

**FEDHCA<sup>2</sup>: Towards Hetero-Client Federated Multi-Task Learning**

001

## Supplementary Material

002

**A. Proof of Theorem 1**

003

**Theorem 1** Given a multi-task model with a shared encoder and task-specific decoders, and a federated learning system consisting of clients with independent encoders and decoders, gradient descent in the shared encoder of MTL is equivalent to averaging parameter aggregation in FL, adding an extra term  $\nabla_{\theta^{(0)}} \langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle$  that maximizes the inner product of gradients  $\hat{g}_i^{(p)}$  and  $\hat{g}_j^{(q)}$  between all pairs of tasks  $i$  and  $j$  in each iteration  $p$  and  $q$ .

004

005

006

007

**Overview** The goal is to establish the relationship between the parameter updates in MTL with a shared encoder and FL with independent encoders. We aim to show that the gradient descent updates in MTL is similar to the parameter aggregation in FL, but MTL inherently reduces gradient conflicts among tasks, an effect not directly achieved by FL.

008

009

010

011

**Analysis** Consider the following scenario, a multi-task model handles  $N$  tasks with a standard multi-decoder architecture consisting of a shared encoder and  $N$  task-specific decoders, trained on a multi-task dataset  $\mathcal{D} = \{(\mathbf{x}_n, \mathbf{y}_n)\}_{n=1}^{|\mathcal{D}|}$ , where  $\mathbf{x}_n$  is the input sample and  $\mathbf{y}_n = \bigcup_{t=1}^N \mathbf{y}_{n,t}$  contains the ground-truth labels for all  $N$  tasks. Meanwhile, a Federated Learning system also handles  $N$  tasks with  $N$  clients, where client  $C_i$  addressing the  $i$ -th task with independent encoder and decoder. The local dataset for client  $C_i$  is separated from the multi-task dataset  $\mathcal{D}_i = \{(\mathbf{x}_n, \mathbf{y}_{n,i})\}_{n=1}^{|\mathcal{D}|}$ . Assume all models are initialized by  $\theta^{(0)}$  and trained for  $M$  iterations before aggregation in FL. Here  $M$  equals the number of iterations in a single epoch multiplies the number of local epochs set by the FL system.

012

013

014

015

016

017

018

Firstly, we could define following annotations to facilitate analysis. In the following proof, we mainly focus on the optimization of encoders in MTL and FL, thus abbreviating the parameters of encoder  $\theta^E$  as  $\theta$  for simplicity. In MTL, the encoder is always shared by all tasks in each mini-batch, we define the loss function for the  $i$ -th task at the  $m$ -th iteration and its corresponding gradient as follows:

019

020

021

022

$$\mathcal{L}_i^{(m)}(\theta^{(m-1)}) = \mathcal{L}_i(\theta^{(m-1)}; \mathcal{B}_i^{(m)}), \quad (1) \quad 023$$

$$g_i^{(m)} = \nabla_{\theta^{(m-1)}} \mathcal{L}_i^{(m)}, \quad (2) \quad 024$$

where  $\theta^{(m-1)}$  is the encoder parameters in the  $(m-1)$ -th iteration,  $\mathcal{B}_i^{(m)}$  is the mini-batch sampled from  $\mathcal{D}$  for task  $i$ .

025

While in FL, each client possesses its specific model, for the same mini-batch  $\mathcal{B}_i^{(m)}$  we have the loss function and gradient:

026

$$\mathcal{L}_i^{(m)}(\theta_i^{(m-1)}) = \mathcal{L}_i(\theta_i^{(m-1)}; \mathcal{B}_i^{(m)}), \quad (3) \quad 027$$

$$g_i^{(m)} = \nabla_{\theta_i^{(m-1)}} \mathcal{L}_i^{(m)}, \quad (4) \quad 028$$

where  $\theta_i^{(m-1)}$  is the encoder parameters of client  $C_i$  in the  $(m-1)$ -th iteration.

029

Since the MTL model  $\theta$  and all FL clients' models  $\theta_i$  are initialized from the same point  $\theta^{(0)}$ , we define

030

$$\hat{g}_i^{(m)} = \nabla_{\theta^{(0)}} \mathcal{L}_i^{(m)}(\theta^{(0)}; \mathcal{B}_i^{(m)}), \quad (5) \quad 031$$

$$\hat{H}_i^{(m)} = \nabla_{\theta^{(0)}}^2 \mathcal{L}_i^{(m)}(\theta^{(0)}; \mathcal{B}_i^{(m)}), \quad (6) \quad 032$$

to represent the derivative and Hessian Matrix of  $\mathcal{L}_i^{(m)}$  of the initial parameters respectively, and are exactly the same for MTL and FL.

033

034

Then we delve deeper into the training procedure of MTL. Since the encoder is shared by all tasks, it is updated by the gradient descents computed from the losses of all tasks, which is formulated as:

035

036

$$\theta^{(m)} = \theta^{(m-1)} - \eta \sum_{i=1}^N g_i^{(m)}, \quad (7) \quad 037$$

038 where  $\eta$  denotes the learning rate.

039 To have an insightful view of gradient descent in MTL, we can perform a Taylor expansion on  $g_i^{(m)}$  assuming  $\eta$  is  
040 sufficiently small, yielding:

$$041 \quad g_i^{(m)} = \nabla_{\theta^{(m-1)}} \mathcal{L}_i^{(m)} \quad (8)$$

$$042 \quad = \nabla_{\theta^{(0)}} \mathcal{L}_i^{(m)} + \nabla_{\theta^{(0)}}^2 \mathcal{L}_i^{(m)} (\theta^{(m-1)} - \theta^{(0)}) + O(\eta^2) \quad (9)$$

$$043 \quad = \hat{g}_i^{(m)} + \hat{H}_i^{(m)} (\theta^{(m-1)} - \theta^{(0)}) + O(\eta^2) \quad (10)$$

$$044 \quad = \hat{g}_i^{(m)} + \hat{H}_i^{(m)} \sum_{k=1}^{m-1} (\theta^{(k)} - \theta^{(k-1)}) + O(\eta^2) \quad (11)$$

$$045 \quad = \hat{g}_i^{(m)} - \eta \hat{H}_i^{(m)} \sum_{k=1}^{m-1} \sum_{j=1}^N g_j^{(k)} + O(\eta^2). \quad (\text{Use Eq. (7)}) \quad (12)$$

