poses uncertainty in approximating the client's label distribution as it can be influenced by diverse attributes in the majority class or a common attribute in the minority class.

## 4   Conclusion

Highly Imbalanced datasets are present in most medical image classifications. This work presents Fed-MAS to deal with this problem. We show that publicly available self-supervised models benefit the FL training procedure more than restricting the optimization direction by incorporating the global attribute imbalance. Future work can explore delayed re-weighting to unleash non-vanishing terms and explore MAS with different local learning strategies in FL settings.

## References

1. Bevan, P.J., Atapour-Abarghouei, A.: Detecting melanoma fairly: Skin tone detection and debiasing for skin lesion classification. In: Kamnitsas, K., Koch, L., Islam, M., Xu, Z., Cardoso, J., Dou, Q., Rieke, N., Tsaftaris, S. (eds.) Domain Adaptation and Representation Transfer. pp. 1–11. Springer Nature Switzerland, Cham (2022)
2. Borgli, H., Thambawita, V.L., Smedsrud, P.H., Hicks, S., Jha, D., Eskeland, S.L., Randel, K.R., Pogorelov, K., Lux, M., Nguyen, D.T.D., Johansen, D., Griwodz, C., Stensland, H.K., Garcia-Ceja, E., Schmidt, P.T., Hammer, H.L., Riegler, M., Halvorsen, P., de Lange, T.: Hyperkvasir, a comprehensive multi-class image and video dataset for gastrointestinal endoscopy. Scientific Data **7** (2019)
3. Cao, K., Wei, C., Gaidon, A., Arechiga, N., Ma, T.: Learning imbalanced datasets with label-distribution-aware margin loss. In: NeurIPS (2019)
4. Chen, H.Y., Tu, C.H., Li, Z., Shen, H.W., Chao, W.L.: On the importance and applicability of pre-training for federated learning. In: ICLR (2023)
5. Chen, Z., Liu, S., Wang, H., Yang, H.H., Quek, T.Q.S., Liu, Z.: Towards federated long-tailed learning. ArXiv **abs/2206.14988** (2022)
6. Collins, K.K., Fields, R.C., Baptiste, D.F., Liu, Y., Moley, J.F., Jeffe, D.B.: Racial differences in survival after surgical treatment for melanoma. Annals of Surgical Oncology **18**, 2925–2936 (2011)
7. Combalia, M., Codella, N.C.F., Rotemberg, V.M., Helba, B., Vilaplana, V., Reiter, O., Halpern, A.C., Puig, S., Malvehy, J.: Bcn20000: Dermoscopic lesions in the wild. ArXiv **abs/1908.02288** (2019)
8. Ding, X., Liu, Z., Li, X.: Free lunch for surgical video understanding by distilling self-supervisions. In: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (eds.) MICCAI 2022. Lecture Notes in Computer Science, vol. 13437, pp. 365–375. Springer (2022)

9. Elbatel, M., Martí, R., Li, X.: Fopro-kd: Fourier prompted effective knowledge distillation for long-tailed medical image recognition. ArXiv **abs/2305.17421** (2023)

10. Galdran, A., Carneiro, G., González Ballester, M.A.: Balanced-mixup for highly imbalanced medical image classification. In: de Bruijne, M., Cattin, P.C., Cotin, S., Padoy, N., Speidel, S., Zheng, Y., Essert, C. (eds.) Medical Image Computing and Computer Assisted Intervention – MICCAI 2021. pp. 323–333. Springer International Publishing, Cham (2021)

11. Hatamizadeh, A., Yin, H., Roth, H.R., Li, W., Kautz, J., Xu, D., Molchanov, P.: Gradvit: Gradient inversion of vision transformers. 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) pp. 10011–10020 (2022)

12. He, K., Fan, H., Wu, Y., Xie, S., Girshick, R.B.: Momentum contrast for unsupervised visual representation learning. CVPR pp. 9726–9735 (2020)

13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR. pp. 770–778 (2016)

