# Multi-Task Distributed Learning using Vision Transformer with Random Patch Permutation

Sangjoon Park, and Jong Chul Ye, *Fellow, IEEE*

*Abstract*—The widespread application of artificial intelligence in health research is currently hampered by limitations in data availability. Distributed learning methods such as federated learning (FL) and shared learning (SL) are introduced to solve this problem as well as data management and ownership issues with their different strengths and weaknesses. The recent proposal of federated split task-agnostic (FᴇSTA) learning tries to reconcile the distinct merits of FL and SL by enabling the multi-task collaboration between participants through Vision Transformer (ViT) architecture, but they suffer from higher communication overhead. To address this, here we present a multi-task distributed learning using ViT with random patch permutation. Instead of using a CNN based head as in FESTA, p-FESTA adopts a randomly permuting simple patch embedder, improving the multi-task learning performance without sacrificing privacy. Experimental results confirm that the proposed method significantly enhances the benefit of multi-task collaboration, communication efficiency, and privacy preservation, shedding light on practical multi-task distributed learning in the field of medical imaging.

*Index Terms*—Federated learning, Split learning, Multi-task learning, Vision Transformer, Privacy preservation

## I. INTRODUCTION

ARTIFICIAL intelligence (AI) has been gaining unprecedented popularity thanks to its potential to revolutionize various fields of data science. Specifically, the deep neural network has attained expert-level performances in the various applications of medical imaging [1], [2].

To enable the AI models to offer precise decision support with robustness, an enormous amount of data are indispensable. However, data collected from volunteer participation of only a few institutions cannot fully meet the amount to guarantee robust performances. Even for the large public datasets, it may inevitably include unquantifiable biases stemming from the limited geographic regions and patient demographics such as ethnicities and races, resulting in performance instability in real-world applications. Especially for the newly emerging disease like Coronavirus disease 19 (COVID-19), this limitation can be exacerbated as it is hard to build a large, well-curated dataset with sufficient diversity promptly.

Therefore, the ability to collaborate between multiple institutions is critical for the successful application of AI in medical imaging, but the rigorous regulations and the ethical restrictions for sharing patient data is an another obstacle to multi-institutional collaborative work. Several formal regulations

and guidelines, such as the United States Health Insurance Portability and Accountability Act (HIPAA) [3] and the European General Data Protection Regulation (GDPR) [4], state the strict regulations regarding the storage and sharing of patient data.

Accordingly, distributed learning methods, which perform learning tasks at edge devices in a distributed fashion, can be effectively utilized in healthcare research. Specifically, distributed learning was introduced to enable the model training with data that reside on the source devices without sharing. Federated learning (FL) is one of these methods that enables distributed clients to collaboratively learn a shared model without sharing their training data [5]. However, it still holds several limitations in that it is heavily dependent on the client-side computation resources for parallel computation and not completely free from privacy concerns with gradient inversion attack [6], [7]. Another distributed learning method, split learning (SL) [8], which splits the network into parts between clients and the server, is a promising method that puts low computational loads at the edge devices; however, it has the disadvantage of high communication overhead between the clients and server [9], and also has limitation in privacy preservation as the private data can be recovered by the malicious attack with feature hijacking and model inversion [10]. In addition, SL show significantly slower convergence compared with FL and shows suboptimal performance under significantly skewed data distribution between clients [11].

Inspired by the modular decomposition structure of Vision Transformer (ViT), a novel distributed learning method dubbed *Federated Split Task-Agnostic learning* (FᴇSTA) was recently proposed for distributed multi-task collaboration using ViT architecture [12]. The FᴇSTA framework, equipped with the shared task-agnostic ViT body on the server-side and multiple task-specific convolutional neural network (CNN) heads and tails on the clients-side, was able to balance the merit of FL and SL, thereby improving the performances of individual tasks under distributed multi-task collaboration setting at a level even better than the single-task expert model trained in a data-centralized manner.

Nevertheless, there remain several critical limitations with the FᴇSTA framework. First, the communication overhead is higher than that of SL and FL, as the model should continuously share features and gradients as well as head and tail parts of the network, which may impose difficulties in practical implementation. Second, we found that the large size head and tail parts in the original FᴇSTA tends to reduce the role of the shared body, resulting in a small improvement compared to the single task learning despite the ViT's the potential for

multi-task learning (MTL). Finally, the FESTA framework was not free from the privacy issue, as the features transmitted to the server body can be hijacked and reverted to the original data by the outside malicious attackers or "honest-but-curious" server in the same manner in SL.

To alleviate these drawbacks, here we introduce $p$-FESTA framework, a *Federated Split Task-Agnostic learning with permutating pure ViT*, which empowers communication efficient MTL with privacy-preservation. Although the overall composition of $p$-FESTA is similar to that of FESTA, instead of using a CNN based head, $p$-FESTA adopts a simple and task non-specific patch embedder like a vanilla ViT, enforcing the self-attention within the transformer architecture to improve the MTL performance. For privacy preservation, we introduce a *Permutation module* that randomly shuffles the order of all patch features ahead of sending them to the server, to prevent either outside attacker or "honest but curious" server from reverting features into original data containing privacy.

The new architectural change gives the $p$-FESTA several unique advantages. First, the communication overhead was reduced significantly by saving features to be used throughout the entire learning process. Furthermore, the benefit of MTL is enhanced by enforcing the head to play a small role and the multi-task body to do heavy lifting. In addition, data privacy is also enhanced with a simple but effective permutation module using the intrinsic property of ViT.

## II. RELATED WORKS

### A. Vision Transformer (ViT)

ViT [13], a recently introduced deep learning model equipped with an exquisite attention mechanism inspired by its successful application in natural language processing (NLP), has demonstrated impressive performances across many vision tasks. The multi-head self-attention in ViT can flexibly attend to a sequence of patches of the image to encode the cue, enabling the model to be robust to nuisances like occlusion, spatial permutation, and adversarial perturbation, and thereby having the model to be more shape-biased like human than CNN-based model [14].

In addition, the modular design of the ViT is straightforward, implying that the components can be easily decomposed into parts: head to project the images patches into embeddings, transformer body to encode the embeddings, and tail to yield task-specific output. This easily decomposable design offers the possibility in the application for MTL. Recall that the motivation of MTL is originated from attempts to mitigate the data insufficiency problem where the numbers of data for individual tasks are limited. MTL can offer the advantage of improving data efficiency, reducing overfitting through shared representation, and faster convergence by leveraging auxiliary knowledge.

