

# Cross-feature Contrastive Loss for Decentralized Deep Learning on Heterogeneous Data

Sai Aparna Aketi  
Purdue University  
USA  
saketi@purdue.edu

Kaushik Roy  
Purdue University  
USA  
kaushik@purdue.edu

## Abstract

*The current state-of-the-art decentralized learning algorithms mostly assume the data distribution to be Independent and Identically Distributed (IID). However, in practical scenarios, the distributed datasets can have significantly heterogeneous data distributions across the agents. In this work, we present a novel approach for decentralized learning on heterogeneous data, where data-free knowledge distillation through contrastive loss on cross-features is utilized to improve performance. Cross-features for a pair of neighboring agents are the features (i.e., last hidden layer activations) obtained from the data of an agent with respect to the model parameters of the other agent. We demonstrate the effectiveness of the proposed technique through an exhaustive set of experiments on various Computer Vision datasets (CIFAR-10, CIFAR-100, Fashion MNIST, Imagenette, and ImageNet), model architectures, and network topologies. Our experiments show that the proposed method achieves superior performance (0.2 – 4% improvement in test accuracy) compared to other existing techniques for decentralized learning on heterogeneous data.*

## 1. Introduction

Every day, substantial volumes of data are generated across the globe, offering the potential to train powerful deep-learning models. Compiling such data for centralized processing is impractical due to communication constraints and privacy concerns. To address this issue, a new interest in developing distributed learning algorithms [1] has emerged. Federated learning (FL) or centralized distributed learning [17] is a popular setting in the distributed machine learning paradigm. In this setting, the training data is kept locally at the edge devices and a global shared model is learned by aggregating the locally computed updates through a coordinating central server. Such a setup requires frequent communication with a central server. This

becomes a potential bottleneck [11] and has led to advancements in decentralized machine learning.

Decentralized learning, a subset of distributed optimization, focuses on learning from data distributed across multiple agents without the need for a central server. It offers many advantages over the traditional centralized approach in core aspects such as data privacy, fault tolerance, and scalability [26]. Research has shown that decentralized learning algorithms can perform comparable to centralized algorithms on benchmark vision datasets [22]. One of the key assumptions to achieve state-of-the-art performance by the decentralized algorithms is that the data is independently and identically distributed (IID) across the agents. In particular, the data is assumed to be distributed in a uniform and random manner across the agents. This assumption does not hold in most real-world settings where the data distributions across the agents are significantly different (non-IID/heterogeneous) [13].

The effect of heterogeneous data in a peer-to-peer decentralized setup is a relatively under-studied problem and an active area of research. Note that, we mainly focus on a common type of non-IID data, widely used in prior works [13, 23]: a skewed distribution of data labels across agents. Recently, several methods were proposed to bridge the performance gap between IID and non-IID data for a decentralized setup. Most of these works either make algorithmic changes to track global information [2, 15, 23, 29, 30] or utilize extra communication rounds [3, 10] to obtain second-order gradient information. In this work, we explore an orthogonal direction of using a data-free knowledge distillation approach to handle heterogeneous data in decentralized setups.

Knowledge distillation methods have been well explored in federated learning (FL) setups with a central server for heterogeneous data [6, 21, 24, 27, 34]. However, these approaches leverage the central server and/or need public dataset access and thus, are not transferable to decentralized setups. In this paper, we propose *Cross-feature Contrastive Loss (CCL)* that improves the performance of decentral-

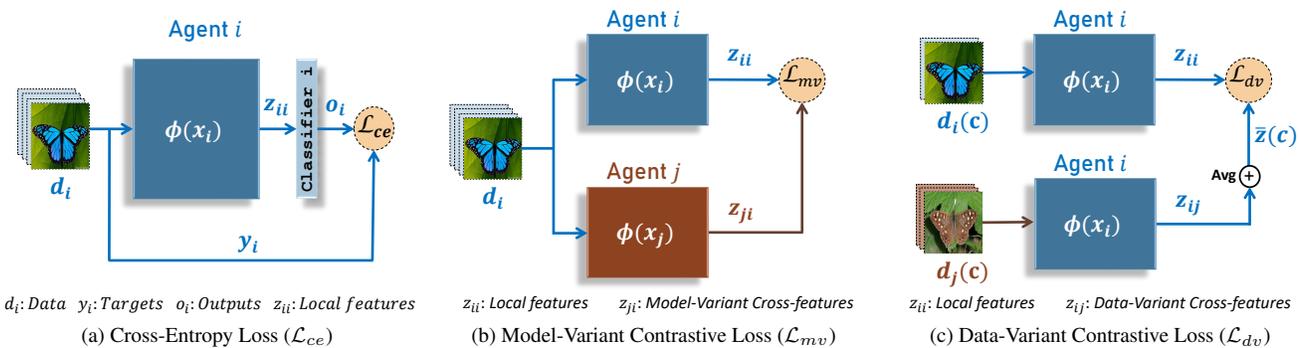


Figure 1. Illustrating different loss components used in the proposed *Cross-features Contrastive Loss*. This illustration describes the loss components with respect to agent  $i$  and assumes that it has only one neighbor  $j$ . The Data-variant Contrastive Loss is only shown for a particular class  $c$  and the same rule will be applied to all the classes in parallel.  $\bar{z}$  is computed at agent  $j$  and then communicated to agent  $i$ .