14. Holste, G., Wang, S., Jiang, Z., Shen, T.C., Shih, G., Summers, R.M., Peng, Y., Wang, Z.: Long-tailed classification of thorax diseases on chest x-ray: A new benchmark study. In: Nguyen, H.V., Huang, S.X., Xue, Y. (eds.) Data Augmentation, Labelling, and Imperfections. pp. 22–32. Springer Nature Switzerland, Cham (2022)

15. Jee Cho, Y., Wang, J., Joshi, G.: Towards understanding biased client selection in federated learning. In: Camps-Valls, G., Ruiz, F.J.R., Valera, I. (eds.) Proceedings of The 25th International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research, vol. 151, pp. 10351–10375. PMLR (28–30 Mar 2022)

16. Jiang, M., Roth, H.R., Li, W., Yang, D., Zhao, C., Nath, V., Xu, D., Dou, Q., Xu, Z.: Fair federated medical image segmentation via client contribution estimation. In: CVPR (2023)

17. Ju, L., Wu, Y., Wang, L., Yu, Z., Zhao, X., Wang, X., Bonnington, P., Ge, Z.: Flexible sampling for long-tailed skin lesion classification. In: Wang, L., Dou, Q., Fletcher, P.T., Speidel, S., Li, S. (eds.) MICCAI 2022. pp. 462–471. Springer Nature Switzerland, Cham (2022)

18. Kang, B., Xie, S., Rohrbach, M., Yan, Z., Gordo, A., Feng, J., Kalantidis, Y.: Decoupling representation and classifier for long-tailed recognition. In: ICLR (2020)

19. Li, Q., He, B., Song, D.: Model-contrastive federated learning. In: CVPR (2021)

20. Li, T., Sahu, A.K., Zaheer, M., Sanjabi, M., Talwalkar, A., Smith, V.: Federated optimization in heterogeneous networks. In: Dhillon, I., Papailiopoulos, D., Sze, V. (eds.) Proceedings of Machine Learning and Systems. vol. 2, pp. 429–450 (2020)

21. Liu, Q., Yang, H., Dou, Q., Heng, P.A.: Federated semi-supervised medical image classification via inter-client relation matching. In: de Bruijne, M., Cattin, P.C., Cotin, S., Padoy, N., Speidel, S., Zheng, Y., Essert, C. (eds.) MICCAI 2021. pp. 325–335. Springer International Publishing, Cham (2021)

22. Loshchilov, I., Hutter, F.: SGDR: Stochastic gradient descent with warm restarts. In: ICLR (2017)

23. McMahan, H.B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A.: Communication-efficient learning of deep networks from decentralized data. In: AISTATS (2017)

24. Mendieta, M., Yang, T., Wang, P., et al.: Local learning matters: Rethinking data heterogeneity in federated learning. In: CVPR. pp. 8397–8406 (2022)

25. Mu, X., Shen, Y., Cheng, K., Geng, X., Fu, J., Zhang, T., Zhang, Z.: Fedproc: Prototypical contrastive federated learning on non-iid data. Future Gener. Comput. Syst. **143**, 93–104 (2021)

26. Nguyen, J., Wang, J., Malik, K., Sanjabi, M., Rabbat, M.: Where to begin? on the impact of pre-training and initialization in federated learning. In: ICLR (2023)
27. Reinke, A., Christodoulou, E., Glocker, B., et al.: Metrics reloaded - a new recommendation framework for biomedical image analysis validation. In: Medical Imaging with Deep Learning (2022)
28. Ren, J., Yu, C., Sheng, S., Ma, X., Zhao, H., Yi, S., Li, H.: Balanced meta-softmax for long-tailed visual recognition. In: Proceedings of Neural Information Processing Systems(NeurIPS) (Dec 2020)
29. Shang, X., Lu, Y., Huang, G., Wang, H.: Federated learning on heterogeneous and long-tailed data via classifier re-training with federated features. In: Raedt, L.D. (ed.) IJCAI. pp. 2218–2224 (7 2022)
30. Tang, K., Tao, M., Qi, J., Liu, Z., Zhang, H.: Invariant feature learning for generalized long-tailed classification. In: ECCV. p. 709–726 (2022)
31. Ogier du Terrail, J., Ayed, S.S., et al.: Flamby: Datasets and benchmarks for cross-silo federated learning in realistic healthcare settings. In: NeurIPS. vol. 35, pp. 5315–5334. Curran Associates, Inc. (2022)
32. Wicaksana, J., Yan, Z., Cheng, K.T.: Fca: Taming long-tailed federated medical image classification by classifier anchoring. ArXiv **abs/2305.00738** (2023)
33. Zhang, J., Li, Z., et al.: Federated learning with label distribution skew via logits calibration. vol. 162, pp. 26311–26329. Proceedings of Machine Learning Research (17–23 Jul 2022)