Specifically, MTL with transformer-based models has emerged as a popular approach to improve the performances of the closely related task in NLP [15], [16]. In this approach, a shared transformer learns several related tasks simultaneously, like sentence classification and word prediction, and the tasks-specific module yields the outcome for each task. As shown in previous literature [16], the model trained with MTL strategy generally shows improved performances in a wide range of tasks. Even though not well been studied as in language, the decomposable design of ViT has unleashed the application of MTL to visual transformer models. In an early approach [17], the ViT was divided into the task-specific head, tail, and shared transformer structures across the tasks, and it was possible to attain a similar generalization performance with fewer training steps, by sharing the transformer model among the related tasks.

### B. Federated Split Task-Agnostic (FESTA) Learning

The main motivation of existing FESTA framework as described in Fig. 1A and Fig. 2B was to devise a framework to maximally exploit the distinct strengths of FL and SL methods and to improve the performances of individual tasks with collaboration between clients performing various tasks.

Let $\mathcal{C} = \bigcup_{k=1}^{K} C_k$ be a group of client sets with different tasks, where $K$ denotes the number of tasks and the client set $C_k$ includes one or more clients having different data sources for the $k$-th task, i.e. $C_k = \{c_1^k, c_2^k, \ldots, c_{N_k}^k : N_k \geq 1\}$. Clients in each client set for the $k$-th task has its own task-specific model architecture for head $\mathcal{H}_c$ and a tail $\mathcal{T}_c$, while the server-side transformer body $\mathcal{B}$ is shared.

For training, the server and each client initialize the weights of each sub-network with random initialization or from the pre-trained parameters. For learning round $i = 1, 2, \ldots R$, individual clients do the forward pass on their task-specific head $\mathcal{H}_c$ using the local training data $\{(x_c^{(i)}, y_c^{(i)})\}_{i=1}^{N_c}$, and send the intermediate feature $h_c^{(i)}$ to the server: $h_c^{(i)} = \mathcal{H}_c(x_c^{(i)})$. The transformer body $\mathcal{B}$, then receives the intermediate features from all clients and gets features $b_c^{(i)}$ in parallel with the forward pass, to send them back to each client $c$: $b_c^{(i)} = \mathcal{B}(h_c^{(i)})$. With the features $b_c^{(i)}$, the task-specific tail in client yield the output $\hat{y}_c^{(i)} = \mathcal{T}_c(b_c^{(i)})$, and forward pass finishes. Back-propagation is performed exactly the opposite way, in order of tail, body, and head. First, loss is calculated in tail as: $\ell_c(y_c^{(i)}, \mathcal{T}_c(\mathcal{B}(\mathcal{H}_c(x_c^{(i)}))))$, where $\ell_c(y, \bar{y})$ denotes the $c$ task-specific loss between the target $y$ and the estimate $\bar{y}$. Then, the gradients are pass from tail, body to head in reverse order to forward-propagation, using the chain-rule.

For multi-task body update, the optimization is performed by fixing the head and tails. For the task-specific head and tail updates, the optimization problem is solved by fixing the Transformer body. In addition, per every "UnifyingRounds", the server aggregates, averages and distributes the head and tail parameters between clients participating in the same task, as in FedAvg [18].

In the previous study, the FESTA along with the MTL was shown to ameliorate the individual performances of the clients in collaboration, while resolving the data governance and ownership issue as well as eliminating the need to transmit the huge weights of the transformer body [12].
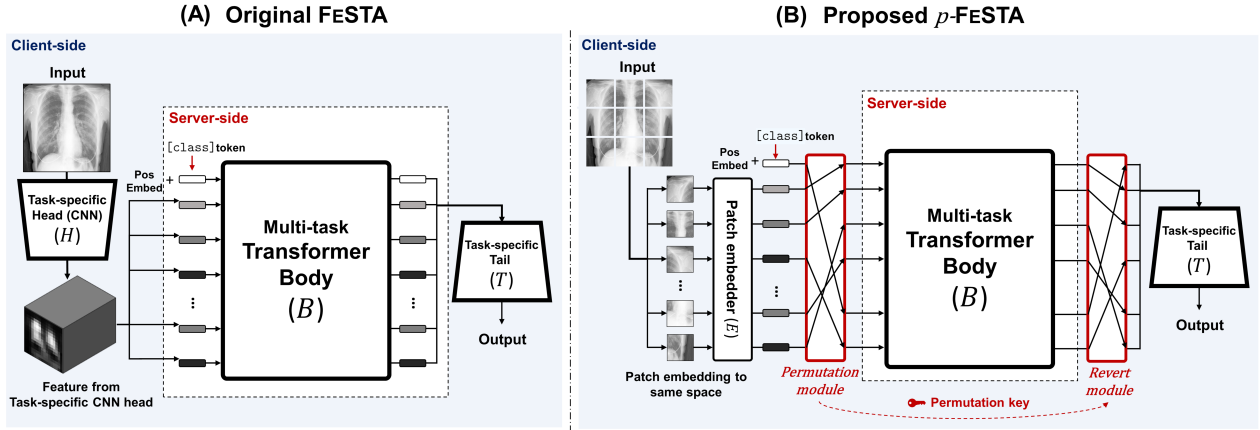
Fig. 1. The model configuration of (A) the original FESTA and (B) the proposed $p$-FESTA frameworks.
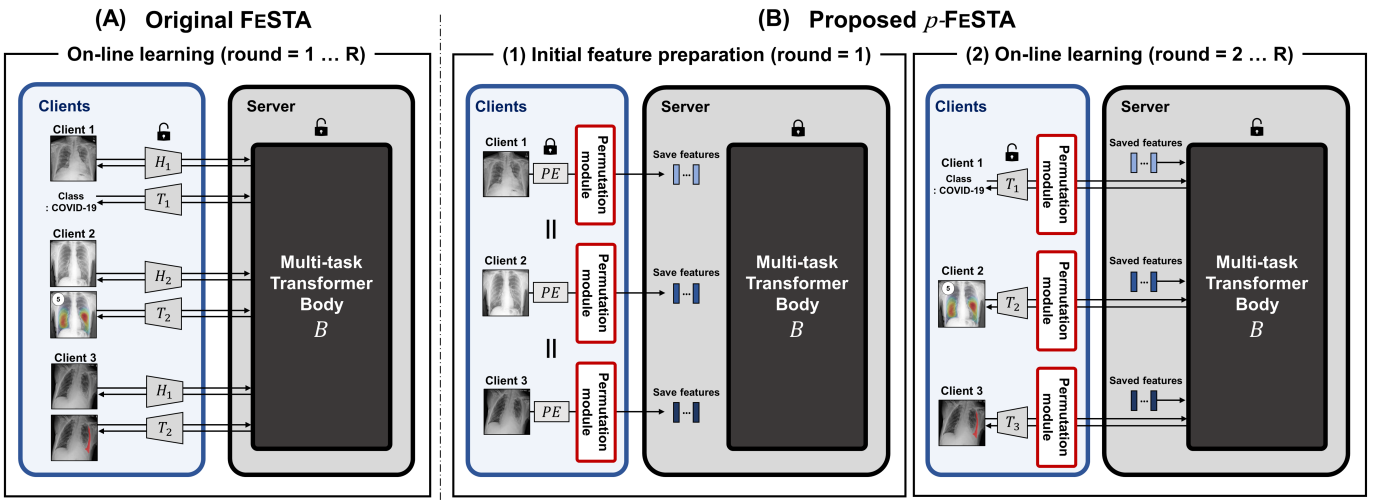


Fig. 2. The learning process of (A) the original FESTA and (B) the proposed $p$-FESTA frameworks.

