Supplementary Material for Towards Multi-modal Transformers in Federated Learning

Guangyu Sun¹[©], Matias Mendieta¹[©], Aritra Dutta²[©], Xin Li²[©], and Chen Chen¹[©]

¹ Center for Research in Computer Vision, University of Central Florida, FL, USA ² Department of Mathematics, University of Central Florida, FL, USA {guangyu, matias.mendieta, aritra.dutta, xin.li}@ucf.edu; chen.chen@crcv.ucf.edu

The supplementary material is structured as follows:

- 1. In Section A we sketch a general framework for transfer MFL (§A.1), then we adopt it to the transfer MFL setting in the vision-language (§A.2). We give the convergence guarantee in §A.3. Finally, in §A.4, we argue why multimodal FL might work with diverse data modalities by providing a simple heuristic on the distribution of the multimodal data.
- 2. Section B elaborates on more experimental details, including (i) implementation details, (ii) visualization of the data partitioning, (iii) stochasticity discussion, (iv) breakdown performance of the results in the main paper, (v) communication analysis, (vi) experiments under imbalanced total client numbers, and (vii) visualizations.
- 3. Section C discusses the potential negative social impact and limitation of the proposed method.

A Theoretical Guarantee

In this section, we start with a general transfer multimodal FL framework.

A.1 A general framework for transfer MFL

In this section, we outline the general algorithmic framework of transfer MFL. Consider the *empirical risk minimization* (ERM) problem:

$$\min_{w \in \mathbb{R}^d} \left[\mathcal{F}(w) := \sum_{j=1}^M F_j(w) = \sum_{j=1}^M \left(\underbrace{\frac{1}{N_j} \sum_{i=1}^{N_j} f_{ji}(w)}_{:=F_j} \right) \right],\tag{1}$$

where $f_{ji}(w) = \frac{1}{|D_j|} \sum_{i \in D_j} \mathbb{E}_{z \sim D_i} l(w; z)$ denotes the loss function evaluated on input datapoint at the *i*th client, with D_j denoting the set of all clients whose data has non-empty *j*th modality, *M* is the total number of modality types, and N_j denotes the number of clients for the *j*th modality. At the beginning of the

2 Sun et al.

training process, the server initializes global models, $w^{(0)}$. The entire process runs for T global communication rounds. At the beginning of each global round, r, each client, i, receives the global model, initializes it to $w_i^0 = w^{(r)}$, and updates its local model in each iteration, t. We compactly write it:

$$w_i^{t+1} = w_i^t - \Omega_{\text{local}} \nabla w_i^t, \tag{2}$$

where Ω_{local} is a block diagonal matrix, where each block is of the size of the parameters of that individual modality type. Each client completes local training (2), for *E* local epochs, communicates the model update to the server, and the server aggregates them at the r^{th} round via:

$$\nabla w^{(r)} = \Omega_{\text{Aggregation}}(w_i^E - w^{(r)}). \tag{3}$$

Finally, the update rule at the server for each round r is given as

$$w^{(r+1)} = w^{(r)} + \Omega_{\text{server}} \nabla w^{(r)}, \qquad (4)$$

where Ω_{server} is a symmetric, $(M + 1) \times (M + 1)$ block band matrix.

A.2 Transfer MFL in our setting

We consider a transfer multi-modal FL setting in the vision-language domain with a total of N clients. Based on the general framework, we describe our problem setup. Although (1) presents a general multimodal loss function, in our case, we are working with 3 modality types, so, M = 3. Let D_v, D_l , and D_{vl} present the dataset for vision, text, and multimodal tasks, respectively. Note that, $D_{vl} = \{d : d = (v, l)\}$ contains vision and text pair as input, and define $D_{vl}(v) := P_v(D_{vl}), D_{vl}(l) := P_l(D_{vl})$, where P_y is the projection operator on the set y. For training the vision models, we use $D_v \cup D_{vl}(v)$; for training the language models, we use $D_l \cup D_{vl}(l)$. At the beginning of the FL process, the server initializes global models, $w^{(0)} := (w^{(0,v)}, w^{(0,vl)}, w^{(0,vl_v)}, w^{(0,vl_l)})^{\top}$ the set of parameters, and $\nabla w^r = (\nabla w_v^r \nabla w_l^r \nabla w_{vl_v}^r \nabla w_{vl_v}^r)^{\top}$.