ized training on heterogeneous data when used along with cross-entropy loss ( $\mathcal{L}_{ce}$ ) at each agent. In particular, at each agent, we introduce two additional contrastive loss terms on cross-features - (a) model-variant contrastive loss ( $\mathcal{L}_{mv}$ ) and (b) data-variant contrastive loss ( $\mathcal{L}_{dv}$ ). Cross-features for a pair of neighboring agents are the features (i.e., last hidden layer activations) obtained from the data of an agent with respect to the model parameters of the other agent. We define two types of cross-features, namely model-variant cross-features and data-variant cross-features. Note that we use *features* as synonymous to the *last hidden layer activations*. Model-variant cross-features are the features obtained from the received neighbors’ model with respect to the local dataset. These cross-features are computed locally at each agent after receiving the neighbors’ model parameters. Communicating the neighbors’ model parameters is a necessary step in any gossip-based decentralized algorithm [22]. Data-variant cross-features are the features obtained from the local model with respect to its neighbors’ datasets. These cross-features are obtained through an additional round of communication.

The  $\mathcal{L}_{mv}$  loss term minimizes the  $L_2$  distance between the local features and model-variant cross-features of the local data at each agent. Whereas  $\mathcal{L}_{dv}$  minimizes the  $L_2$  distance between the local- and data variant cross-feature representations of the same class. Figure. 1 provides an illustration of the loss terms introduced by the proposed framework. We validate the performance of the proposed framework through an exhaustive set of experiments on various vision datasets, model architectures, and graph topologies. We show that the proposed framework achieves superior performance on heterogeneous data compared to the current state-of-the-art method. We also report the communication and compute overheads required for proposed *CCL* as compared to D-PSGD.

**Contributions:** In summary, we make the following contributions.

- We present a novel data-free knowledge distillation-based loss called *Cross-feature Contrastive Loss (CCL)* for decentralized machine learning on heterogeneous data.
- Through an exhaustive set of experiments, we show the advantages of our framework against the current state-of-the-art methods.
- We also report the communication and compute overheads incurred by the proposed framework.

## 2. Background

In this section, we provide the background on decentralized learning algorithms with peer-to-peer connections. Figure. 2 illustrates a decentralized setup with 5 agents connected in a ring topology.

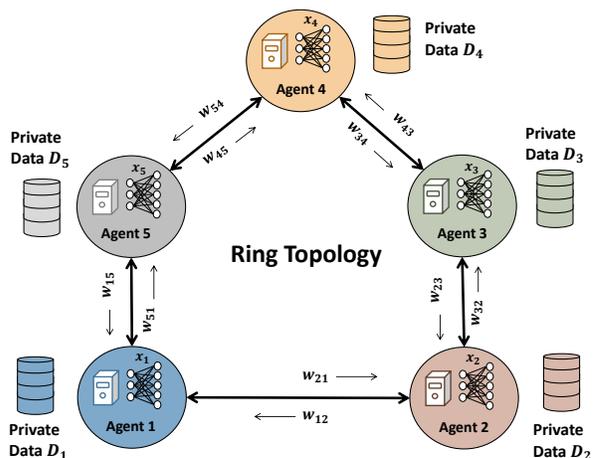


Figure 2. Decentralized training setup with 5 agents connected in a ring topology. Each agent has its own private dataset and a local model.

The main goal of decentralized machine learning is to learn a global model using the knowledge extracted from the locally generated and stored data samples across  $n$  edge devices/agents while maintaining privacy constraints. In particular, we solve the optimization problem of minimizing global loss function  $f(x)$  distributed across  $n$  agents as given in equation. 1. Note that  $f_i$  is a local loss function (for example, cross-entropy loss  $\mathcal{L}_{ce}$ ) defined in terms of the data sampled ( $d_i$ ) from the local dataset  $D_i$  at agent  $i$  with model parameters  $x_i$ .

$$\min_{x \in \mathbb{R}^d} f(x) = \frac{1}{n} \sum_{i=1}^n f_i(x), \quad (1)$$

and  $f_i(x) = \mathbb{E}_{d_i \in D_i} [F_i(x; d_i)] \quad \forall i$

This is typically achieved by combining stochastic gradient descent [7] with global consensus-based gossip averaging [32]. The communication topology in this setup is modeled as a graph  $G = ([n], E)$  with edges  $\{i, j\} \in E$  if and only if agents  $i$  and  $j$  are connected by a communication link exchanging the messages directly. We represent  $\mathcal{N}_i$  as the neighbors of  $i$  including itself. It is assumed that the graph  $G$  is strongly connected with self-loops i.e., there is a path from every agent to every other agent. The adjacency matrix of the graph  $G$  is referred to as a mixing matrix  $W$  where  $w_{ij}$  is the weight associated with the edge  $\{i, j\}$ . Note that, weight 0 indicates the absence of a direct edge between the agents. We assume that the mixing matrix is doubly stochastic and symmetric, similar to all previous works in decentralized learning [22, 23]. Further, the initial models and all the hyperparameters are synchronized at the beginning of the training.