# Appendix for "Fed-MAS"

Marawan Elbatel[1,2], Hualiang Wang[1], Robert Martí[2], Huazhu Fu[3], and Xiaomeng Li[1]

[1] The Hong Kong University of Science and Technology
[2] Computer Vision and Robotics Institute, University of Girona
[3] Institute of High Performance Computing (IHPC), Agency for Science, Technology and Research (A*STAR), Singapore

**Table 6.** HyperKvasir FL methods with local BRM [28].

| Method | All | | B-Acc | |
|--------|-----|-----|-------|-----|
| | IID | non-IID | IID | non-IID2 |
| FedAvg | 58.92 | 58.24 | 60.28±0.6 | 59.15±1.3 |
| FedProx | 59.37 | 58.86 | 60.47±1.3 | 59.64±2.0 |
| Moon | 59.45 | 58.72 | 60.69±0.9 | 59.66±0.8 |
| Ours | **61.05** | **60.57** | **62.08±0.2** | **61.61±1.4** |

**Table 4.** HyperKvasir $\lambda_f$ ablation.

| Method | IID | | | non-IID | | |
|--------|-----|-----|-----|---------|-----|-----|
| | $\lambda_f = 0$ | $\lambda_f = 1$ | $\lambda_f = 3$ | $\lambda_f = 0$ | $\lambda_f = 1$ | $\lambda_f = 3$ |
| Fed-MAS | 60.28 | 61.43 | **62.08** | 59.15 | 61.08 | **61.61** |

**Table 5.** HyperKvasir $f_\xi$ ablation non-IID.

| $f_\xi$ | All | B-Acc |
|---------|-----|-------|
| CLIP-ViTB/32 | 60.34 | 61.39±2.1 |
| MoCo-RN50 | 60.57 | 61.61±1.0 |

**Table 7.** Using a plug-in cRT [18] on HyperKvasir on non-IID.

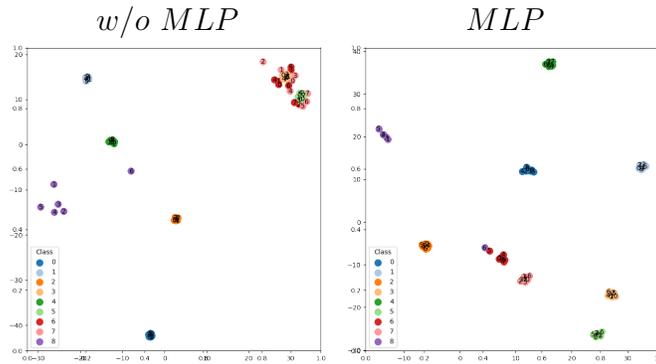| Method + cRT | All | B-Acc |
|--------------|-----|-------|
| Decoupling [18] | 54.21 | 55.6 |
| BSM-FL [28] | 62.67 | 63.11 |
| **Ours** | **65.05** | **65.11** |



**Fig. 4.** Feature representation with and without the learnable projector $MLP_\theta$. We sample a subset of head (0,1,2), medium (3,4,5), and tail (6,7,8) classes for feature visualization across different clients. Each point represents the mean feature output for each class (color) in each client (point).