For local training, $\Omega_{\text{local}} = \text{diag}(\eta_{v} \mathbf{I}_{|\mathbf{D}_{v}|} \ \eta_{l} \mathbf{I}_{|\mathbf{D}_{v}|} (\eta_{v} \mathbf{I}_{|\mathbf{D}_{v}|(v)|})$, where $\eta_{v}, \eta_{l}, \eta_{vl} \geq 0$ are stepsizes for respective modality type. Note that, |x| denotes the dimension of an arbitrary vector, x, and $I_{|x|}$ is the identity matrix in the space $\mathbb{R}^{|x| \times |x|}$. Let n_{v}, n_{l} and n_{vl} clients be sampled uniformly at random from each $\{N_{j}\}_{j=1}^{3}$, respectively, in each training round. Hence, after E local epochs, (3) is given as:

$$\begin{cases} \nabla w_v^r := \frac{1}{n_v} \sum_{i \in n_v} (w_{v_i}^E - w_v^r), \nabla w_l^r := \frac{1}{n_l} \sum_{i \in n_l} (w_{l_i}^E - w_l^r), \\ \begin{pmatrix} \nabla w_{vl_v}^r \\ \nabla w_{vl_v}^r \end{pmatrix} := \frac{1}{n_{vl}} \sum_{i \in n_{vl}} \begin{pmatrix} w_{vl_v}^E - w_{vl_v}^r \\ w_{vl_i}^E - w_{vl_l}^r \end{pmatrix}. \end{cases}$$

Our setting allows different structures of Ω_{server} that can be adapted to consider cross-modal contribution. Specifically, we consider Ω_{server} as follows: For

i odd, $(\Omega_{\text{server}})_{ii} = A_v$, and for *i* even, $(\Omega_{\text{server}})_{ii} = B_l$. The superdiagonal blocks, j > i, are given as: (*i*) For *i* odd and *j* even, $(\Omega_{\text{server}})_{ij} = 0$; for *j* odd, $(\Omega_{\text{server}})_{ij} = B_{vl}$; (*ii*) For *i* even and *j* odd, $(\Omega_{\text{server}})_{ii} = 0$, for *j* even, $(\Omega_{\text{server}})_{ij} = D_{vl}$. Constructing Ω_{server} as above allows us to leverage the participation and interaction of each modality over the *vanilla* multimodal FedAvg. For *vanilla* FedAvg [4,7], Ω_{server} is a block diagonal matrix and $n_v = n_l = n_{vl}$.

A.3 Convergence guarantee

Based on the convergence of FedAvg in [4,7], in this section, we will comment on the convergence of general transfer MFL. For ease of notation we consider, at each round, S_j be the number of clients sampled for the the j^{th} modality.

Assumptions. We require the following assumptions.

Assumption 1 (Global minimum) For each $j \in [M]$, there exists w^* such that, $F_j(w^*) = F_j^* \leq F_j(w)$, for all $w \in \mathbb{R}^d$.

Assumption 2 (β -Smoothness) The loss function $f_{ij} : \mathbb{R}^d \to \mathbb{R}$ at each node is β -smooth, i.e. $f_{ij}(y) \leq f_{ij}(x) + \nabla f_{ij}(x)^\top (y-x) + \frac{\beta}{2} ||y-x||^2$ for all $x, y \in \mathbb{R}^d$.

Remark 1. The above assumption implies that F_j is β -smooth for all j.

Assumption 3 For each $j \in [M]$, there exist constants $G_j \ge 0, B_j \ge 1$, such that for all $x \in \mathbb{R}^d$, the stochastic noise, $\xi_{i,t}$ follows

$$\frac{1}{N_j} \sum_{i=1}^{N_j} \|\nabla f_{ij}(x)\|^2 \le G_j^2 + B_j^2 \|F_j(x)\|^2.$$

Assumption 4 (Bounded variance) For each $j \in [M]$, let $g_{ji}(w) := \nabla f_{ji}(w, z_{i(k)})$ be the unbiased stochastic gradient of f_{ji} with bounded variance. That is, there exists, $\sigma_j \geq 0$ such that, $\mathbb{E}_{z_{i(k)}} \left[\|g_{ji}(w) - \nabla f_{ji}(w)\|^2 \right] \leq \sigma_j^2$, for all w, i, where $z_{i(k)}$ is the k^{th} sample data at the i^{th} node.