---

**Algorithm 1** Decentralized Learning with *DSGDm* [22]

---

**Input:** Each agent  $i \in [1, N]$  initializes model weights  $x_i^{(0)}$ , step size  $\eta$ , momentum coefficient  $\beta$ , and mixing matrix  $W = [w_{ij}]_{i,j \in [1, N]}$ .

Each agent simultaneously implements the TRAIN() procedure

1. **procedure** TRAIN()
  2.   **for**  $k=0, 1, \dots, K-1$  **do**
  3.      $d_i^k \sim D_i$
  4.      $g_i^k = \nabla_x F_i(d_i^k; x_i^k)$
  5.      $m_i^k = \beta m_i^{(k-1)} + g_i^k$
  6.      $x_i^{k+\frac{1}{2}} = x_i^k - \eta m_i^k$
  7.     SENDRECEIVE( $x_i^{k+\frac{1}{2}}$ )
  8.      $x_i^{(k+1)} = \sum_{j \in \mathcal{N}_i} w_{ij} x_j^{k+\frac{1}{2}}$    // gossip averaging
  9. **return**
- 

Algorithm. 1 describes the flow of Decentralized Stochastic Gradient Descent with momentum (DSGDm).

There are three main stages in any traditional decentralized learning method - (a) Local update, (b) Communication, and (c) Gossip averaging. At every iteration, each agent computes the gradients using local data and updates its model parameters as shown in line 6 of Alg. 1. Then these updated model parameters are communicated to the neighbors as shown in line 7 of Alg. 1. Finally, in the gossip averaging step, the local model parameters are averaged with the received model parameters of the neighbors using the mixing weights (shown in line 8 of Alg. 1). The convergence of the DSGDm algorithm assumes the data distribution across the agents to be Independent and Identically Distributed (IID).

### 3. Related Work

Decentralized Parallel Stochastic Gradient Descent (DSGD) [22] is the first work to show that decentralized algorithms can converge at the same rate as their centralized counterparts [8]. DSGD algorithm combines Stochastic Gradient Descent (SGD) with a gossip averaging algorithm [32]. A momentum version of DSGD referred to as Decentralized Momentum Stochastic Gradient Descent (DSGDm) was proposed in [5]. Further, Stochastic Gradient Push (SGP) [4] extends the scope of DSGD to directed and time-varying graphs. Recently, a unified framework for analyzing gossip-based decentralized SGD methods and the best-known convergence guarantees was presented in [16]. However, all of these above-mentioned algorithms assume the data distribution to be IID.

One of the core challenges in decentralized learning is tackling data that is not identically distributed among agents. A wide range of algorithms were proposed in the literature to deal with the heterogeneous data. The Methods such as Gradient Tracking [15], and Momentum Tracking [29] track the global gradient and use it for the local update. This reduces variation in the local gradients across the agents and hence is more robust to heterogeneous data. Similarly, CGA [10] and NGC [3] also improve the performance by reducing the variation local gradient by utilizing the cross-gradient information. However, all these techniques incur  $2 \times$  communication overhead. The authors in [23] introduce Quasi-Global Momentum (QGM), a decentralized learning method that mimics the global synchronization of momentum buffer to mitigate the difficulties of decentralized learning on heterogeneous data. Recently, RelaySGD was presented in [30] that replaces the gossip averaging step with RelaySum [33]. RelaySGD improves the performance of heterogeneous data by utilizing the delayed information in the RelaySum step. However, this technique only works on a spanning tree and the improvements do not scale well with the graph size.

Knowledge distillation methods are well established for heterogeneous data in federated learning setups with a cen-

tral server. However, there are only a few methods [20, 27] that explore knowledge distillation for decentralized learning on heterogeneous data. Decentralized federated learning via mutual knowledge transfer (Def-KT) [20] replaces gossip averaging with mutual knowledge transfer-based model fusion. In the Def-KT method, only a subset of agents are trained at a time while the other agents participate in model fusion. This deviates from the standard decentralized setup we use where all agents are trained parallelly. In-Distribution Knowledge Distribution (IDKD) proposed in [27] uses a public dataset and an Out-Of-Distribution detector to homogenize the data across decentralized agents. Orthogonal to these methods, we explore data-free knowledge distillation across agents through the proposed *Cross-feature Contrastive Loss*. We compare the proposed *CCL* method with QGM [23] and RelaySGD [30], the current state-of-the-art methods in decentralized learning on heterogeneous data that do not incur any communication overhead or public dataset access.