046 After  $M$  iterations, we can calculate the overall change of encoder parameters by combining Eq. (7) and Eq. (12):

$$047 \quad \Delta\theta_{MTL} = \theta^{(M)} - \theta^{(0)} \quad (13)$$

$$048 \quad = -\eta \sum_{m=1}^M \sum_{i=1}^N g_i^{(m)} \quad (14)$$

$$049 \quad = -\eta \sum_{m=1}^M \sum_{i=1}^N \left( \hat{g}_i^{(m)} - \eta \hat{H}_i^{(m)} \sum_{k=1}^{m-1} \sum_{j=1}^N g_j^{(k)} \right) + O(\eta^3) \quad (15)$$

$$050 \quad = -\eta \sum_{m=1}^M \sum_{i=1}^N \hat{g}_i^{(m)} + \eta^2 \sum_{m=1}^M \sum_{i=1}^N \left( \hat{H}_i^{(m)} \sum_{k=1}^{m-1} \sum_{j=1}^N g_j^{(k)} \right) + O(\eta^3). \quad (16)$$

051 Similarly, we consider the training procedure of client  $C_i$  in FL, where the encoder  $\theta_i$  is updated independently by the  
052 gradients computed from the loss of its own task, which is formulated as:

$$053 \quad \theta_i^{(m)} = \theta_i^{(m-1)} - \eta g_i^{(m)}. \quad (17)$$

054 Conducting Taylor expansion on  $g_i^{(m)}$  yields:

$$055 \quad g_i^{(m)} = \nabla_{\theta_i^{(m-1)}} \mathcal{L}_i^{(m)} \quad (18)$$

$$056 \quad = \nabla_{\theta_i^{(0)}} \mathcal{L}_i^{(m)} + \nabla_{\theta_i^{(0)}}^2 \mathcal{L}_i^{(m)} (\theta_i^{(m-1)} - \theta_i^{(0)}) + O(\eta^2) \quad (19)$$

$$057 \quad = \hat{g}_i^{(m)} + \hat{H}_i^{(m)} (\theta_i^{(m-1)} - \theta_i^{(0)}) + O(\eta^2) \quad (20)$$

$$058 \quad = \hat{g}_i^{(m)} + \hat{H}_i^{(m)} \sum_{k=1}^{m-1} (\theta_i^{(k)} - \theta_i^{(k-1)}) + O(\eta^2) \quad (21)$$

$$059 \quad = \hat{g}_i^{(m)} - \eta \hat{H}_i^{(m)} \sum_{k=1}^{m-1} g_i^{(k)} + O(\eta^2). \quad (\text{Use Eq. (17)}) \quad (22)$$

060 After  $M$  iterations, we can calculate the overall change of the client encoder parameters by combining Eq. (17) and  
061 Eq. (22):

$$062 \quad \Delta\theta_i = \theta_i^{(M)} - \theta_i^{(0)} \quad (23)$$

$$063 \quad = -\eta \sum_{m=1}^M g_i^{(m)} \quad (24)$$

$$= -\eta \sum_{m=1}^M \left( \hat{g}_i^{(m)} - \eta \hat{H}_i^{(m)} \sum_{k=1}^{m-1} g_i^{(k)} \right) + O(\eta^3) \quad (25) \quad 064$$

$$= -\eta \sum_{m=1}^M \hat{g}_i^{(m)} + \eta^2 \sum_{m=1}^M \left( \hat{H}_i^{(m)} \sum_{k=1}^{m-1} g_i^{(k)} \right) + O(\eta^3). \quad (26) \quad 065$$

FL server typically aggregates client models by performing a weighted sum of client model parameters, such as FedAvg [14]. The aggregation is formulated as: 066  
067

$$\tilde{\theta} = \frac{1}{N} \sum_{i=1}^N \theta_i, \quad (27) \quad 068$$

here we simplify it with identical weights for all clients. 069

Consider its change from the initial weights: 070

$$\Delta \tilde{\theta}_{FL} = \frac{1}{N} \sum_{i=1}^N \Delta \theta_i \quad (28) \quad 071$$

$$= \frac{1}{N} \sum_{i=1}^N \left( -\eta \sum_{m=1}^M \hat{g}_i^{(m)} + \eta^2 \sum_{m=1}^M \left( \hat{H}_i^{(m)} \sum_{k=1}^{m-1} g_i^{(k)} \right) \right) + O(\eta^3) \quad (\text{Use Eq. (26)}) \quad (29) \quad 072$$

$$= \frac{1}{N} \left( -\eta \sum_{m=1}^M \sum_{i=1}^N \hat{g}_i^{(m)} + \eta^2 \sum_{m=1}^M \sum_{i=1}^N \left( \hat{H}_i^{(m)} \sum_{k=1}^{m-1} g_i^{(k)} \right) \right) + O(\eta^3). \quad (30) \quad 073$$

If we regard the optimizer as capable of automatically scaling the learning rate  $\eta$ , we can neglect the coefficient  $1/N$  in Eq. (30). Hence, we could find that the first term  $-\eta \sum_{m=1}^M \sum_{i=1}^N \hat{g}_i^{(m)}$  exists in the parameter update of both MTL (Eq. (16)) and FL (Eq. (30)), showcasing their similarity in optimization of multiple tasks. Furthermore, we can calculate the difference between them: 074  
075  
076  
077

$$\Delta \theta_{MTL} - \Delta \tilde{\theta}_{FL} \approx \eta^2 \sum_{m=1}^M \sum_{i=1}^N \left( \hat{H}_i^{(m)} \sum_{k=1}^{m-1} \sum_{j=1}^N g_j^{(k)} - \hat{H}_i^{(m)} \sum_{k=1}^{m-1} g_i^{(k)} \right) \quad (31) \quad 078$$

$$= \eta^2 \sum_{m=1}^M \sum_{i=1}^N \left( \hat{H}_i^{(m)} \sum_{k=1}^{m-1} \sum_{j=1, j \neq i}^N g_j^{(k)} \right) \quad (32) \quad 079$$

$$= \eta^2 \sum_{m=1}^M \sum_{k=1}^{m-1} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \hat{H}_i^{(m)} g_j^{(k)}. \quad (33) \quad 080$$

From Eq. (12) we know the approximation that  $g_i^{(m)} = \hat{g}_i^{(m)} + O(\eta)$ , thus there is 081

$$\Delta \theta_{MTL} - \Delta \tilde{\theta}_{FL} = \eta^2 \sum_{m=1}^M \sum_{k=1}^{m-1} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \hat{H}_i^{(m)} \hat{g}_j^{(k)}. \quad (34) \quad 082$$

Since the sequence of mini-batches is randomly shuffled in every local epoch, the ordering of iterations within a local epoch is actually randomized. Consequently, we can compute the expectation of Eq. (34) as follows: 083  
084

$$\mathbb{E}[\Delta \theta_{MTL} - \Delta \tilde{\theta}_{FL}] = \eta^2 \mathbb{E} \left[ \sum_{m=1}^M \sum_{k=1}^{m-1} \sum_{i=1}^N \sum_{j=1, j \neq i}^N \hat{H}_i^{(m)} \hat{g}_j^{(k)} \right] \quad (35) \quad 085$$

$$= \frac{\eta^2}{2} \sum_{p=1}^M \sum_{q=1, q \neq p}^M \sum_{i=1}^N \sum_{j=1, j \neq i}^N \hat{H}_i^{(p)} \hat{g}_j^{(q)} \quad (36) \quad 086$$