Assumption 5 The eigenvalues of the symmetric matrix, Ω_{server} are nonnegative.

Remark 2. The above assumption implies that there exists an orthogonal matrix P such that, $\Omega_{\text{server}} = P \Lambda P^{\top}$, where Λ is a diagonal matrix of nonnegative eigenvalues of Ω_{server} .

Remark 3. Based on the previous remark, the update rule (4) can be rewritten as:

$$P^{\top}w^{(r+1)} = P^{\top}w^{(r)} + P^{\top}P\Lambda P^{\top}\nabla w^{(r)}.$$

Consider the change of variable, $\tilde{w}^{(r)} := P^{\top} w^{(r)}$ and hence the above becomes:

$$\tilde{w}^{(r+1)} = \tilde{w}^{(r)} + \Lambda \nabla \tilde{w}^{(r)}.$$
(5)

4 Sun et al.

Finally, we are all set to give our main convergence result based on the vanilla FedAvg framework; for more details see [4,7].

Theorem 1. For each $j \in [M]$, let F_j satisfies Assumptions 1-5. Then

$$\mathbb{E}\left[\|\nabla F_j(\tilde{w}^{(T)})\|^2\right] \le O\left(\frac{\beta\sqrt{(F_j(\tilde{w}^{(0)}) - F_j^{\star})}}{\sqrt{TES_j}}\right)$$

Remark 4. From (1), we have $\mathcal{F}(\tilde{w}) = \sum_{j=1}^{M} F_j(\tilde{w})$. Hence the boundedness of each $\mathbb{E}\left[\|\nabla F_j(\tilde{w}^{(T)})\| \right]$ guarantee the boundedness of $\mathbb{E}\left[\|\nabla \mathcal{F}(\tilde{w}^{(T)})\| \right]$.

A.4 Distribution of the data

In this subsection, we discuss the perspective of multi-modal learning from the learning of the joint distribution of the modalities. In general, we could (and should) not assume independence among the different modalities at hand, and thus, their joint distribution is not simply the product of the marginal distributions. The training datasets then should be samples that reflect the same distributions of each modality feature. Let \mathcal{D} be the joint *unknown* distribution of the input data of two modalities. Let the datasets, D_v and D_l have \mathcal{D}_v and \mathcal{D}_l as their marginal probability distributions of modalities v and l respectively. The availability of the dataset D_{vl} that follows the distribution \mathcal{D} makes it possible to learn the joint distribution when v and l modalities are not independent; see Figure 1. The learning of a joint density model over the space of multimodal inputs is likely to yield a better generalization in various applications [10]. Intuitively, this explains the possibility that multi-modal learning could improve performance when modalities are jointly used in training the parameters even for the individual modality model.

B Numerical Experiments

This section serves as an addendum to the numerical experiments in the original paper.

B.1 Implementation Details

General setup. We use AdamW [6] as the optimizer with a learning rate of 0.0001 with a decay of 0.99 every epoch for local training. The batch size is set as 112 in most cases. All the experiments are implemented under the PyTorch framework and run on $4 \times$ Nvidia A5000 GPUs.



Fig. 1: Multi-modal FL in the vision-language domain with collaboration from different modalities.



Fig. 2: Visualization of the data partitioning of different datasets: CIFAR-100 [5], AG News [13], Flickr10k [9].

CreamFL [11]. In the original CreamFL, there is public data in the server on which the global model can be directly trained. However, we assume no training data on the server following the traditional FL setting. Therefore, we replace the centralized training on the server with an aggregation of the client models. We use 500 samples from the MS-COCO [2] dataset for knowledge distillation and set the optimal distillation and local contrastive weights as 1 and 1e - 7, respectively, after a parameter search.

FedIoT [14]. We follow the original design of FedIoT by applying a factor of 100 to the multi-modal models during aggregation of the transformer blocks.

6 Sun et al.

B.2 Visualization of the data partitioning

We perform non-IID data partitioning to simulate the client data. For CIFAR-100 and AG NEWS datasets, we partition samples of each class with a random Dirichlet distribution with a given α . For Flickr10k, we apply a non-IID number of training samples due to a lack of class labels, following [3]. A visualization of the number of samples on each client is shown in Fig. 2.