## III. METHOD

### A. $p$-FESTA

Nonetheless, the FESTA framework still has several drawbacks. First, the communication cost can be higher since the features and gradients should be continuously exchanged between the server and clients like in SL but the head and tail weights should also be aggregated between the clients as in FL. Accordingly, the total communication costs are inevitably higher than SL, and even higher than FL depending on the network size. Second, as shown in our ablation study without the transformer body, the CNN head and tail themselves already have strong representation capacity, which may diminish the role of the transformer body between head and tail. Third, privacy concerns may arise as there is no privacy-preserving method from the model inversion attack on the feature transmitted from client to server.

The proposed $p$-FESTA is a framework devised to mitigate these shortcomings. As shown in Fig. 1B and Fig. 2B, the overall composition of $p$-FESTA is similar to that of FESTA, which decomposes networks into head $\mathcal{H}$, body $\mathcal{B}$ and tail $\mathcal{T}$. However, unlike previous FESTA, we do not use the CNN head tailored for each task. Having the CNN head to

be powerful enough to play a major role in the task hinders the shared transformer from being an important component as there remains little room to improve with this additional module. Instead, we adopted a simple and task non-specific patch embedder like a vanilla ViT, enforcing the self-attention within the transformer architecture to do the heavy lift.

Unfortunately, the use of patch embedding in a vanilla ViT may be prone to outside attackers that attempts to invert the patch embedder to obtain the original images. To address this, here we propose a novel *permutation module* as depicted in Fig. 1B to prevent either outside attacker or "honest but curious" server from reverting features into original data containing privacy. Specifically, this *Permutation module* randomly shuffles the order of all patch features ahead of sending them to the server, and stores the key to reverse the permutation on the client-side. Then, the transformer body $\mathcal{B}$ in the server does a forward pass with the permutated patch features and sends the encoded features back to the clients. Finally, the client reverses the permutation with the saved key and yields the final output by passing the reverted features to the task-specific tail $\mathcal{T}_k$. The back-propagation is performed in the exact opposite way, where the same *Permutation module* to forward-propagation is
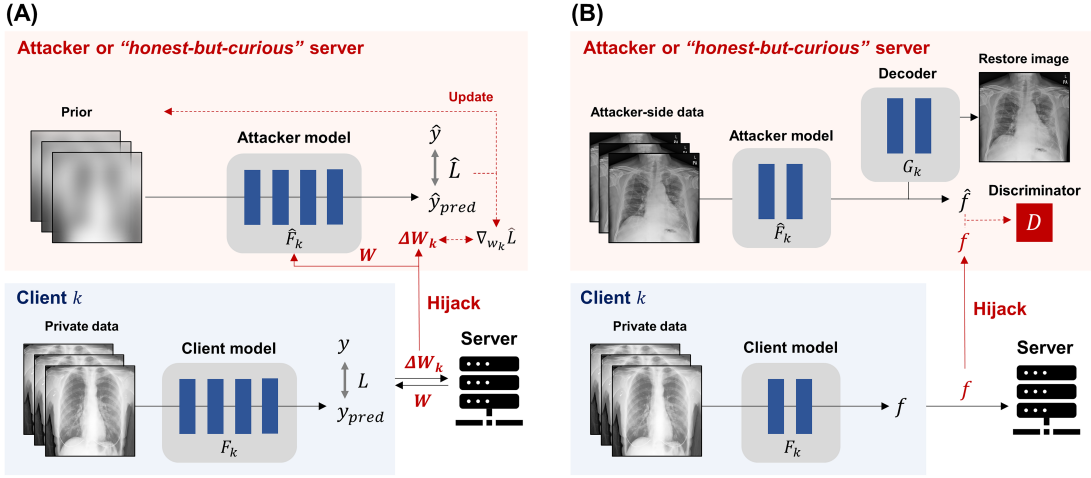
Fig. 3. Possible privacy attacks in (A) federated learning and (B) split learning.
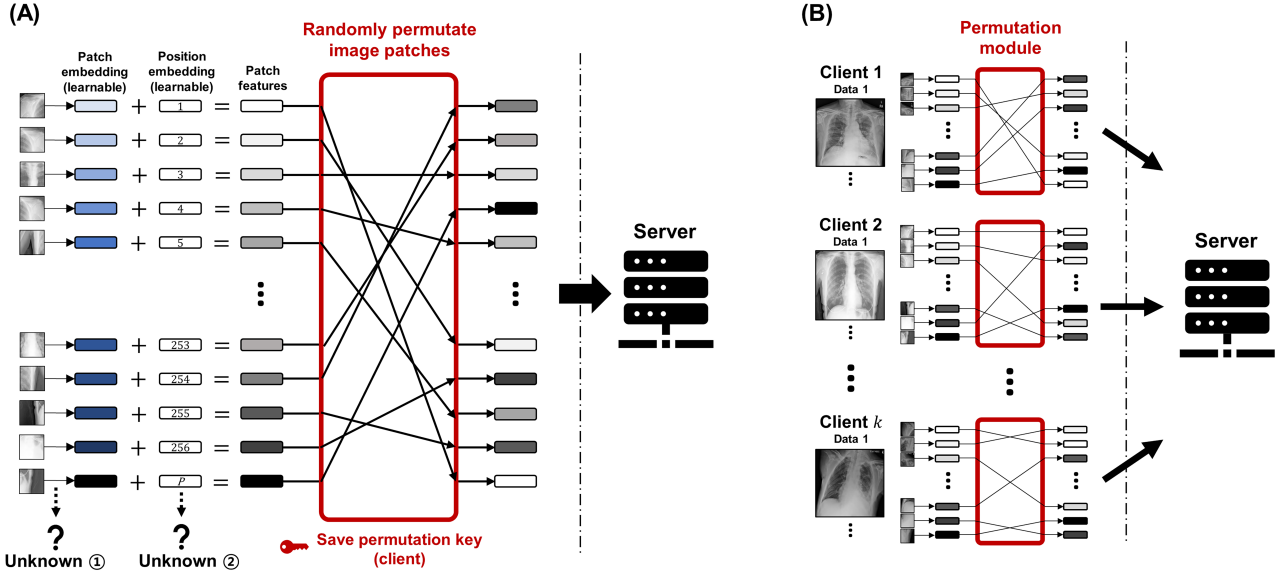


Fig. 4. (A) The *Permutation module* to enhance privacy and (B) the illustration of different permutation patterns for each data of each client.

utilized.