087

$$= \frac{\eta^2}{2} \sum_{p=1}^M \sum_{q=p+1}^M \sum_{i=1}^N \sum_{j=i+1}^N \left( \hat{H}_i^{(p)} \hat{g}_j^{(q)} + \hat{H}_j^{(p)} \hat{g}_i^{(q)} + \hat{H}_i^{(q)} \hat{g}_j^{(p)} + \hat{H}_j^{(q)} \hat{g}_i^{(p)} \right), \quad (37)$$

088

089

where the restriction  $k < m$  is relaxed under the expectation, note that the indexes are changed in Eq. (37) for re-organizing the terms in summation. From the definition that  $\hat{H}_i = \nabla_{\theta^{(0)}} \hat{g}_i$ , we can take a closer look at the inner term:

090

$$\hat{H}_i^{(p)} \hat{g}_j^{(q)} + \hat{H}_j^{(q)} \hat{g}_i^{(p)} = (\nabla_{\theta^{(0)}} \hat{g}_i^{(p)}) \hat{g}_j^{(q)} + (\nabla_{\theta^{(0)}} \hat{g}_j^{(q)}) \hat{g}_i^{(p)} \quad (38)$$

091

$$= \nabla_{\theta^{(0)}} \langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle, \quad (39)$$

092

093

where  $\langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle$  is the inner product of gradient of task  $i$  and  $j$  at initial weight, computed with mini-batch  $\mathcal{B}_i^{(p)}$  and  $\mathcal{B}_j^{(q)}$ . Similarly, we have

094

$$\hat{H}_i^{(q)} \hat{g}_j^{(p)} + \hat{H}_j^{(p)} \hat{g}_i^{(q)} = (\nabla_{\theta^{(0)}} \hat{g}_i^{(q)}) \hat{g}_j^{(p)} + (\nabla_{\theta^{(0)}} \hat{g}_j^{(p)}) \hat{g}_i^{(q)} \quad (40)$$

095

$$= \nabla_{\theta^{(0)}} \langle \hat{g}_i^{(q)}, \hat{g}_j^{(p)} \rangle. \quad (41)$$

096

Then we can bring Eq. (39) and Eq. (41) back to Eq. (37):

097

$$\mathbb{E}[\Delta\theta_{MTL} - \Delta\tilde{\theta}_{FL}] = \frac{\eta^2}{2} \sum_{p=1}^M \sum_{q=p+1}^M \sum_{i=1}^N \sum_{j=i+1}^N \left( \nabla_{\theta^{(0)}} \langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle + \nabla_{\theta^{(0)}} \langle \hat{g}_i^{(q)}, \hat{g}_j^{(p)} \rangle \right) \quad (42)$$

098

$$= \frac{\eta^2}{2} \sum_{p=1}^M \sum_{q=1, q \neq p}^M \sum_{i=1}^N \sum_{j=i+1}^N \left( \nabla_{\theta^{(0)}} \langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle \right). \quad (43)$$

099

100

This can be viewed as the summation of gradients of the inner products in Eq. (39) between all  $N(N-1)/2$  pairs of tasks and all  $M(M-1)$  pairs of mini-batches.

101

102

To understand the effect of Eq. (43) in the optimization process of MTL, we can have a quick review at the gradient descent algorithm:

103

$$\theta \leftarrow \theta - \eta \nabla_{\theta} \mathcal{L}(\theta), \quad (44)$$

104

which can minimize the loss  $\mathcal{L}(\theta)$ :

105

$$\Delta\theta = -\eta \nabla_{\theta} \mathcal{L}(\theta) \Leftrightarrow \min \mathcal{L}(\theta). \quad (45)$$

106

Then the term  $\nabla_{\theta^{(0)}} \langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle$  in parameter update of MTL is equivalent to maximizing the inner product:

107

$$\nabla_{\theta^{(0)}} \langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle \Leftrightarrow \max \langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle, \quad (46)$$

108

since its coefficient  $\eta^2/2$  in Eq. (43) is positive.

109

110

111

From existing works in multi-object and multi-task optimization [7, 11, 19, 22], we know that the inner product  $\langle \hat{g}_i^{(p)}, \hat{g}_j^{(q)} \rangle$  is a measurement of accordance between gradients  $\hat{g}_i^{(p)}$  and  $\hat{g}_j^{(q)}$ , and maximizing this inner product in parameter update is equal to reducing the conflict of gradients. This complete the proof.

112

113

Therefore, the parameter sharing mechanism in MTL can implicitly help mitigate the conflict of gradients across different tasks, yet the parameter aggregation in FL could be limited in this function.

## 114 B. Algorithm Pseudo-codes

115 We provide detailed illustrations of the proposed Hyper  
116 Conflict-Averse Aggregation scheme from Section 3.3 and  
117 Hyper Cross Attention Aggregation scheme from Section  
118 3.4 in Algorithm 1 and Algorithm 2, respectively.

---

### Algorithm 1 Hyper Conflict-Averse Aggregation

---

**Input:** Previous round encoder parameters  $\theta^{E,(r-1)} = \{\theta_1^{E,(r-1)}, \dots, \theta_N^{E,(r-1)}\}$ , current round encoder updates  $\Delta\theta^{E,(r)} = \{\Delta\theta_1^{E,(r)}, \dots, \Delta\theta_N^{E,(r)}\}$ , a hyper-parameter  $c \in [0, 1)$ , Hyper Aggregation Weights for encoders  $\alpha = \{\alpha_1, \dots, \alpha_N\}$

**Output:** Personalized encoder parameters  $\theta^{E,(r)}$

- 1: Define  $\Delta\bar{\theta}^E = \frac{1}{N} \sum_{i=1}^N \Delta\theta_i^{E,(r)}$ ,  $\phi = c^2 \|\Delta\bar{\theta}^E\|^2$
- 2: Solve for  $w^*$  and  $\lambda^*$

$$\min_w F(w) = U_w^\top \Delta\bar{\theta}^E + \sqrt{\phi} \|U_w\| \quad (47)$$

$$\text{where } U_w = \frac{1}{N} \sum_{i=1}^N w_i \Delta\theta_i^{E,(r)} \quad (48)$$

- 3: Compute aggregated update

$$\tilde{U} = \Delta\bar{\theta}^E + U_{w^*}/\lambda^* = \Delta\bar{\theta}^E + \frac{\sqrt{\phi}}{\|U_w\|} U_w$$

- 4: **for**  $i \in \{1, \dots, N\}$  **do**
- 5:     Compute personalized update with Hyper Aggregation Weights

$$\theta_i^{E,(r)} = \theta_i^{E,(r-1)} + \Delta\theta_i^{E,(r)} + \alpha_i \tilde{U}$$