B.3 Stochasticity discussion

To study the impact of the stochasticity in the experiments, we additionally conduct experiments with two additional random seeds besides the original seed 1 reported in the main paper and report the standard deviation (STD) of each method. As shown in Table 2, FedCola consistently outperforms all the comparison methods with a significant gap meanwhile holding the smallest STD, indicating the effectiveness and robustness of our proposed method.

B.4 Breakdown performance

To provide more details for the reported results under each setting in the main paper, a breakdown performance with image-to-text top-1 recall (i2t R@1), text-to-image top-1 recall (t2i R@1) under both the 1k and 5k test image settings are given in Table 4 for Flickr and Table 5 for COCO Captions.

B.5 Communication analysis

In §6.1 in the main paper, we study the communication trade-off of the proposed complementary local training. We further propose the communication costs of the comparison methods as a reference. Specifically, we report the size of the total download communication on one image client and one text client. An extended version is shown in Table 1. It shows that even FedCola with collaborative aggregation only (CA-only) can outperform all comparison methods without additional communication overhead.

Table	1:	Com	municat	ion	cost	and	perfor-
mance	of e	each	method	on	Flick	r	

Method	Comm. Cost (MB)	$R@1_{sum}$
FedAvg	208.81	81.08
FedProx	208.81	78.55
CreamFL	211.74	74.83
FedIoT	208.81	85.51
FedCola (CA-only)	208.81	90.09
FedCola (Attn)	262.95	91.73
FedCola	371.26	91.96

Further, when more communication budget is acceptable, FedCola (Attn) can provide a better trade-off between communication cost and performance, while the original FedCola can provide the highest performance. Table 3: Flickr performance under imbalanced to-

Table 2: Performance on Flickrunder different random seeds. Fed-Cola has the lowest standard devi-ation (STD).

Math a d						
Method	1	42	2024	51D ↓		
FedAvg	81.08	79.14	82.04	1.48		
FedProx	78.55	77.86	81.69	2.04		
CreamFL FedIoT	74.83 85.51	$75.94 \\ 80.16$	$78.34 \\ 81.10$	$1.79 \\ 2.86$		
FedCola	91.96	90.80	93.21	1.21		

1k Test Image 5k Test Image $R@1_{sum}$ Setting Method i2t R@1 t2i R@1 i2t R@1 t2i R@1 FedAv 79.04 31.5822.7414.749.98 29.24 20.51 13.64 8.76 72.15 FedProx More Total Image Clients 21.34 73.98 CreamFL 29.5813.849.22FedIoT 32 76 23.36 15.68 10.5382.33 FedCola 37.16 26.07 18.64 12.4694.33 FedAvg 15.48 32.90 23.34 10.39 82.11 FedProx 20.02 14.60 5.64 48.14 7.88 More Total Text Clients CreamFL 30.38 21.8613.829.56 75.62FedIoT 31.88 22.8514.82 10.29 79.84

25.76

17.62

12.06

91.68

36.24

B.6 Imbalanced client scenario

In the main paper, we reported the performance when the number of participating clients is imbalanced and the number of total clients is the same as the default setting, considering the total client numbers in each type of client will only impact the uni-aggregation before the collaboration. To provide more experimental results, we report the performance under there are more image clients ($N_v = 16$ increased from 12) and more text clients ($N_l = 16$ increased from 12) in Table 3. As expected, FedCola still outperforms all comparison methods under such settings.

tal client numbers

FedCola

B.7 Visualization

The smoothness of the parametric loss space has been utilized as a significant indicator of the model generalizabilty [1,8,12]. To illustrate that FedCola learns a more generalized global model, we visualize the loss space on 256 training samples of FedAvg (Fig. 3a) and FedCola (Fig. 3b) when the weights of the model are perturbed along the direction of the top Hessian eigenvectors. The loss landscape of FedCola is significantly smoother than FedAvg, indicating that with the help of the proposed framework, a more generalized global model can be obtained. Additionally, we further conduct visualizations with Linear Discriminant Analysis (LDA) at the feature level, as shown in Fig. 3c. By computing the distance between the feature centers, we find the gaps between uni-modal and multi-modal datasets are reduced under FedCola.