The availability of *Permutation module* attributes to an intriguing property of ViT that all the components composing the transformer body, such as multi-head self-attention, feed-forward network, and layer normalization, is fundamentally "permutation equivariant" [14]. They are processed independently in a patch-based manner and the order of the patch does not affect the outcome, and therefore, the transformer body can be trained without any performance degradation. In addition, as the orders of patches are completely shuffled, it is infeasible for a malicious attacker to successfully revert the original image. How the *Permutation module* can provide privacy protection from the malicious attacker will be described in more detail in the following.

### B. Protecting Privacy with Permutation Module

For FL, privacy is improved by the ephemeral and focused nature of the federated aggregation, averaging, and distribution

of the model updates, assuming that the model updates are considered to be less informative than the original data. However, recent studies have thrown doubt to the false sense of security, showing that private data can be uncovered faithfully only with these local model updates [19]–[21]. In detail, given the access to the global model $W$ and the client's model update $\Delta W$, the attacker can optimize the input image from a prior to produce a gradient that matches the client's model update as illustrated in Fig. 3A. However, this type of attack is infeasible for the proposed $p$-FESTA method, since only the tail part of the entire model is aggregated and distributed by the server to the clients. For instance, for COVID-19 classification, the task-specific tail is a simple linear classifier, with which the original data with privacy cannot be uncovered.

SL protects privacy in a different way. As the name suggests, SL split the entire model into client-side and server-side sub-networks and does not send the models between the server and clients. Instead, the features and gradients are transmitted back

and forth between the server and clients, and it can be the prey of the malicious attacker [10]. As described in Fig. 3B, when clients send the intermediate features $f$ to the server, the attackers may hijack the features, and instead of running the remainder of the SL model, they train three components: encoder $\hat{F}$, decoder $G$, and discriminator $D$ with their own data. The $D$ is trained to discriminate between the hijacked feature $f$ and the feature $\hat{f}$ encoded by $\hat{F}$, which enforces $\hat{f}$ to be in the same feature space to $f$. Simultaneously, $G$ learns to decode $\hat{f}$ into the image with minimal error. Then, well-trained $G$ can also be used to decode the hijacked feature $f$ to original data faithfully.

The feature space hijacking is also possible for our $p$-FESTA. To make the matter worse, the head part of our model is relatively simple and can be easy prey for an attacker. This is why we introduce a novel *Permutation module* to protect privacy as described in Fig. 4A. The *Permutation module* randomly shuffles the order of all patch features. In the implementation, the permutations of each data of each client are all different without any regularity as shown in Fig. 4B, resulting in innumerable patterns for all data. With these random permutations, even if a malicious attacker or server steals patch features to uncover private data, the parameter of position embedding, an unknown learnable variable, cannot be inferred as they have no information about the original order of patches. It is also infeasible to inverse the patch features to image patches since the added position embedding, which is unknown to the attacker, should be subtracted first for inversion. This makes a contradiction that the attacker should already know an "unknown" to infer the other "unknown", making the inversion attack a type of underdetermined problem.

### C. Training Procedure

The learning process of the $p$-FESTA is akin to the original FESTA, but dissimilar in several aspects. Instead of task-specific head $\mathcal{H}_k$ for each task $k$, the task non-specific patch embedder $\mathcal{H}$ prepares the patch embeddings $h_c$ for each client $c$ at the beginning and sends them to the server after passing them through the *Permutation module*. The server then saves the received patch embeddings $h_c$ on its side and uses them throughout the remainder of the learning process in order to update the body $\mathcal{B}$ and tail $\mathcal{T}_k$ parts of the model. Consequently, the overall communication costs can be significantly reduced compared to the original FESTA, as the communications to send the intermediate feature $h_c$ or to update the head $\mathcal{H}$ are no more required.

As can be seen, the head part of the model, the patch embedder, cannot be updated in this configuration. However, fixing the parameters of the patch embedder did not bring about any performance exacerbation thanks to the simple structure to just embed the image patches into the same vector spaces. Having them trainable rather slightly decreased the performances by resulting in the discrepancy in embedding between tasks. The experimental results will be provided in the ablation study of Section III-H. The detailed process of the proposed $p$-FESTA is formally described in Algorithm 1.

---

**Algorithm 1:** Proposed $p$-FESTA algorithm

---

1 **Function** `ServerMain`:
2     Initialize server body weight client head/tail weights
      **for tasks** $k \in \{1, 2, \dots K\}$ **do in parallel**
3         **for clients** $c \in C_k$ **do in parallel**
4             $h_c \leftarrow$ `ClientHead`$(c)$
5             Save patch embedding $h_c$ in server memory

6     **for rounds** $i = 1, 2, \dots R$ **do**
7         **for tasks** $k \in \{1, 2, \dots K\}$ **do in parallel**
8             **for clients** $c \in C_k$ **do in parallel**
9                 **if** $i = 1$ **or** $(i - 1) \in$ UnifyingRounds **then**
10                    $\lfloor$ Set client $w_{\mathcal{T}_c}^{(i)} \leftarrow \bar{w}_{\mathcal{T}, k}$
11                    Load $h_c^{(i)}$ by batch from server memory & $b_c^{(i)} \leftarrow \mathcal{B}(h_c^{(i)})$
12                    $\frac{\partial L_c^{(i)}}{\partial b_c^{(i)}} \leftarrow$ `ClientTail`$(c, b_c^{(i)})$ & Backprop.
13                    $w_{\mathcal{T}_c}^{(i+1)} \leftarrow$ `ClientUpdate`$(c)$

14        Update body
          $$w_{\mathcal{B}}^{(i+1)} \leftarrow w_{\mathcal{B}}^{(i)} - \frac{\eta}{K} \sum_{k=1}^{K} \sum_{c \in C_k} \frac{\partial L_c^{(i)}}{N_k \partial w_{\mathcal{B}}^{(i)}}$$
15        **if** $i \in$ UnifyingRounds **then**
16            **for tasks** $k \in \{1, 2, \dots K\}$ **do**
17                Update $\bar{w}_{\mathcal{T}, k} \leftarrow \frac{1}{N_k} \sum_{c \in C_k} w_{\mathcal{T}_c}^{(i+1)}$