- 6: **end for**
- 

---

### Algorithm 2 Hyper Cross Attention Aggregation

---

**Input:** Previous round decoder parameters  $\theta^{D,(r-1)} = \{\theta_1^{D,(r-1)}, \dots, \theta_K^{D,(r-1)}\}$ , current round decoder updates  $\Delta\theta^{D,(r)} = \{\Delta\theta_1^{D,(r)}, \dots, \Delta\theta_K^{D,(r)}\}$ , Hyper Aggregation Weights for decoders  $\beta = \{\beta_1, \dots, \beta_K\}$

**Output:** Personalized decoder parameters  $\theta^{(r)}$

- 1: **for** each layer  $l$  in decoder **do**
- 2:      $V_l = [\Delta\theta_{1,l}^{D,(r)}, \dots, \Delta\theta_{K,l}^{D,(r)}]^\top$
- 3:     **for**  $i \in \{1, \dots, K\}$  **do**
- 4:         Compute Cross Attention
- 5:          $\tilde{A}_{i,l} = \text{Softmax}(\Delta\theta_{i,l}^{D,(r)} V_l^\top / \sqrt{d}) V_l$
- 6:         Compute personalized update with Hyper Aggregation Weights

$$\theta_{i,l}^{D,(r)} = \theta_{i,l}^{D,(r-1)} + \Delta\theta_{i,l}^{D,(r)} + \beta_{i,l} \tilde{A}_{i,l}$$

- 6:     **end for**
  - 7: **end for**
- 

## 119 C. Additional Implementation Details

**Data augmentation.** Our methodology follows the established data augmentation procedures in previous studies [4, 12, 21]. To augment the training images, we use random scaling with factors ranging from 0.5 to 2.0, random cropping to the specified input resolutions (which are  $512 \times 512$  for the PASCAL-Context dataset [15] and  $448 \times 576$  for the NYUD-v2 dataset [20], in accordance with Swin Transformer’s requirements), random horizontal flipping, and random color jittering. Surface normal labels are corrected for horizontal flipping, and depth labels are corrected for scaling. Additionally, we perform image normalization during both the training and evaluation phases.

**Loss functions and weights.** In alignment with established practices in the field [4, 12, 21], our approach employs distinct loss functions for each specific task. For semantic segmentation and human parts segmentation, we utilize the cross-entropy loss. For saliency detection we use the balanced cross-entropy loss. We adopt the  $\mathcal{L}_1$  loss for surface normal estimation and depth estimation. For edge detection, we use the weighted binary cross-entropy loss, assigning a weight of 0.95 to positive pixels and 0.05 to negative ones. For the losses in multi-task clients, we employ a weighted sum of the individual losses to ensure task balancing. The respective loss weights for SemSeg, Parts, Normals, Sal, Edge, and Depth tasks are set at 1, 2, 10, 5, 50, and 1.

**Implementation.** Our models are optimized using the AdamW optimizer [10] with an initial learning rate of  $1e-4$  and a weight decay rate of  $1e-4$ . Additionally, we employ a cosine decay learning rate scheduler [9] complemented by a 5-epoch warm-up phase. The Hyper Aggregation Weights are clamped between 0 to 1, and are initially set to 0.01 and updated via the SGD optimizer, which operates at a learning rate of  $1e-2$ , a weight decay rate of  $1e-4$ , and a momentum of 0.9. The process of fine-tuning these hyper-parameters is demonstrated in subsequent sections, specifically in Tab. 4, Tab. 5, and Tab. 6. For the Hyper Conflict-Averse Aggregation hyper-parameter  $c$ , we explore a range of values including  $\{0.2, 0.4, 0.6, 0.8\}$ , with the results of this exploration presented in Tab. 7 and Tab. 8.

To evaluate the performance of our method, we compare with representative works including two traditional FL approaches FedAvg [14] and FedProx [5], four pFL methods FedPer [1], Ditto [6], FedAMP [3], FedBABU [16] and two FMTL methods FedSTA [17] and MaT-FL [2]. These methods are adapted for HC-FMTL setting, with adjustments made to the algorithm implementations to fit this novel setting while striving to retain as much of the original algorithms’ integrity as possible. For both FedAvg and FedProx, we apply their aggregation mechanisms to the encoders and decoders separately. The update strategies of FedAMP and Ditto are similarly extended to both encoders and decoders. In the case of FedPer, FedBABU, FedSTA and MaT-FL, ag-

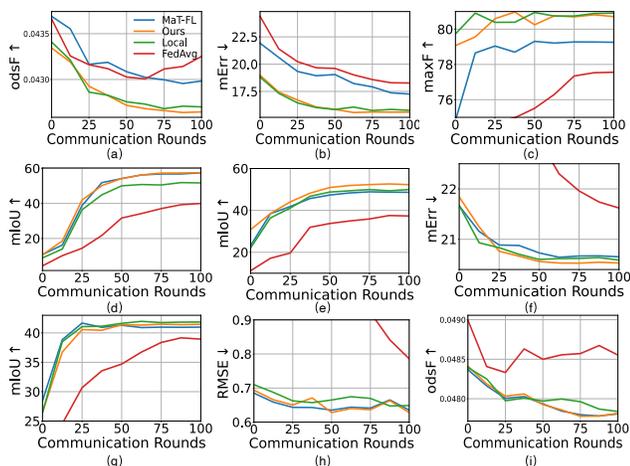


Figure 1. Evaluation results during training, using PASCAL-Context for five single-task clients and NYUD-v2 for one multi-task client. (a) Edge from PASCAL-Context (hereinafter called P) on single-task client. (b) Normals from P on single-task client. (c) Sal from P on single-task client. (d) SemSeg from P on single-task client. (e) Parts from P on single-task client. (f) Normals from NYUD-v2 (hereinafter called N) on multi-task client. (g) SemSeg from N on multi-task client. (h) Depth from N on multi-task client. (i) Edge from N on multi-task client.

gregation is confined to the encoders.

Regarding hyperparameters, we adhere to the default configurations as specified in the original publications or source code of these methods. Specifically, for FedProx, the hyperparameter  $\mu$  is set to 0.01. For Ditto, the regularization parameter  $\lambda$  is set to 0.1. Within FedAMP, the hyperparameters  $\lambda$ ,  $\alpha_k$ , and  $\sigma$  are uniformly assigned a value of 1. For MaT-FL, the number of clusters  $K$  is set to 2.

**Metric evaluation.** To evaluate edge detection results, we use the SEISM package [18] and set the maximum allowed mis-localization of the optimal dataset F-measure (odsF) [13] to 0.0075 and 0.011 for PASCAL-Context and NYUD-v2, respectively.

## D. Additional Experimental Results

**Evaluation results on all tasks.** As the supplementary of Fig. 4 in our paper, Fig. 1 shows that FEDHCA<sup>2</sup> converges faster to a better result on most tasks of PASCAL-Context and NYUD-v2, compared to local training baseline, FedAvg and MaT-FL.