C Potential Negative Societal Impact and Limitation

Potential Negative Societal Impact. The effectiveness of FedCola, like any machine learning model, is contingent on the data it's trained on. Given that data distribution in FL settings can be highly non-uniform and biased towards certain demographics or modalities, there's a risk of amplifying existing biases

Setting	Method	1k Test Image		5k Test Image					1k Test Image		5k Test Image		
		i2t R@1	t2 i $R@1$	i2t $R@1$	t2 i $R@1$	$R@1_{sum}$	Setting	Method	i2t $R@1$	t2 i $R@1$	i2t R@1	t2 i $R@1$	R@1 _{sum}
FedAvg FedProx Default CreamFI FedIoT FedIoT	FedAvg FedProx	32.84 31.36	22.90 22.41	15.32 14.84	10.02 9.94	81.08 78.55	Default	FedAvg FedProx	36.98 37.56	29.28 28.46	$16.76 \\ 16.68$	12.40 12.46	95.42 95.16
	CreamFL FedIoT	30.2 34.42	21.34 23.87	13.82 16.34	9.46 10.88	74.83 85.51		CreamFL FedIoT	37.60 38.62	28.64 29.97	$16.68 \\ 17.16$	12.34 12.65	95.26 98.40
	FedCola	35.68	26.14	18.10	12.04	91.96		FedCola	41.02	31.62	18.74	13.72	105.10
More -	FedAvg FedProx	32.5 31.1	23.34 22.06	15.40 13.9	10.46 9.26	81.70 76.33	More Heterogeneity	FedAvg FedProx	37.46 37.66	29.11 28.86	$16.40 \\ 16.90$	12.35 12.20	95.32 95.62
	CreamFL FedIoT	31.48 33.02	22.59 23.73	$15.74 \\ 15.94$	$10.19 \\ 10.57$	80.00 83.28		CreamFL FedIoT	35.76 37.66	28.11 29.47	$15.74 \\ 16.52$	11.80 12.24	91.41 95.89
	FedCola	36.26	26.06	17.54	11.96	91.82		FedCola	39.62	30.37	17.72	13.12	100.83
Less Participation FedAvg FedProv CreamF FedIoT	FedAvg FedProx	25.84 25.94	18.95 19.02	$11.96 \\ 10.74$	8.06 7.64	64.82 63.33	Less Participation	FedAvg FedProx	32.68 31.20	26.12 25.13	14.22 13.28	10.90 10.27	83.91 79.88
	CreamFL FedIoT	26.94 25.18	19.54 18.13	11.9 11.02	8.47 7.62	66.85 61.94		CreamFL FedIoT	31.58 31.62	25.00 25.35	12.60 13.44	9.79 10.24	78.97 80.65
	FedCola	34.94	25.48	16.60	11.84	88.85		FedCola	40.12	30.47	18.28	13.43	102.30
More FedAv FedPro Image FedIo [*] FedCol	FedAvg FedProx	31.22 31.46	22.68 22.84	14.42 14.90	9.96 10.05	78.28 79.25	More Image	FedAvg FedProx	38.26 37.46	29.33 28.67	17.22 16.80	12.47 12.46	97.28 95.39
	CreamFL FedIoT	32.02 33.22	23.18 23.40	14.74 15.78	10.36 10.34	80.31 82.74		CreamFL FedIoT	36.80 36.86	28.66 29.06	16.02 16.78	12.17 12.34	93.65 95.04
	FedCola	35.42	25.8	17.76	12.26	91.24		FedCola	40.58	31.05	19.26	13.33	104.22
More Text Fedd Creat Fedd Fedd Fedd Fedd Fedd	FedAvg FedProx	31.94 31.20	22.55 22.25	$15.20 \\ 14.44$	10.00 9.70	79.69 77.59	More Text	FedAvg FedProx	38.00 36.96	$28.95 \\ 28.70$	17.26 16.92	12.48 12.38	96.69 94.96
	CreamFL FedIoT	31.96 31.46	23.20 22.22	$15.12 \\ 14.56$	10.47 9.77	80.75 78.02		CreamFL FedIoT	36.68 37.74	28.46 29.47	15.62 17.22	12.06 12.61	92.81 97.04
	FedCola	35.48	25.50	17.40	11.72	90.10		FedCola	39.82	30.32	17.96	12.86	100.96
Fewer C: Image-Text F F	FedAvg FedProx	24.92 24.28	18.01 17.50	10.70 10.22	7.49 7.19	61.12 59.19	Fewer Image-Text	FedAvg FedProx	30.30 29.22	23.78 23.32	12.02 11.32	8.99 9.22	75.10 73.08
	CreamFL FedIoT	23.20 24.68	17.12 17.76	9.64 10.42	7.15 7.47	57.12 60.34		CreamFL FedIoT	29.60 30.80	23.78 23.56	12.18 12.42	9.12 9.36	74.69 76.14
	FodCola	34.06	24.28	16.18	11 16	85.68		FedCola	37.78	28.08	16.80	11.74	94.40