18 **Function** `ClientHead`$(c)$:
19     $x_c \leftarrow$ All data on client $c$
20     $h_c \leftarrow \mathcal{H}(x_c)$
21     Randomly permutate patch embedding $h_c$
22     **return** $h_c$
23 **Function** `ClientTail`$(c, b_c)$:
24     $y_c \leftarrow$ Current batch of label from client $c$
25     $L_c \leftarrow \ell_c(y_c, \mathcal{T}_c(b_c))$ & Backprop.
26     **return** $\frac{\partial L_c}{\partial b_c}$
27 **Function** `ClientUpdate`$(c)$:
28     Backprop. tail, body & $w_{\mathcal{T}_c} \leftarrow w_{\mathcal{T}_c} - \eta \frac{\partial L_c}{\partial w_{\mathcal{T}_c}}$
29     **return** $w_{\mathcal{T}_c}$

---

## IV. RESULTS

### A. Implementation Details

As for the head part, we used the task non-specific patch embedder consisting of the convolution layer with a kernel size of $16 \times 16$ and stride of 16, input channel of 3, and output channel of 768. For the server-side body, the transformer encoder of the ViT-base model, consisting of 12 encoder layers and 12 attention heads, was used. For the tail part, the network architectures specialized to yield the task-specific output were adopted. For the COVID-19 classification task, we used a simpler linear classifier. For severity prediction, the mapping module with five up-sizing convolution layers was adopted as proposed in [22]. For pneumothorax segmentation, the decoder

TABLE I
DATA PARTITIONING FOR COVID-19 CLASSIFICATION

| Class | CNUH (Test) | KNUH (Client #1) | BIMCV (Client #2) |
|---|---|---|---|
| Normal | 417 | 400 | 93 |
| Other infection | 58 | 400 | - |
| COVID-19 | 81 | 293 | 782 |
| Total | 556 | 1093 | 875 |

TABLE II
DATA PARTITIONING FOR COVID-19 SEVERITY PREDICTION

| Severity | CNUH (Test) | YNU (Client #3) | Brixia (Client #4) |
|---|---|---|---|
| 1 | 26 | 63 | 261 |
| 2 | 11 | 59 | 443 |
| 3 | 8 | 25 | 414 |
| 4 | 7 | 35 | 866 |
| 5 | 12 | 18 | 745 |
| 6 | 17 | 86 | 1536 |
| Total | 81 | 286 | 4265 |

TABLE III
DATA PARTITIONING FOR PNEUMOTHORAX SEGMENTATION

| Data | Subset #1 (Test) | Subset #2 (Client #5) | Subset #3 (Client #6) |
|---|---|---|---|
| Total | 1000 | 4840 | 4839 |

(KNUH, client #1) and 875 from public data (BIMCV [26], client #2) were used for training and 556 CXRs from another hospital (CNUH) were used as the external test set in COVID-19 classification task as shown in Table I. Similarly, for the COVID-19 severity prediction task, 286 CXRs from a local hospital (YNU, client #3) and 4,265 from public data (Brixia [27], client #4) were used as the training, and 81 CXRs data of another hospital (CNUH) were used as the external test set as provided in Table II. For pneumothorax segmentation, we used the Society for Imaging Informatics in Medicine and the American College of Radiology (SIIM-ACR) Pneumothorax Segmentation Challenge [28] dataset consisting of 10,679 CXR images. The randomly selected 1,000 CXR images were used as the test set, and the remaining 9679 CXR images were randomly split with a 1:1 ratio (4840 and 4839 CXRs) into two clients (client #5 and client #6) to emulate the participation of two hospitals as in Table III. For practical simulation of collaboration between hospitals, we allocated non-overlapping data sources to each client except for the pneumothorax segmentation task where the exact sources of the data can not be estimated. Overall, six clients participated in the MTL scenario, two clients per task. For this study, Institutional Review Board approvals of each participating hospital were obtained and informed consent was waived.

Considering the sizes and compositions of each client, collaboration for the COVID-19 classification task can be regarded as the collaboration between all clients having small data with a substantial imbalance in data distribution. Likewise, the collaboration for COVID-19 severity can be considered to be the simulation of an imbalance in data size between the participants, one client has scanty data while the other client has relatively sufficient data, in addition to the differences in data composition. Finally, the clients for pneumothorax segmentation emulate the situation in which each participating clients have relatively sufficient and homogeneous data with similar sizes.

When viewed in terms of the relevance between tasks, the COVID-19 classification and severity prediction task can be considered to be highly correlated tasks, while the pneumothorax segmentation task may be regarded as a less relevant task.

part of U-Net [23] was used.

We simulated the distributed MTL between the institutions participating in three different CXR tasks: classification, severity prediction of COVID-19, and pneumothorax segmentation. As in [12], the model was first initialized with pre-trained weights for the CheXpert dataset. We minimized the binary cross-entropy (BCE) losses for each class for the classification task. The severity of COVID-19 was predicted and evaluated in an array-based manner as suggested in [24]. Specifically, BCE losses for each six location arrays of the lung were used for the optimization in severity prediction. Finally, for the pneumothorax segmentation, we minimized the binary cross-entropy loss combined with dice and focal losses. The SGD optimizer was used for the classification and severity prediction tasks, while the Adam optimizer was utilized for the segmentation task, with a learning rate of 0.0001 and a warm-up constant learning rate scheduler for all tasks. The batch size was 4 per client, and the warm-up step was 500. The total training round was 6,000 for all clients, and the tail weights are averaged every 100 local iterations. To adjust the scale of gradients, the 1:2:10 gradient scaling was applied for classification, severity prediction, and segmentation, respectively.

The FL, SL, FESTA and $p$-FESTA was simulated on the modification of Flower (licensed under an Apache-2.0 license) [25] framework. All experiments were performed with Python 3.8 and Pytorch 1.8 on Nvidia RTX 3090, 2080 Ti.

### B. Practical Simulation for Multi-task Collaboration

One of the paramount motivations for FL in medical imaging is to make a robust model leveraging the dispersed and small-sized datasets from multiple institutions while avoiding data governance. Therefore, we assume the FL scenario in which the data of several clients are scanty.