**Impact of the number of clients.** To assess the effectiveness of FEDHCA<sup>2</sup> across varying client counts, we conduct tests by scaling the number of clients per task by factors of 2 and 4, with the datasets evenly split. In our paper, we demonstrate the consistent superior performance and the overall growth trend of our method in Fig. 5. Here, the detailed results of Fig. 5 are illustrated in Tab. 1.

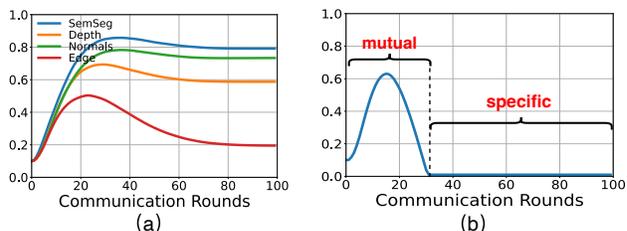


Figure 2. Hyper Aggregation Weights  $\alpha$  for encoders of the client models, using NYUD-v2 for four single-task clients and PASCAL-Context for one multi-task client. (a) Weights of four single-task clients. (b) Weights of the multi-task client which differs in two stages.

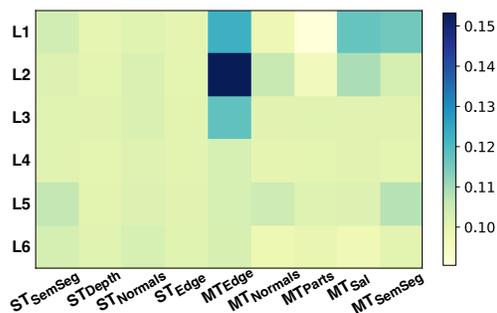


Figure 3. Learned Hyper Aggregation Weights  $\beta$  across decoders for different tasks, spanning layers from L1 to L6, using NYUD-v2 for four single-task clients and PASCAL-Context for one multi-task client.

**Impact of different FMTL scenarios.** To further verify the necessity of introducing our new setting as we do in the paper, we conduct experiments comparing two scenarios: 1) each client handles a single task, and 2) HC-FMTL encompasses both single-task and multi-task clients. These experiments are carried out on the PASCAL-Context, as a supplementary of the NYUD-v2 used and illustrated in Tab. 4 in the paper. As Tab. 2 illustrates, similar results are obtained that while FEDHCA<sup>2</sup> improves upon the local baseline in the single-task client scenario, integrating the multi-task client results in a greater enhancement.

**Impact of different backbones.** We also conduct a series of experiments using Swin-T, Swin-S, and Swin-B [8] as backbones within the HC-FMTL benchmark setting to evaluate the consistent performance of FEDHCA<sup>2</sup>. The results, detailed in Tab. 3, confirm our expectations: as the complexity of the backbone increases from Swin-T to Swin-B, we observe a corresponding improvement in performance across all tasks, further validating the effectiveness of FEDHCA<sup>2</sup>.

**Interaction between tasks.** We investigate the dynamic learning process of Hyper Aggregation Weights for both encoders and decoders, aiming to understand their role in facilitating personalized aggregation for different clients. In Fig. 6 and Fig. 7 in our paper, we show the evolution

223 of weights  $\alpha$  for encoders and the last learned weights  $\beta$   
224 across decoders on the first benchmark setting described in  
225 Sec. 4.1. Here, similar results are obtained on the second  
226 benchmark setting as depicted in Fig. 2 and Fig. 3.

227 **Hyper-parameter tuning.** Additional experiments are car-  
228 ried out to assess the impact of different initialization val-  
229 ues for Hyper Aggregation Weights  $\alpha$  and  $\beta$ . As depicted  
230 in Tab. 4 and Tab. 5, we explore a range from 0.01 to 1,  
231 across two HC-FMTL benchmark settings. Optimal results  
232 on task metrics within PASCAL-Context are achieved when  
233 the initialization value is set to 0.1, which also yield the best  
234 average per-task performance drop across both settings. We  
235 also investigate the effect of the learning rate for Hyper Ag-  
236 gregation Weights, as illustrated in Tab. 6. When the learn-  
237 ing rate varies from 0.005 to 0.1, FEDHCA<sup>2</sup> achieves the  
238 best overall performance at a learning rate of 0.01. More-  
239 over, as indicated in Tab. 7 and Tab. 8, we experiment  
240 with various values of the hyper-parameter  $c$  in the Hyper  
241 Conflict-Averse Aggregation across the two benchmark set-  
242 tings, scaling from 0.2 to 0.8. In both instances, setting  $c$  to  
243 0.4 results in the best results.

## 244 E. Qualitative results

245 To intuitively compare the proposed FEDHCA<sup>2</sup> with exist-  
246 ing methods, we visualize the task predictions of local train-  
247 ing, FedAvg, MaT-FL and our model with an example from  
248 the PASCAL-Context dataset in Fig. 4 and an example from  
249 the NYUD-v2 dataset in Fig. 5. All methods are trained on  
250 the first benchmark scenario. Our model obviously gen-  
251 erates better details with less error, especially in semantic  
252 segmentation and human parts segmentation in PASCAL-  
253 Context and semantic segmentation in NYUD-v2. Our re-  
254 sults also more closely resemble the ground truths in other  
255 tasks.

Table 1. Comparison to representative methods with the number of clients scaling to 2 and 4 times, using PASCAL-Context for five single-task clients and NYUD-v2 for one multi-task client. ‘ $\Delta_m$ ’ is calculated w.r.t. corresponding local baseline of 1C, 2C and 4C.