#### 4. Cross-feature Contrastive Loss

We propose the *Cross-feature Contrastive Loss (CCL)* which aims to improve the performance of decentralized learning on heterogeneous data. *CCL* introduces the concept of cross-features.

---

##### Algorithm 2 Decentralized Learning with *CCL*

---

**Input:** Each agent  $i \in [1, n]$  initializes model weights  $x_i^{(0)}$ , step size  $\eta$ , momentum coefficient  $\beta$ , averaging rate  $\gamma$ , contrastive loss coefficients  $\lambda_m, \lambda_d$ , mixing matrix  $W = [w_{ij}]_{i,j \in [1,n]}$ , number of classes  $C$ ,  $\mathcal{N}_i$  represents neighbors of  $i$  including itself.

Each agent simultaneously implements TRAIN() procedure

1. **procedure** TRAIN()
  2.   **for**  $k = 0, 1, \dots, K - 1$  **do**
  3.     SENDRECEIVE( $x_k^i$ )
  4.      $d_i^k \sim D_i$
  5.     **for each neighbor**  $j \in \mathcal{N}_i$  **do**
  6.        $z_{ji}^k = \phi(x_j^k; d_i^k)$
  7.       Compute  $\bar{z}_{ji}^k(c)$  i.e., the class-wise sum
  8.       SENDRECEIVE( $\{\bar{z}_{ji}^k(c), \text{count}(c)\}_{c=1}^C$ )
  9.     **end**
  10.     $\bar{z}^k(c) = \frac{1}{|c|} \sum_j \bar{z}_{ji}^k(c) \quad \forall c \in [1, C]$
  11.     $CCL_i = \lambda_d \mathcal{L}_{dv}(z_{ii}^k, \bar{z}^k) + \lambda_m \mathcal{L}_{mv}(z_{ii}^k, \{z_{ji}\}_{j \in \mathcal{N}_i})$
  12.     $g_i^k = \nabla_x [\mathcal{L}_{ce}(x_i^k, d_i^k) + CCL_i]$
  13.     $m_i^k = \beta \hat{m}_i^{(k-1)} + g_i^k$
  14.     $x_i^{(k+1)} = (\sum_{j \in \mathcal{N}_i} w_{ij} x_j^k) - \eta m_i^k$
  15.     $\hat{m}_i^k = \beta \hat{m}_i^{(k-1)} + (1 - \beta) \frac{x_i^k - x_i^{(k+1)}}{\eta}$
  16.    **end**
  17. **return**
- 

**Cross-features ( $z_{ij}$ ):** For an agent  $i$  with model parameters  $x_i$  connected to neighbor  $j$  that has local dataset  $D_j$ , the cross-features are the last layer hidden representation obtained from the model parameters  $x_i$ , evaluated on mini-batch  $d_j$  sampled from dataset  $D_j$ .

$$z_{ij} = \phi(x_i; d_j) \quad (2)$$

Note that  $\phi$  represents the neural network up to the last hidden layer (excluding the classifier) and all the definitions are provided with respect to an agent  $i$ .  $z_{ii}$  represents the local feature representation i.e.,  $\phi(x_i, d_i)$ . We define two types of cross-features, namely model-variant and data-variant cross-features. *Model-variant cross-features* ( $\{z_{ji} | j \in \mathcal{N}_i\}$ ) are the features obtained from the received neighbors' model  $x_j$  with respect to the local dataset  $d_i$ . These cross-features are computed locally at each agent after receiving the neighbors' model parameters. *Data-variant cross-features* ( $\{z_{ij} | j \in \mathcal{N}_i\}$ ) are the features obtained from the local model  $x_i$  with respect to its neighbors' datasets  $d_j$ . These cross-features are received through an additional communication round.

Inspired by knowledge distillation methods, we introduce two different contrastive loss terms on cross-features. (a) **Model-variant contrastive loss ( $\mathcal{L}_{mv}$ ):** At each agent  $i$ ,  $\mathcal{L}_{mv}$  minimizes the  $L_2$  distance between the local feature representation  $z_{ii}$  and the model-variant cross-features  $z_{ji}$  for each data-point  $q \in d_i$ .

$$\mathcal{L}_{mv}(z_{ii}, \{z_{ji}\}_{j \in \mathcal{N}_i}) = \sum_{j \in \mathcal{N}_i} \frac{1}{|d_i|} \sum_{q \in d_i} \|z_{ii}^q - z_{ji}^q\|_2^2 \quad (3)$$

The model-variant contrastive loss ensures that the model parameters on the neighboring agents are similar by enforcing the models to have similar representations for a given input sample. This reduces the variation in model parameters across the agents caused by the data heterogeneity.