For COVID-19 classification and severity prediction, we used both publicly available datasets and private data collected from local institutions. Overall, 1093 CXR from a local hospital

### C. Performance Metrics

To evaluate the classification performance, the area under the receiver operating characteristic curve (AUC) was used. For the severity prediction task, the mean squared error (MSE) of prediction was used as in the previous work [22]. To evaluate the segmentation accuracy, the Dice coefficient was calculated to measure the intersection of the segmentation results and ground truth annotations. All experiments were performed

TABLE IV
PERFORMANCE COMPARISON WITH OTHER DISTRIBUTED LEARNING METHODS

| Methods | Classification AUC | | | | Severity MSE | Segmentation Dice |
|---|---|---|---|---|---|---|
| | Average | Normal | Others | COVID-19 | | |
| Data centralized | 0.671 (0.051) | 0.735 (0.071) | 0.777 (0.045) | 0.500 (0.051) | 1.592 (0.081) | 0.793 (0.005) |
| Federated learning | 0.601 (0.036) | 0.597 (0.146) | 0.483 (0.068) | 0.722 (0.023) | 2.159 (0.188) | 0.789 (0.001) |
| Split learning | 0.546 (0.024) | 0.522 (0.067) | 0.534 (0.050) | 0.583 (0.013) | 2.546 (0.414) | 0.790 (0.000) |
| FeSTA (STL) | 0.718 (0.047) | 0.680 (0.088) | 0.677 (0.032) | 0.795 (0.036) | **1.318 (0.125)** | 0.801 (0.011) |
| p-FeSTA (STL) | 0.696 (0.022) | 0.739 (0.093) | 0.557 (0.118) | 0.790 (0.045) | 1.717 (0.148) | 0.803 (0.004) |
| FeSTA (MTL) | 0.780 (0.019) | 0.785 (0.009) | 0.793 (0.100) | 0.761 (0.034) | 1.416 (0.048) | 0.796 (0.013) |
| **p-FeSTA (MTL)** | **0.884 (0.008)** | **0.906 (0.004)** | **0.890 (0.011)** | **0.857 (0.014)** | 1.361 (0.057) | **0.808 (0.003)** |

Values are presented in mean (standard deviation) of three repeats with different seed.

repeatedly with three different seeds to exclude the coincidence of getting over- or underestimated results.

### D. Comparison Results

Table IV shows a comparison of the proposed $p$-FeSTA with data centralized training, other distributed learning, and original FeSTA methods. For a fair comparison, all other methods underwent the same pre-training step as the proposed method. The same model architectures were used for all other distributed learning methods except for the original FeSTA methods, in which the task-specific CNN head is a key part of the method. For original FeSTA methods, DenseNet-121 equipped with PCAM operation [29] tailored for CXR classification were used as the head instead of the simple patch embedder as in our previous work [22].

The single-task models trained with either $p$-FeSTA showed a similar order of magnitude in performances of three tasks compared with data centralized training, surpassing those of FL and SL. The improvements with the $p$-FeSTA over FL and SL were noticeable, especially for the classification and severity prediction tasks, where the data insufficiency and imbalance problems are prominent. Note that the slightly better performance of the single-task model trained with the original FESTA than the $p$-FeSTA for the task of severity quantifiacation, which is even better than data centralized learning, may attribute to the more expressive task-specific CNN head tailored for CXR tasks.

As shown in Table IV, the model obtained with MTL between three tasks using $p$-FeSTA significantly outperforms the single-task counterparts and all other distributed learning methods. Note that the performance gain with the MTL over the single-task model is more prominent in the $p$-FeSTA than the previous FeSTA. Even when compared with the MTL model obtained with FeSTA, $p$-FeSTA showed similar or slightly better performance in severity prediction and segmentation, but substantially outperformed the previous one in the classification task, providing generally superior performances.

The fact that the benefit of MTL is formidable in classification and severity prediction is intriguing, as they are the tasks in which scanty data with skewed distribution are problematic. On the contrary, the performance improvement was modest for pneumothorax segmentation where each participating clients have a relatively large number of data with even distribution. Moreover, the close relevance between COVID-19 classification

and severity prediction might have further enhanced the benefit of MTL to those tasks, compared with the relatively less relevant task of pneumothorax segmentation.

### E. Communication Costs between Server and Clients

In this section, we provide the estimated communication costs between the server and clients. Given the number of data as $D$, the batch size as $B$, the rounds between aggregation and distribution by the server as $n$, and the transmission of features, gradients, and the head, body, tail parameters as $F$, $G$, and $P_h$, $P_b$, $P_t$, the communication costs $T$ of each distributed learning strategies for a total of $R$ rounds between the server and one client can be formulated as follow:

$$T_{\text{FL}} = \frac{2R}{n}(P_h + P_b + P_t), \quad (1)$$

$$T_{\text{SL}} = 2BR(F + G), \quad (2)$$

$$T_{\text{FESTA}} = 2BR(F + G) + \frac{2R}{n}(P_h + P_t), \quad (3)$$

$$T_{p\text{-FESTA}} = DF + BR(F + G) + \frac{2RP_t}{n}, \quad (4)$$

where the constant 2 is multiplied to take account of the both-way transmissions between server and client. Note that the cost for features and gradients transmissions are not multiplied by 2 in $p$-FeSTA to reflect no transmission of features and gradients to the head during the learning process.

TABLE V
COMMUNICATION COSTS OF THE DISTRIBUTED LEARNING METHODS

| | Total | Features/gradients | Parameters |
|---|---|---|---|
| **Classification** | | | |
| FL | 10456.152M | - | 10456.152M |
| SL | 9474.048M | 9474.048M | - |
| FESTA | 11390.423M | 9474.048M | 1916.375M |
| $p$-FESTA | 4880.648M | 4844.890M | 35.758M |
| **Severity** | | | |
| FL | 11090.794M | - | 11090.794M |
| SL | 9474.048M | 9474.048M | - |
| FESTA | 12025.065M | 9474.048M | 2551.017M |
| $p$-FESTA | 5435.649M | 4765.249M | 670.401M |
| **Segmentation** | | | |
| FL | 11160.985M | - | 11160.985M |
| SL | 9474.048M | 9474.048M | - |
| FESTA | 12095.256M | 9474.048M | 2621.208M |
| $p$-FESTA | 5899.113M | 5158.520M | 740.592M |

TABLE VI
ABLATION STUDIES FOR THE PROPOSED *p*-FESTA

| Methods | Classification AUC | | | | Severity MSE | Segmentation Dice |
|---|---|---|---|---|---|---|
| | Average | Normal | Others | COVID-19 | | |
| Proposed | 0.884 (0.008) | 0.906 (0.004) | 0.890 (0.011) | 0.857 (0.014) | **1.361 (0.057)** | 0.808 (0.003) |
| w learnable head | 0.890 (0.001) | 0.909 (0.014) | 0.895 (0.005) | **0.866 (0.013)** | 1.545 (0.386) | 0.789 (0.000) |
| w/o permutation | **0.890 (0.010)** | **0.909 (0.002)** | **0.904 (0.023)** | 0.858 (0.008) | 1.461 (0.064) | **0.809 (0.002)** |
| w/o position embedding | 0.827 (0.028) | 0.831 (0.035) | 0.786 (0.049) | 0.862 (0.007) | 1.942 (0.112) | 0.798 (0.004) |

Values are presented in mean (standard deviation) of three repeats with different seed.