Scale	Method	PASCAL-Context (ST)					NYUD-v2 (MT)				$\Delta_m\% \uparrow$
		SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normal mErr $\downarrow$	Edge odsF $\uparrow$	SemSeg mIoU $\uparrow$	Depth RMSE $\downarrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	
1C	Local	51.69	49.94	<b>80.91</b>	15.76	71.95	41.86	0.6487	20.59	76.46	0.00
	FedAvg [14]	39.98	37.33	77.56	18.27	69.17	38.94	0.7858	21.62	75.77	-11.76
	FedProx [5]	44.42	38.10	77.26	18.03	69.39	39.19	0.8068	21.52	76.03	-10.68
	FedPer [1]	54.51	46.56	78.85	16.95	71.00	<b>44.02</b>	0.6467	21.19	<b>76.61</b>	-1.11
	FedAMP [3]	55.98	52.05	80.79	15.74	72.02	41.67	0.6428	20.54	76.40	1.47
	MaT-FL [2]	57.45	48.63	79.26	17.26	71.23	40.99	0.6352	20.65	76.59	-0.46
	FEDHCA <sup>2</sup>	<b>57.55</b>	<b>52.30</b>	80.71	<b>15.60</b>	<b>72.08</b>	41.47	<b>0.6281</b>	<b>20.53</b>	76.50	<b>2.18</b>
2C	Local	42.21	47.22	78.64	16.62	70.66	36.77	0.7300	21.96	76.06	0.00
	FedAvg [14]	14.99	29.69	76.68	18.34	69.07	33.01	0.7081	22.15	75.49	-13.95
	FedProx [5]	14.93	30.01	77.03	18.36	68.90	32.47	0.7837	22.14	75.43	-15.20
	FedPer [1]	33.68	42.32	77.72	17.34	70.11	<b>40.99</b>	0.6694	21.93	76.02	-1.89
	FedAMP [3]	<b>45.83</b>	46.78	78.31	16.62	70.62	36.45	0.7093	22.00	75.53	0.91
	MaT-FL [2]	35.23	43.23	78.14	18.11	70.40	39.91	0.6607	<b>20.85</b>	76.16	-1.30
	FEDHCA <sup>2</sup>	45.80	<b>49.52</b>	<b>79.73</b>	<b>15.98</b>	<b>71.61</b>	39.66	<b>0.6603</b>	20.92	<b>76.21</b>	<b>4.70</b>
4C	Local	21.85	34.84	76.59	17.57	69.17	31.49	0.7433	23.55	74.53	0.00
	FedAvg [14]	9.32	17.22	75.76	18.99	68.42	21.89	0.7725	23.37	74.86	-16.82
	FedProx [5]	7.76	13.93	74.38	20.64	68.04	14.78	0.8147	25.35	<b>75.87</b>	-24.04
	FedPer [1]	15.38	31.21	76.00	18.21	<b>76.00</b>	34.27	0.6877	23.35	75.20	-2.96
	FedAMP [3]	<b>25.48</b>	<b>38.38</b>	76.44	17.65	69.22	31.98	0.7908	23.56	74.45	2.36
	MaT-FL [2]	20.10	31.74	76.68	19.06	69.48	33.77	0.7547	22.87	74.74	-1.77
	FEDHCA <sup>2</sup>	22.36	36.54	<b>79.10</b>	<b>16.49</b>	70.72	<b>36.20</b>	<b>0.6510</b>	<b>21.30</b>	75.62	<b>6.36</b>

Table 2. Comparison between different settings. ‘ST+Local’ and ‘ST+Ours’ denote the setting with five single-task clients on PASCAL-Context, trained with local baseline and FEDHCA<sup>2</sup>, respectively. ‘ST+MT+Ours’ denotes the setting in Tab. 1 trained with our framework. ‘ $\Delta_m$ ’ is calculated w.r.t. ‘ST+Local’ baseline.

Method	PASCAL-Context (ST)					$\Delta_m\% \uparrow$
	SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normal mErr $\downarrow$	Edge odsF $\uparrow$	
ST+Local	51.69	49.94	<b>80.91</b>	15.76	71.95	0.00
ST+Ours	<b>57.76</b>	51.38	80.75	15.66	<b>72.08</b>	3.05
ST+MT+Ours	57.55	<b>52.30</b>	80.71	<b>15.60</b>	<b>72.08</b>	<b>3.40</b>

Table 3. Comparison between different backbones on the first benchmark scenario.

Method	PASCAL-Context (ST)					NYUD-v2 (MT)			
	SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normal mErr $\downarrow$	Edge odsF $\uparrow$	SemSeg mIoU $\uparrow$	Depth RMSE $\downarrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$
Swin-T	57.55	52.30	80.71	15.60	72.08	41.47	0.6281	20.53	76.50
Swin-S	62.09	56.42	81.35	15.70	73.20	44.80	0.6247	20.44	76.94
Swin-B	65.30	60.73	81.76	15.53	74.65	47.26	0.6055	19.85	77.58

Table 4. Comparison between different initialized values of Hyper Aggregation Weights  $\alpha$  and  $\beta$  on the first benchmark scenario.

Init	PASCAL-Context (ST)					NYUD-v2 (MT)				$\Delta_m\% \uparrow$
	SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normal mErr $\downarrow$	Edge odsF $\uparrow$	SemSeg mIoU $\uparrow$	Depth RMSE $\downarrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	
0.01	57.03	51.65	80.69	15.64	72.02	41.21	0.6339	20.62	<b>76.54</b>	1.67
0.1	57.55	52.30	80.71	15.60	72.08	<b>41.47</b>	<b>0.6281</b>	<b>20.53</b>	76.50	<b>2.18</b>
0.5	57.15	<b>52.59</b>	80.66	15.55	<b>72.14</b>	40.93	0.6348	20.57	76.45	1.91
1	<b>58.68</b>	51.91	<b>80.79</b>	<b>15.52</b>	72.04	40.96	0.6396	20.61	76.52	2.02

Table 5. Comparison between different initialized values of Hyper Aggregation Weights  $\alpha$  and  $\beta$  on the second benchmark scenario.

Init	NYUD-v2 (ST)				PASCAL-Context (MT)					$\Delta_m\% \uparrow$
	SemSeg mIoU $\uparrow$	Depth RMSE $\downarrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	
0.01	34.84	0.7126	23.24	74.98	<b>66.21</b>	<b>55.17</b>	83.29	14.08	71.95	0.62
0.1	<b>34.95</b>	<b>0.7018</b>	23.19	75.03	65.81	55.01	83.18	14.08	71.97	<b>0.75</b>
0.5	34.85	0.7120	<b>23.15</b>	<b>75.12</b>	65.62	55.02	83.33	14.09	<b>71.99</b>	0.57
1	34.74	0.7113	<b>23.15</b>	74.99	65.45	55.01	<b>83.37</b>	<b>14.02</b>	71.95	0.55

Table 6. Comparison between different learning rates of Hyper Aggregation Weights  $\alpha$  and  $\beta$  on the first benchmark scenario.

lr	PASCAL-Context (ST)					NYUD-v2 (MT)				$\Delta_m\% \uparrow$
	SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normal mErr $\downarrow$	Edge odsF $\uparrow$	SemSeg mIoU $\uparrow$	Depth RMSE $\downarrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	
0.005	<b>58.18</b>	52.06	<b>80.73</b>	15.63	72.08	41.03	0.6337	20.61	<b>76.61</b>	2.00
0.01	57.55	52.30	80.71	<b>15.60</b>	72.08	<b>41.47</b>	<b>0.6281</b>	<b>20.53</b>	76.50	<b>2.18</b>
0.05	57.26	<b>52.46</b>	80.64	15.64	72.09	41.05	0.6338	20.64	76.48	1.84
0.1	57.90	52.13	80.60	15.62	<b>72.10</b>	40.54	0.6332	20.55	76.52	1.85

Table 7. Comparison between different values of hyper-parameter  $c$  in Hyper Conflict-Averse Aggregation on the first benchmark scenario.