(b) **Data-variant contrastive loss ( $\mathcal{L}_{dv}$ ):** To compute this loss, we first generate the neighborhood's representation  $\bar{z}(c)$  for each class  $c \in [1, C]$  using the data-variant cross-features  $z_{ij}$ 's as shown in Equation. 4. Now at every agent  $i$ ,  $\mathcal{L}_{dv}$  minimizes the  $L_2$  distance between the local representation  $z_{ii}$  and the neighborhood's representation of the same class.

$$\bar{z}(c) = \frac{1}{|c|} \sum_{j \in \mathcal{N}_i} \sum_{q \in d_j} z_{ij}^q \mathbb{1}_c(z_{ij}^q) \quad \forall c \in [1, C] \quad (4)$$

$$\mathcal{L}_{dv}(z_{ii}, \bar{z}) = \frac{1}{|d_i|} \sum_{q \in d_i} \sum_{c=1}^C \|z_{ii}^q - \bar{z}(c)\|_2^2 * \mathbb{1}_c(z_{ii}^q)$$

Here  $|c|$  represents the total number of samples in the set  $\{z_{ij} | j \in \mathcal{N}_i\}$  that belongs to class  $c$  and  $|d_i|$  is the mini-batch size.  $\mathbb{1}_c(z)$  is an indicator function that outputs 1 if  $z$

belongs to class  $c$ . Since  $\mathcal{L}_{dv}(z_{ii}, \bar{z})$  only uses the averaged representation of data-variant cross-features for each class, we sum these cross-features class-wise and communicate this sum along with the class count to the respective neighbors. The data-variant contrastive loss ensures that the feature representations of a particular class are similar across the agents reducing the disparities caused due to data heterogeneity.

Algorithm. 2 presents the decentralized learning algorithm with the proposed *CCL*. Each agent minimizes the contrastive loss terms along with the traditional cross-entropy loss as shown in Equation. 5.  $\lambda_m$  and  $\lambda_d$  are the hyper-parameters that weigh the model-variant and data-variant contrastive loss respectively.

$$\mathcal{L}_i = \mathcal{L}_{ce}(x_i, d_i) + \lambda_m \mathcal{L}_{mv}(z_{ii}, \{z_{ji}\}_{\forall j}) + \lambda_d \mathcal{L}_{dv}(z_{ii}, \bar{z}) \quad (5)$$

It has been shown in [23] that quasi-global momentum works better than local momentum for decentralized learning with heterogeneous data. Hence, we employ quasi-global momentum with the proposed loss as described in Algorithm. 2.

## 5. Experiments

In this section, we analyze the performance of the proposed Cross-feature Contrastive Loss and compare it with the baselines DSGDm-N [22], RelaySGD [30], and the current state-of-the-art QG-DSGDm-N [23]. The source code is available at [https://github.com/aparna-aketi/Cross\\_feature\\_Contrastive\\_Loss](https://github.com/aparna-aketi/Cross_feature_Contrastive_Loss)

### 5.1. Experimental Setup

The efficiency of the proposed algorithm is demonstrated through experiments on a diverse set of datasets, model architectures, graph topologies, and graph sizes. We present the analysis on – (a) **Datasets:** CIFAR-10 [18], CIFAR-100 [18], Fashion MNIST [31], Imagenette [14], and ImageNet [9]. (b) **Model architectures:** ResNet-20, ResNet-18 [12], LeNet-5 [19] and, MobileNet-V2 [28]. All the models except LeNet-5 use Evonorm [25] as the activation-normalization layer as it is shown to be better suited for decentralized learning on heterogeneous data [23]. LeNet-5 has no normalization layers. (c) **Graph topologies:** Ring graph with 2 peers per agent, Dyck graph with 3 peers per agent, and Torus graph with 4 peers per agent (refer Figure 3). (d) **Number of agents:** 8 to 40 agents.

For the decentralized setup, we use an undirected ring, undirected Dyck graph, and undirected torus graph topologies with a uniform mixing matrix. The undirected ring topology for any graph size has 3 peers per agent including itself and each edge has a weight of  $\frac{1}{3}$ . The undirected Dyck topology with 32 agents has 4 peers per agent including itself and each edge has a weight of  $\frac{1}{4}$ . The undirected

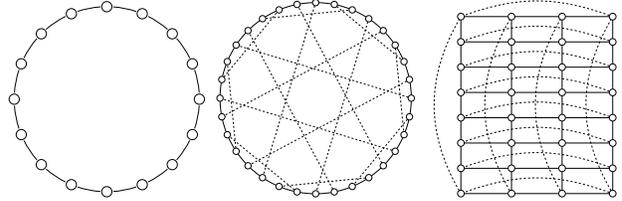


Figure 3. Ring (left), Dyck (center), and Torus (right).

torus topology with 32 agents has 5 peers per agent including itself and each edge has a weight of  $\frac{1}{5}$ . RelaySGD baseline only works on the spanning trees. Therefore, for a fair comparison, we use an undirected chain topology (spanning tree of ring topology) for all the RelaySGD experiments. We use the Dirichlet distribution to generate disjoint non-IID data across the agents similar to [23]. The partitioned data across the agents is fixed, non-overlapping, and never shuffled across agents during the training. The degree of heterogeneity is regulated by the value of  $\alpha$  – the smaller the  $\alpha$ , the larger the non-IIDness across the agents.