**Table 8.** Flamby-ISIC [31] results on the first fold with the global model (gFL) and the local models (pFL) with ImageNet Weight Initialization. MOON [19] and FedProx [20] are reported with local BRM [28].

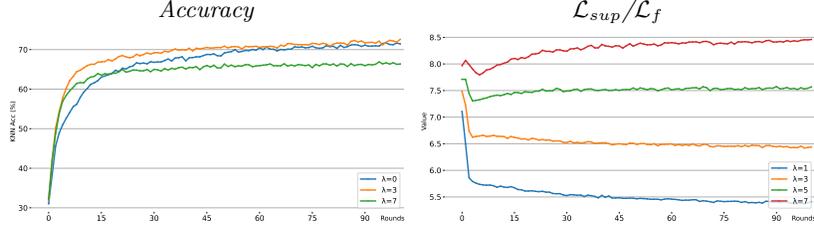| Method | Metric | | Method | Metric | |
|---|---|---|---|---|---|
| | gFL | pFL | | gFL | pFL |
| MOON ($\mu = 0.01$) | 72.13 | 80.03 | FedProx ($\mu = 0.1$) | 72.47 | 79.82 |
| MOON ($\mu = 0.1$) | 72.45 | 79.64 | FedProx ($\mu = 0.01$) | 73.25 | 79.82 |
| MOON ($\mu = 1$) | 73.12 | 79.46 | FedProx ($\mu = 0.001$) | 73.52 | 79.70 |
| FedLC [33] | 68.07 | 78.60 | BSM-FL [28] | 72.83 | 79.79 |
| [28] w/ KD ($\lambda_f$=1) | 72.26 | 79.66 | [28] w/ KD ($\lambda_f$=3) | 72.85 | 80.06 |
| **Fed-MAS** ($\lambda_f$=1) | 72.94 | 82.73 | **Fed-MAS** ($\lambda_f$=3) | **74.12** | **83.28** |



**Fig. 5.** Higher Value of $\lambda_f$ ($\lambda_f = 7$) causes task deviation. $\lambda = 3$ show faster convergence (Acc.), and make $L_{sup}/L_f$ ratio consistent on a toy dataset (CIFAR-100 non-iid).

---

**Algorithm 1** Pseudocode for Fed-MAS.

---

1: **Notations** total number of clients (C), server (S), total communication rounds (R), local epochs (E), learning rate ($\eta$), and a set of client's data sliced into batches of size B ($\mathcal{B}$).

2: **<u>ServerExecution:</u>**

3: Init $\theta_{glob}^1$

4: **for** *each round* $r = 1, ..., R$ **do**

5:     **for** *client* $c \in C$ *in parallel* **do**

6:         $\theta_c, RF_c \leftarrow$**LocalUpdate**$(\theta_{glob}^r)$;

7:     $\theta_{glob}^{r+1} \leftarrow$**DLMA**$(RF_c, \theta_c', c = 1 \text{ to C})$; // Eq. (5)

8: **Return** $\theta_{glob}^R$

9: **<u>LocalUpdate</u>** $(\theta_{glob})$:

10: Init $\hat{w}_k = 0$;

11: Init $w_k = \mathcal{L}_k^{\theta_{glob}}$;

12: **for** *each local epoch* $e = 1, ..., E$ **do**

13:     **for** *each batch* $b \in \mathcal{B}$ **do**

14:         $\mathcal{L}_{total} = \mathcal{L}_{sup} + \lambda_f \mathcal{L}_f$; // Eq. (4)

15:         $\theta' \leftarrow \theta' - \eta \nabla \mathcal{L}_{total}$;

16:         $\hat{w}_k \leftarrow \hat{w}_k + \mathcal{L}_f(b_k)$; // running distillation loss mean for each class k

17: $RF = \sum_{k=1}^{K} w_k \hat{w}_k$;  // RF $\uparrow \approx$ divergence $\theta_{glob}, f_\xi \uparrow$

18: **Return** $\theta', RF$

---