 Table 4: Flickr Breakdown Performance

 Table 5: COCO Breakdown Performance



Fig. 3: Visualization of the parametric loss landscape with Hessian eigenvectors ϵ_0 and ϵ_1 and the extracted features for each resulting global multi-modal model.

or creating new ones. This can lead to unfair models that perform inequitably across different groups or modalities.

Limitations. FedCola currently does not address system heterogeneity, representing a limitation in the present framework. We propose to explore this aspect in future research.

References

1. Chen, X., Hsieh, C.J.: Stabilizing differentiable architecture search via perturbation-based regularization. In: International conference on machine learn-

ing. pp. 1554–1565. PMLR (2020)

- Chen, X., Fang, H., Lin, T.Y., Vedantam, R., Gupta, S., Dollár, P., Zitnick, C.L.: Microsoft coco captions: Data collection and evaluation server. arXiv preprint arXiv:1504.00325 (2015)
- Hahn, S.J., Jeong, M., Lee, J.: Connecting low-loss subspace for personalized federated learning. In: Proceedings of the 28th ACM SIGKDD Conference on Knowledge Discovery and Data Mining. p. 505–515. KDD '22, Association for Computing Machinery, New York, NY, USA (2022)
- Karimireddy, S.P., Kale, S., Mohri, M., Reddi, S., Stich, S., Suresh, A.T.: Scaffold: Stochastic controlled averaging for federated learning. In: International conference on machine learning. pp. 5132–5143. PMLR (2020)
- 5. Krizhevsky, A.: Learning multiple layers of features from tiny images pp. 32–33 (2009)
- Loshchilov, I., Hutter, F.: Decoupled weight decay regularization. In: International Conference on Learning Representations (2019)
- McMahan, H.B., Moore, E., Ramage, D., Hampson, S., y Arcas, B.A.: Communication-efficient learning of deep networks from decentralized data. In: Artificial intelligence and statistics. pp. 1273–1282 (2017)
- Mendieta, M., Yang, T., Wang, P., Lee, M., Ding, Z., Chen, C.: Local learning matters: Rethinking data heterogeneity in federated learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 8397– 8406 (2022)
- Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In: Proceedings of the IEEE international conference on computer vision. pp. 2641–2649 (2015)
- Potamianos, G., Marcheret, E., Mroueh, Y., Goel, V., Koumbaroulis, A., Vartholomaios, A., Thermos, S.: Audio and visual modality combination in speech processing applications. In: The Handbook of Multimodal-Multisensor Interfaces: Foundations, User Modeling, and Common Modality Combinations-Volume 1, pp. 489–543 (2017)
- Yu, Q., Liu, Y., Wang, Y., Xu, K., Liu, J.: Multimodal federated learning via contrastive representation ensemble. In: The Eleventh International Conference on Learning Representations (2022)
- Zela, A., Elsken, T., Saikia, T., Marrakchi, Y., Brox, T., Hutter, F.: Understanding and robustifying differentiable architecture search. In: International Conference on Learning Representations (2020)
- 13. Zhang, X., Zhao, J., LeCun, Y.: Character-level convolutional networks for text classification. Advances in neural information processing systems **28** (2015)
- Zhao, Y., Barnaghi, P., Haddadi, H.: Multimodal federated learning on iot data. In: 2022 IEEE/ACM Seventh International Conference on Internet-of-Things Design and Implementation (IoTDI). pp. 43–54 (2022)