The initial learning rate is either set to 0.1 (CIFAR-10, CIFAR-100) or 0.01 (Fashion MNIST, Imagenette) and is decayed by a factor of 10 after 50% and 75% of the training. We use a weight decay of  $1e^{-4}$  and a mini-batch size of 32 per agent in all the experiments. We use the Nesterov version of the Quasi-Global Momentum with a momentum coefficient of 0.9. The stopping criteria is a fixed number of epochs. The experiments on CIFAR-10 are run for 200 epochs, CIFAR-100 and Imagenette for 100 epochs, and Fashion MNIST for 50 epochs. Note, DSGDm-N indicates Decentralized Stochastic Gradient Descent with Nesterov momentum, QG-DSGDm-N and the proposed *CCL* uses a DSGD optimizer with Nesterov version of Quasi-Global Momentum. DSGDm-N, QG-DSGDm-N, and RelaySGD utilized the cross-entropy loss whereas our framework uses the proposed *Cross-feature Contrastive Loss* along with the cross-entropy loss. We use grid search on the set  $\{1, 0.1, 0.01, 0.001\}$  to obtain the hyper-parameters  $\lambda_m, \lambda_d$  for *CCL*. We report the test accuracy of the consensus model averaged over three randomly chosen seeds. All our experiments were conducted on a system with an NVIDIA A40 card with 4 GPUs. A detailed description of the setup and hyperparameters for each experiment is presented in Appendix. A.1.

### 5.2. Results

We evaluate the effectiveness of the *Cross-feature Contrastive Loss* through an exhaustive set of experiments. We show that the proposed loss terms improve the accuracy as compared to simple cross-entropy loss.

Table. 1 presents the average test accuracy for training

Table 1. Test accuracy of different decentralized algorithms evaluated on CIFAR-10, distributed with different degrees of heterogeneity (non-IID) trained on ResNet-20 over ring topologies. The results are averaged across all agents over three seeds where the standard deviation is indicated. We also include the results of the IID baseline as DSGDm-N (IID) where the local data is randomly partitioned independent of  $\alpha$ .

Agents ( $n$ )	Method	ResNet-20	
		$\alpha = 0.1$	$\alpha = 0.01$
16	DSGDm-N (IID)	89.61 $\pm$ 0.43	
	DSGDm-N [22]	80.60 $\pm$ 0.50	58.78 $\pm$ 9.63
	RelaySGD [30]	73.81 $\pm$ 1.97	34.33 $\pm$ 2.42
	QG-DSGDm-N [23]	85.95 $\pm$ 1.64	77.16 $\pm$ 7.02
	CCL (ours)	<b>86.63</b> $\pm$ 0.52	<b>81.29</b> $\pm$ 0.36
32	DSGDm-N (IID)	88.13 $\pm$ 0.32	
	DSGDm-N [22]	76.46 $\pm$ 1.32	53.08 $\pm$ 5.12
	RelaySGD [30]	72.22 $\pm$ 2.58	38.16 $\pm$ 1.34
	QG-DSGDm-N [23]	84.91 $\pm$ 0.56	75.70 $\pm$ 0.80
	CCL (ours)	<b>85.25</b> $\pm$ 0.52	<b>77.60</b> $\pm$ 1.58

ResNet-20 architecture on the CIFAR-10 dataset with varying degrees of heterogeneity over the ring topology of 16 and 32 agents. We observe that *CCL* outperforms QG-DSGDm-N for all models, graph sizes, and degree of heterogeneity with a significant performance gain varying from 0.34 – 4.13%. We also notice that the proposed framework has less variation in accuracy across various initial weight initialization (seeds) compared to the baselines. In our settings, we find that Relay-SGD with local momentum performs worse than DSGDm-N and doesn’t scale with graph size. Note that *Cross-feature Contrastive Loss* is an orthogonal technique to RelaySGD and can be used in synergy with RelaySGD.

Table 2. Test accuracy of CIFAR-10 dataset trained on ResNet-20 over various graph topologies with 32 agents.

Graph	Method	$\alpha = 0.1$	$\alpha = 0.01$
Dyck	DSGDm-N (IID)	88.89 $\pm$ 0.10	
	QG-DSGDm-N	86.20 $\pm$ 0.38	78.18 $\pm$ 4.01
	CCL (ours)	<b>86.78</b> $\pm$ 0.41	<b>80.63</b> $\pm$ 1.54
Torus	DSGDm-N (IID)	88.86 $\pm$ 0.31	
	QG-DSGDm-N	87.75 $\pm$ 0.39	81.74 $\pm$ 0.87
	CCL (ours)	<b>88.14</b> $\pm$ 0.36	<b>82.30</b> $\pm$ 0.56

To demonstrate the scalability and generalizability of *CCL*, We present the results on various graph topologies and datasets. Firstly, we train the CIFAR-10 dataset on ResNet-20 over the Dyck and Torus graphs with 32 agents to illustrate the impact of connectivity on the proposed framework. We observe a 0.4 – 2.5% performance gains with varying connectivity (or spectral gap) as seen in Table. 2. This

shows that the gains from the proposed technique are more pronounced in graphs with less connectivity such as ring graphs. We then evaluate *CCL* on various image datasets such as Fashion MNIST and Imagenette and on challenging datasets such as CIFAR-100 and ImageNet. The proposed *CCL* outperforms QG-DSGDm-N by 0.2 – 2.4% across various datasets as shown in Table. 3. Further, Table. 4 shows that the proposed method can achieve an average improvement of 1.1% on the large-scale ImageNet dataset.