$c$	PASCAL-Context (ST)					NYUD-v2 (MT)				$\Delta_m\% \uparrow$
	SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normal mErr $\downarrow$	Edge odsF $\uparrow$	SemSeg mIoU $\uparrow$	Depth RMSE $\downarrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	
0.2	57.12	52.15	80.60	15.63	<b>72.12</b>	<b>41.77</b>	0.6300	20.63	76.53	2.02
0.4	57.55	<b>52.30</b>	<b>80.71</b>	<b>15.60</b>	72.08	41.47	<b>0.6281</b>	<b>20.53</b>	76.50	<b>2.18</b>
0.6	<b>57.60</b>	51.94	80.64	15.62	72.10	40.87	0.6337	20.66	76.46	1.76
0.8	57.18	52.07	80.56	15.65	72.05	40.94	0.6353	20.61	<b>76.54</b>	1.69

Table 8. Comparison between different values of hyper-parameter  $c$  in Hyper Conflict-Averse Aggregation on the second benchmark scenario.

$c$	NYUD-v2 (ST)				PASCAL-Context (MT)					$\Delta_m\% \uparrow$
	SemSeg mIoU $\uparrow$	Depth RMSE $\downarrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	SemSeg mIoU $\uparrow$	Parts mIoU $\uparrow$	Sal maxF $\uparrow$	Normals mErr $\downarrow$	Edge odsF $\uparrow$	
0.2	34.47	0.7135	23.20	75.07	65.63	<b>55.09</b>	<b>83.36</b>	14.08	71.97	0.42
0.4	<b>34.95</b>	<b>0.7018</b>	<b>23.19</b>	75.03	65.81	55.01	83.18	14.08	71.97	<b>0.75</b>
0.6	34.48	0.7135	23.22	74.96	<b>66.02</b>	55.05	83.35	14.09	<b>72.01</b>	0.45
0.8	34.56	0.7105	23.25	<b>75.17</b>	65.83	54.89	83.32	<b>14.06</b>	71.94	0.48

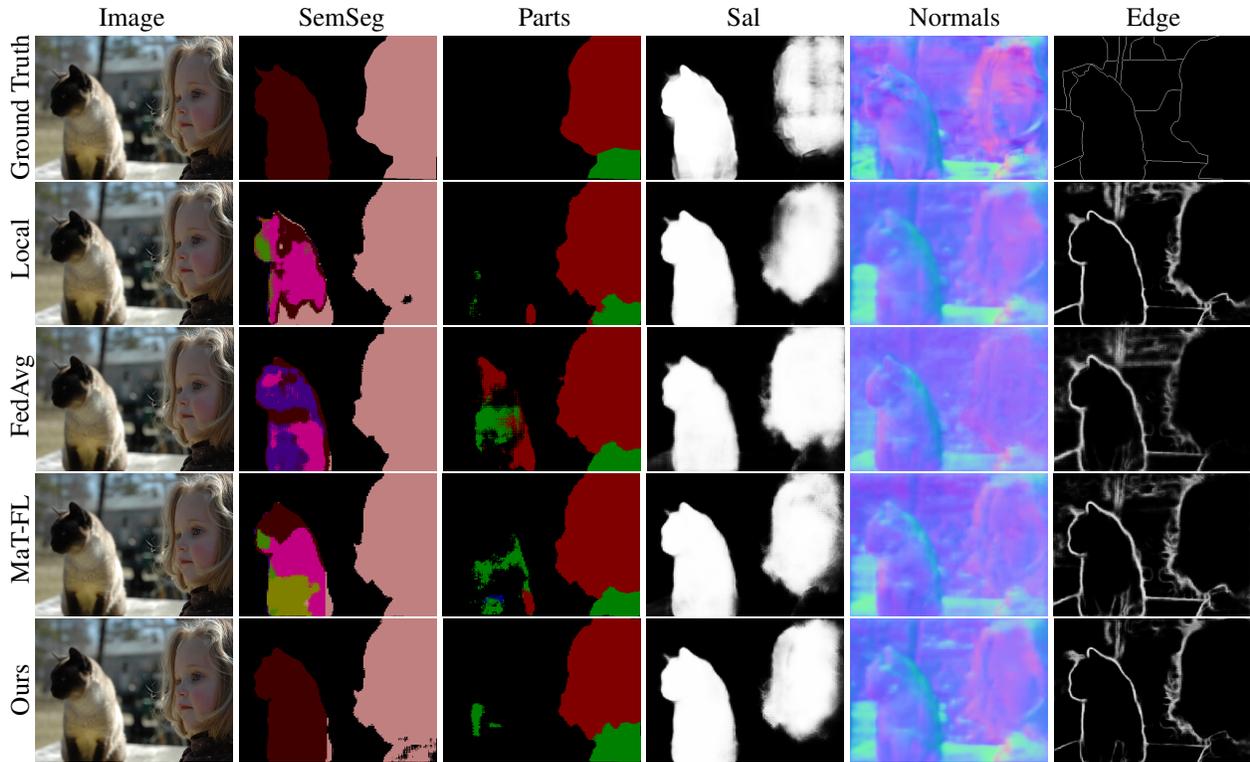


Figure 4. Qualitative results compared with representative methods on PASCAL-Context dataset.