Numerically, the communication cost for each distributed learning method in our experimental setting can be calculated as in Table V. As one of the critical drawbacks, the communication cost of FESTA is larger than SL and even higher than FL. On the contrary, the proposed *p*-FESTA substantially lessens the communication burden by saving the head features at the beginning and using them throughout the entire learning process on the server-side. In our experimental setting, the total communication overhead of the proposed *p*-FESTA is less than half of the previous FESTA, and also significantly lower than SL as well as FL.

### F. Ablation Studies

Results of ablation studies are suggested in Table VI.

*1) Fixed Head:* We first performed an ablation to verify whether fixing the head parametersand the patch embedder does not harm the performance. Compared with the proposed method with a fixed head, the same model trained using the learnable head showed similar or even slightly worse performance for severity prediction and segmentation, which may attribute to the overfitting to training data. Therefore, we concluded that the patch embedder can be fixed during all the learning rounds without concerns of performance degradation.

*2) Permutation Module:* We next ablated the *Permutation module* to verify whether our method is indeed "permutation equivariant". For this proposition to be true, the performance should be the same regardless of the presence of *Permutation module*. As expected, the performances with and without the *Permutation module* were in the same order of magnitude, with the differences falling within standard deviations, proving that the permutation does not affect the performance of the transformer model.

*3) Position Embedding:* In the proposed method, the position embedding takes two roles, first provides position information to yield the final output in the tail, and second adds an unknown parameter to prevent an attacker from uncovering patch features into image patches. We performed the ablation study to confirm that the position embedding is necessary for optimal performance in addition to privacy preservation. The model trained without the position embedding showed slightly lower performances than that with the position embedding in all tasks, suggesting that the position embedding is indispensable for the best performance as well as privacy preservation.

TABLE VII
COMPARISON OF THE DISTRIBUTED LEARNING METHODS.

| | FL | SL | FESTA | *p*-FESTA |
|---|---|---|---|---|
| Model averaging | O | X | O | O |
| Client-side learning | Parallel | Sequential | Parallel | Parallel |
| Model split | X | O | O | O |
| Communication cost | High | High | High | Low |
| Benefit of MTL | X | X | Small | Large |
| Privacy protection | X | X | X | O |

## V. DISCUSSION

In this work, we introduced a significantly improved federated task-agnostic learning framework with permutating pure ViT, dubbed *p*-FESTA, which resolves the major drawbacks of our previous FESTA framework, leveraging the intrinsic properties of the ViT. The newly proposed *p*-FESTA substantially reduces the communication overhead between server and clients as well as enhances the performance with the authentic multi-task training in the same embedding space, while offering better privacy preservation. Table VII summarizes the comparison between the proposed *p*-FESTA, original FESTA and other distributed learning methods.

One of the most tackling problems of the previous FESTA method was the communication overhead between server and clients since it requires the feature and gradient transmission the same as in SL as well as the server-side aggregation and distribution of the heads and tails parameters for each client. This configuration increases the communication cost to be inevitably larger than SL and even larger than FL according to the network sizes of each model component. To mitigate the problem, we configured the head part to be a simpler structure like a patch embedder so that the pre-trained head can sufficiently show the best performance without the additional training of the head that requires communications back and forth between server and clients. Consequently, the feature from the head could be stored on the server-side at the beginning and used throughout the entire learning process, reducing the overall communication cost to approximately half of other distributed learning methods.

Moreover, having the head part be a common patch embedder provides another advantage of embedding the image features of different tasks to be in the same embedding spaces. This results in the increasing role of following multi-task transformer and facilitates learning better shared representation, compared with our previous method where task-specific CNN heads embed

the features into different embedding spaces for each task and confine the role of transformer resultingly.

Concerns may arise that using a simpler head structure and saving features in the server-side memory could unleash privacy leakage, which is especially important for medical data. To alleviate this concern, we utilized the intriguing properties of the ViT, the "permutation equivariance" of the self-attention mechanism. The patch features were randomly permutated ahead of transmission to the server, making the problem underdetermined to the attacker while not deteriorating the performance of the transformer.

The merit of our method can be maximized in "data-hungry" collaboration. As shown in the experiments, the performance gain was more prominent for the classification task and severity prediction tasks where the data of each client are scanty. Given that one of the important motivations of distributed learning is to enable building a robust model without data centralization with the participation of many clients having limited data, this potential gain in data-hungry collaboration will further incentivize the widespread application.

Nevertheless, our study is not free of limitations. First, even though we simulated the practical collaboration between hospitals on the customized Flower framework, the robustness to the other tackling factor such as straggler-resilience was not verified [6], [30]. Considering that connection instability becomes a common problem in online learning, it should be resolved technically ahead of real-world implementation. Second, we did not consider other types of attacks for distributed learning, such as model poisoning or data poisoning [31]–[33], which is beyond the scope of this work. For defense against these types of malicious attacks, the existing methods [34]–[36] can be utilized along with our framework. Future work might verify the robustness for these types of malicious attacks.

## VI. Conclusion

In this paper, we proposed the novel *p*-FeSTA framework with pure ViT, which elicits the synergy of MTL among heterogeneous tasks as well as reduces the communication overhead significantly compared to the existing FeSTA. In addition, we also enhanced the privacy using the *Permutation module* in a way specific to ViT. We believe that our work is a step toward facilitating distributed learning among the institutions wanting to participate in different tasks, mitigating the major drawbacks of the existing methods.

## Acknowledgement

## References

[1] A. Hosny, C. Parmar, J. Quackenbush, L. H. Schwartz, and H. J. Aerts, "Artificial intelligence in radiology," *Nature Reviews Cancer*, vol. 18, no. 8, pp. 500–510, 2018.

[2] M. K. K. Niazi, A. V. Parwani, and M. N. Gurcan, "Digital pathology and artificial intelligence," *The lancet oncology*, vol. 20, no. 5, pp. e253–e261, 2019.

[3] P. F. Edemekong, P. Annamaraju, and M. J. Haydel, "Health insurance portability and accountability act," 2018.

[4] C. J. Hoofnagle, B. van der Sloot, and F. Z. Borgesius, "The european union general data protection regulation: what it is and what it means," *Information & Communications Technology Law*, vol. 28, no. 1, pp. 65–98, 2019.

[5] J. Konečnỳ, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, "Federated learning: Strategies for improving communication efficiency," *arXiv preprint arXiv:1610.05492*, 2016.

[6] T. Li, A. K. Sahu, A. Talwalkar, and V. Smith, "Federated learning: Challenges, methods, and future directions," *IEEE Signal Processing Magazine*, vol. 37, no. 3, pp. 50–60, 2020.