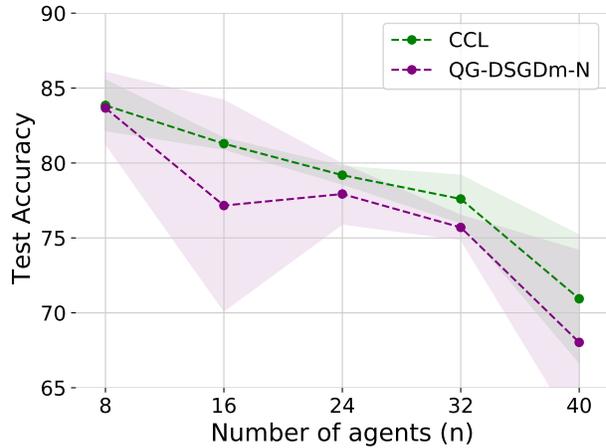


Figure 4. Test accuracy for the CIFAR-10 dataset trained on ResNet-20 architecture over varying sizes of ring topology with a skew of  $\alpha = 0.01$ .

We then evaluate the scalability in decentralized settings by training CIFAR-10 on varying the size of the ring topology between 8 and 40 as shown in Figure. 4. We observe that the proposed *CCL* framework consistently outperforms the QG-DSGDm-N baseline over varying graph sizes by an average improvement of 2%. In summary, the proposed *Cross-feature Contrastive Loss* makes the decentralized training more robust to heterogeneity in the data distribution and has superior performance to all the comparison methods with an average improvement of 1.3%. Additional results are presented in Appendix. A.2.

### 5.3. Analysis

In this section, we analyze the various aspects of the proposed *CCL* terms such as the effect of the contrastive loss on IID vs. non-IID data, the choice of the loss function, and the contribution of each loss term.

The proposed *CCL* framework minimizes three different loss functions namely  $\mathcal{L}_{ce}$ ,  $\mathcal{L}_{mv}$ , and  $\mathcal{L}_{dv}$  where as the baseline methods (DSGDm-N, QG-DGSDm-N) only focus on  $\mathcal{L}_{ce}$ . We hypothesize that the cross-entropy loss  $\mathcal{L}_{ce}$  alone does not capture the data heterogeneity across the agents. To address this we add two different contrastive loss terms explicitly representing the data skew. Figure. 5 measures

Table 3. Test accuracy of different decentralized algorithms evaluated on various datasets, distributed with different degrees of heterogeneity over 16 agents ring topology with 16 agents.

Method	Fashion MNIST (LeNet-5)		CIFAR-100 (ResNet-20)		Imagenette (MobileNet-V2)	
	$\alpha = 0.1$	$\alpha = 0.01$	$\alpha = 0.1$	$\alpha = 0.01$	$\alpha = 0.1$	$\alpha = 0.01$
DSGDm-N (IID)	90.95 $\pm$ 0.09		59.72 $\pm$ 1.00		74.17 $\pm$ 0.83	
QG-DSGDm-N	88.92 $\pm$ 0.50	87.15 $\pm$ 0.64	52.33 $\pm$ 3.59	44.12 $\pm$ 6.85	65.94 $\pm$ 1.17	51.47 $\pm$ 2.67
CCL (ours)	<b>90.21</b> $\pm$ 0.34	<b>87.42</b> $\pm$ 0.78	<b>54.20</b> $\pm$ 0.86	<b>46.49</b> $\pm$ 4.19	<b>66.14</b> $\pm$ 0.84	<b>52.87</b> $\pm$ 5.15