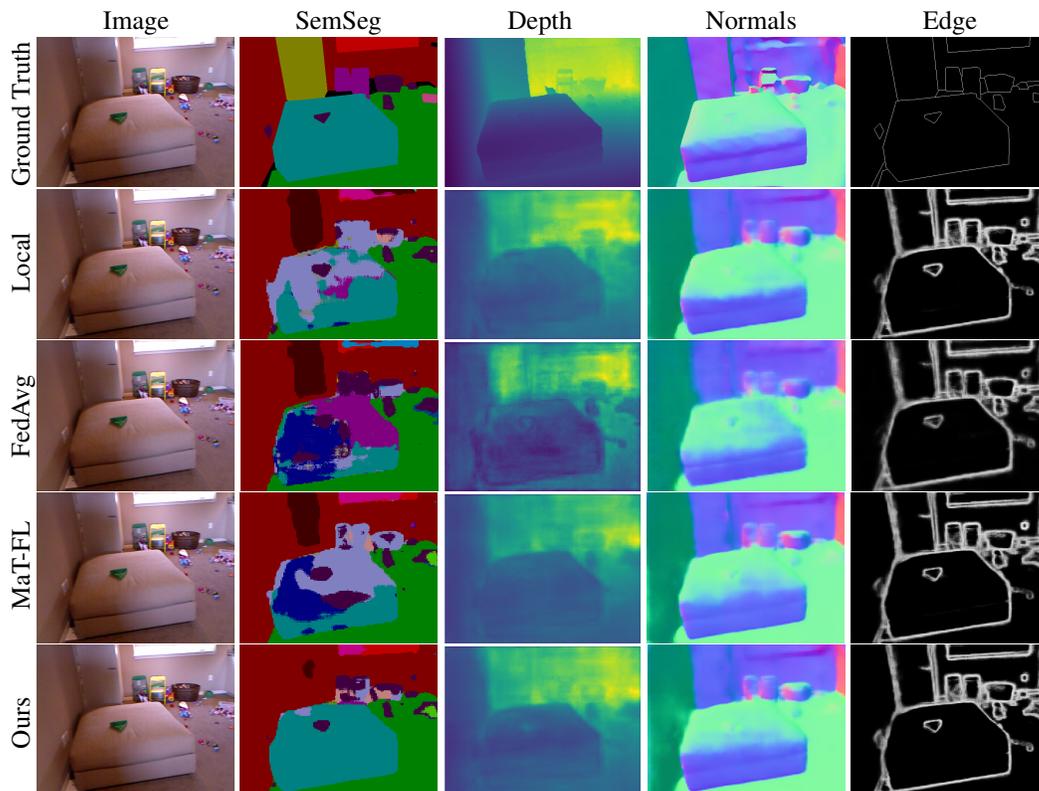


Figure 5. Qualitative results compared with representative methods on NYUD-v2 dataset.

256

**References**257  
258  
259  
260  
261  
262  
263  
264  
265  
266  
267  
268  
269  
270  
271  
272  
273  
274  
275  
276  
277  
278  
279  
280  
281  
282  
283  
284  
285  
286  
287  
288  
289  
290  
291  
292  
293  
294  
295  
296  
297  
298  
299  
300  
301  
302  
303  
304  
305  
306  
307  
308  
309  
310  
311

- [1] Manoj Ghuhana Arivazhagan, Vinay Aggarwal, Aaditya Kumar Singh, and Sunav Choudhary. Federated learning with personalization layers. *arXiv preprint arXiv:1912.00818*, 2019. 5, 8
- [2] Ruisi Cai, Xiaohan Chen, Shiwei Liu, Jayanth Srinivasa, Myungjin Lee, Ramana Kompella, and Zhangyang Wang. Many-task federated learning: A new problem setting and a simple baseline. In *CVPR*, pages 5037–5045, 2023. 5, 8
- [3] Yutao Huang, Lingyang Chu, Zirui Zhou, Lanjun Wang, Jiangchuan Liu, Jian Pei, and Yong Zhang. Personalized cross-silo federated learning on non-iid data. In *AAAI*, pages 7865–7873, 2021. 5, 8
- [4] Menelaos Kanakis, David Bruggemann, Suman Saha, Stamatios Georgoulis, Anton Obukhov, and Luc Van Gool. Reparameterizing convolutions for incremental multi-task learning without task interference. In *ECCV*, pages 689–707, 2020. 5
- [5] Tian Li, Anit Kumar Sahu, Manzil Zaheer, Maziar Sanjabi, Ameet Talwalkar, and Virginia Smith. Federated optimization in heterogeneous networks. In *MLSys*, 2020. 5, 8
- [6] Tian Li, Shengyuan Hu, Ahmad Beirami, and Virginia Smith. Ditto: Fair and robust federated learning through personalization. In *ICML*, pages 6357–6368, 2021. 5
- [7] Bo Liu, Xingchao Liu, Xiaojie Jin, Peter Stone, and Qiang Liu. Conflict-averse gradient descent for multi-task learning. *NeurIPS*, 34:18878–18890, 2021. 4
- [8] Ze Liu, Yutong Lin, Yue Cao, Han Hu, Yixuan Wei, Zheng Zhang, Stephen Lin, and Baining Guo. Swin transformer: Hierarchical vision transformer using shifted windows. In *ICCV*, pages 10012–10022, 2021. 6
- [9] Ilya Loshchilov and Frank Hutter. Sgdr: Stochastic gradient descent with warm restarts. In *ICLR*, 2017. 5
- [10] Ilya Loshchilov and Frank Hutter. Decoupled weight decay regularization. In *ICLR*, 2019. 5
- [11] Linhao Luo, Yumeng Li, Buyu Gao, Shuai Tang, Sinan Wang, Jiancheng Li, Tanchao Zhu, Jiancai Liu, Zhao Li, and Shirui Pan. MAMDR: A model agnostic learning framework for multi-domain recommendation. In *ICDE*, pages 3079–3092, 2023. 4
- [12] Kevis-Kokitsi Maninis, Ilija Radosavovic, and Iasonas Kokkinos. Attentive single-tasking of multiple tasks. In *CVPR*, pages 1851–1860, 2019. 5
- [13] David R Martin, Charless C Fowlkes, and Jitendra Malik. Learning to detect natural image boundaries using local brightness, color, and texture cues. *IEEE TPAMI*, 26(5): 530–549, 2004. 6
- [14] Brendan McMahan, Eider Moore, Daniel Ramage, et al. Communication-efficient learning of deep networks from decentralized data. In *AISTATS*, pages 1273–1282, 2017. 3, 5, 8
- [15] Roozbeh Mottaghi, Xianjie Chen, Xiaobai Liu, Nam-Gyu Cho, Seong-Whan Lee, Sanja Fidler, Raquel Urtasun, and Alan Yuille. The role of context for object detection and semantic segmentation in the wild. In *CVPR*, pages 891–898, 2014. 5
- [16] Jaehoon Oh, Sangmook Kim, and Se-Young Yun. Fedbabu: Toward enhanced representation for federated image classification. In *ECCV*, 2022. 5 312  
313  
314
- [17] Sangjoon Park, Gwanghyun Kim, Jeongsol Kim, Boah Kim, and Jong Chul Ye. Federated split task-agnostic vision transformer for covid-19 cxr diagnosis. *NeurIPS*, 34:24617–24630, 2021. 5 315  
316  
317  
318
- [18] Jordi Pont-Tuset and Ferran Marques. Supervised evaluation of image segmentation and object proposal techniques. *IEEE TPAMI*, 38(7):1465–1478, 2015. 6 319  
320  
321
- [19] Ozan Sener and Vladlen Koltun. Multi-task learning as multi-objective optimization. *Advances in neural information processing systems*, 31, 2018. 4 322  
323  
324
- [20] Nathan Silberman, Derek Hoiem, Pushmeet Kohli, and Rob Fergus. Indoor segmentation and support inference from rgb-d images. In *ECCV*, pages 746–760, 2012. 5 325  
326  
327
- [21] Hanrong Ye and Dan Xu. Inverted pyramid multi-task transformer for dense scene understanding. In *ECCV*, pages 514–530, 2022. 5 328  
329  
330
- [22] Tianhe Yu, Saurabh Kumar, Abhishek Gupta, Sergey Levine, Karol Hausman, and Chelsea Finn. Gradient surgery for multi-task learning. *NeurIPS*, 33:5824–5836, 2020. 4 331  
332  
333