[7] P. M. Mammen, "Federated learning: Opportunities and challenges," *arXiv preprint arXiv:2101.05428*, 2021.

[8] P. Vepakomma, O. Gupta, T. Swedish, and R. Raskar, "Split learning for health: Distributed deep learning without sharing raw patient data," *arXiv preprint arXiv:1812.00564*, 2018.

[9] C. Thapa, M. A. P. Chamikara, S. Camtepe, and L. Sun, "Splitfed: When federated learning meets split learning," *arXiv preprint arXiv:2004.12088*, 2020.

[10] G. Gawron and P. Stubbings, "Feature space hijacking attacks against differentially private split learning," *arXiv preprint arXiv:2201.04018*, 2022.

[11] Y. Gao, M. Kim, S. Abuadbba, Y. Kim, C. Thapa, K. Kim, S. A. Camtepe, H. Kim, and S. Nepal, "End-to-end evaluation of federated learning and split learning for internet of things," *arXiv preprint arXiv:2003.13376*, 2020.

[12] S. Park, G. Kim, J. Kim, B. Kim, and J. C. Ye, "Federated split vision transformer for covid-19cxr diagnosis using task-agnostic training," *arXiv preprint arXiv:2111.01338*, 2021.

[13] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly *et al.*, "An image is worth 16x16 words: Transformers for image recognition at scale," *arXiv preprint arXiv:2010.11929*, 2020.

[14] M. M. Naseer, K. Ranasinghe, S. H. Khan, M. Hayat, F. Shahbaz Khan, and M.-H. Yang, "Intriguing properties of vision transformers," *Advances in Neural Information Processing Systems*, vol. 34, 2021.

[15] S. Chen, Y. Zhang, and Q. Yang, "Multi-task learning in natural language processing: An overview," *arXiv preprint arXiv:2109.09138*, 2021.

[16] X. Liu, P. He, W. Chen, and J. Gao, "Multi-task deep neural networks for natural language understanding," *arXiv preprint arXiv:1901.11504*, 2019.

[17] H. Chen, Y. Wang, T. Guo, C. Xu, Y. Deng, Z. Liu, S. Ma, C. Xu, C. Xu, and W. Gao, "Pre-trained image processing transformer," *arXiv preprint arXiv:2012.00364*, 2020.

[18] 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.

[19] J. Geiping, H. Bauermeister, H. Dröge, and M. Moeller, "Inverting gradients-how easy is it to break privacy in federated learning?" *Advances in Neural Information Processing Systems*, vol. 33, pp. 16 937–16 947, 2020.

[20] L. Zhu, Z. Liu, and S. Han, "Deep leakage from gradients," *Advances in Neural Information Processing Systems*, vol. 32, 2019.

[21] B. Zhao, K. R. Mopuri, and H. Bilen, "idlg: Improved deep leakage from gradients," *arXiv preprint arXiv:2001.02610*, 2020.

[22] S. Park, G. Kim, Y. Oh, J. B. Seo, S. M. Lee, J. H. Kim, S. Moon, J.-K. Lim, and J. C. Ye, "Multi-task vision transformer using low-level chest x-ray feature corpus for covid-19 diagnosis and severity quantification," *Medical Image Analysis*, vol. 75, p. 102299, 2022.

[23] O. Ronneberger, P. Fischer, and T. Brox, "U-net: Convolutional networks for biomedical image segmentation," in *International Conference on Medical image computing and computer-assisted intervention*. Springer, 2015, pp. 234–241.

[24] D. Toussie, N. Voutsinas, M. Finkelstein, M. A. Cedillo, S. Manna, S. Z. Maron, A. Jacobi, M. Chung, A. Bernheim, C. Eber *et al.*, "Clinical and chest radiography features determine patient outcomes in young and

middle-aged adults with COVID-19," *Radiology*, vol. 297, no. 1, pp. E197–E206, 2020.

[25] D. J. Beutel, T. Topal, A. Mathur, X. Qiu, T. Parcollet, P. P. de Gusmão, and N. D. Lane, "Flower: A friendly federated learning research framework," *arXiv preprint arXiv:2007.14390*, 2020.

[26] M. De La Iglesia Vayá, J. M. Saborit, J. A. Montell, A. Pertusa, A. Bustos, M. Cazorla, J. Galant, X. Barber, D. Orozco-Beltrán, F. García-García *et al.*, "Bimcv COVID-19+: a large annotated dataset of rx and ct images from COVID-19 patients," *arXiv preprint arXiv:2006.01174*, 2020.

[27] A. Signoroni, M. Savardi, S. Benini, N. Adami, R. Leonardi, P. Gibellini, F. Vaccher, M. Ravanelli, A. Borghesi, R. Maroldi *et al.*, "Bs-net: Learning covid-19 pneumonia severity on a large chest x-ray dataset," *Medical Image Analysis*, vol. 71, p. 102046, 2021.

[28] SIIM-ACR, "SIIM-ACR Pneumothorax Segmentation," https://www.kaggle.com/c/siim-acr-pneumothorax-segmentation, 2019.

[29] W. Ye, J. Yao, H. Xue, and Y. Li, "Weakly supervised lesion localization with probabilistic-cam pooling," *arXiv preprint arXiv:2005.14480*, 2020.

[30] A. Reisizadeh, I. Tziotis, H. Hassani, A. Mokhtari, and R. Pedarsani, "Straggler-resilient federated learning: Leveraging the interplay between statistical accuracy and system heterogeneity," *arXiv preprint arXiv:2012.14453*, 2020.

[31] L. Lyu, H. Yu, and Q. Yang, "Threats to federated learning: A survey," *arXiv preprint arXiv:2003.02133*, 2020.

[32] V. Tolpegin, S. Truex, M. E. Gursoy, and L. Liu, "Data poisoning attacks against federated learning systems," in *European Symposium on Research in Computer Security*. Springer, 2020, pp. 480–501.

[33] E. Bagdasaryan, A. Veit, Y. Hua, D. Estrin, and V. Shmatikov, "How to backdoor federated learning," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 2938–2948.

[34] L. Lyu, H. Yu, X. Ma, L. Sun, J. Zhao, Q. Yang, and P. S. Yu, "Privacy and robustness in federated learning: Attacks and defenses," *arXiv preprint arXiv:2012.06337*, 2020.

[35] Y. Zhao, J. Chen, J. Zhang, D. Wu, J. Teng, and S. Yu, "Pdgan: A novel poisoning defense method in federated learning using generative adversarial network," in *International Conference on Algorithms and Architectures for Parallel Processing*. Springer, 2019, pp. 595–609.

[36] S. Li, Y. Cheng, W. Wang, Y. Liu, and T. Chen, "Learning to detect malicious clients for robust federated learning," *arXiv preprint arXiv:2002.00211*, 2020.