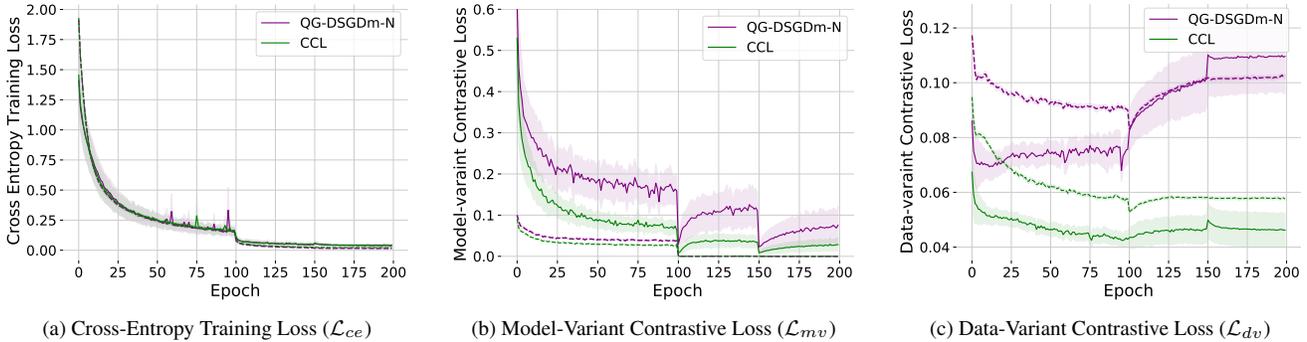


Figure 5. Comparing various training loss terms for IID (dashed lines) and non-IID (solid lines) partitions of CIFAR-10 trained on ResNet-20 over a ring topology of 16 agents. We use  $\alpha = 10$  for IID data and  $\alpha = 0.1$  for non-IID data.

Table 4. Test accuracy of ImageNet trained on ResNet-18 architecture over ring topology with 16 agents.

Graph	Method	$\alpha = 1$	$\alpha = 0.1$
Ring	DSGDm-N (IID)	65.62 $\pm$ 0.03	
	QG-DSGDm-N	64.09 $\pm$ 1.49	58.11 $\pm$ 3.81
	CCL (ours)	<b>64.64</b> $\pm$ 1.09	<b>59.82</b> $\pm$ 1.75

the different training losses for both IID and non-IID distribution of the CIFAR-10 dataset. We observe that the training cross-entropy loss (Fig. 5a) for IID and non-IID data converges to zero even though there is a huge gap in the validation loss. However, Fig. 5b shows that the model-variation contrastive loss for the baseline is much higher in non-IID settings compared to IID and hence is a good measure of data-heterogeneity. On the other hand, data-variant contrastive loss measures the variation in class representations across agents. Fig. 5c shows that this variation is relatively stable throughout the training process for QG-DSGDm-N with IID Data. However, for QG-DSGDm-N with the non-IID setting, a significant increase in the variation of class representations across agents is evident. Note that the absolute value of the data-variant contrastive loss for the non-IID setting is lower than the IID setting because of limited shared class samples among neighboring agents (say 6 out of 32 in a mini-batch). The proposed

CCL framework explicitly minimizes the model-variant and data-variant contrastive loss. Fig. 5b, 5c show that the proposed framework reduces the model variance and stabilizes the variance in class representations across agents resulting in better performance on heterogeneous data.

Table 5. Test accuracy of CIFAR-10 dataset trained with different contrastive loss functions on ResNet-20 architecture over a ring topology.

Loss function	Agents	$\alpha = 0.1$	$\alpha = 0.01$
L1 Loss	16	85.76 $\pm$ 1.74	80.43 $\pm$ 2.70
MSE Loss		<b>86.16</b> $\pm$ 0.67	81.29 $\pm$ 0.36
Cosine Loss		86.02 $\pm$ 0.78	<b>82.36</b> $\pm$ 0.93
L1 Loss	32	<b>85.76</b> $\pm$ 0.32	76.13 $\pm$ 2.59
MSE Loss		85.25 $\pm$ 0.52	<b>77.60</b> $\pm$ 1.58
Cosine Loss		85.71 $\pm$ 0.27	75.71 $\pm$ 3.73

Cross-feature Contrastive Loss reduces the similarity between local feature representations and cross-feature representations. The similarity between the two representations can be determined using various measures. We explore three different similarity measures (or the loss functions) for the CCL namely L1 loss, Mean Square Error (MSE) loss, and Cosine loss during the training phase. L1 loss measures the  $L_1$  distance between the local and cross-features, MSE loss measures the  $L_2$  distance, and Cosine loss measures

the cosine distance. We observe that on average MSE loss provides better improvements as shown in Table. 5.

Table 6. Test accuracy of CIFAR-10 dataset trained with different components of contrastive loss on ResNet-20 over a ring topology.

$\mathcal{L}_{ce}$	$\mathcal{L}_{mv}$	$\mathcal{L}_{dv}$	$\alpha = 0.1$	$\alpha = 0.01$
✓	x	x	85.95 ± 1.64	77.16 ± 7.02
✓	✓	x	<b>86.63</b> ± 0.52	80.55 ± 1.61
✓	x	✓	85.67 ± 1.58	77.78 ± 4.01
✓	✓	✓	86.16 ± 0.67	<b>81.29</b> ± 0.36

Cross-feature Contrastive Loss introduces two loss terms namely model-variant contrastive loss and data-variant contrastive loss. We evaluate the contribution of each of these loss terms to the improved accuracy in Table. 6. For lower skew ( $\alpha = 0.1$ ), we observe that the accuracy improvements can be mostly attributed to the addition of model-variant contrastive loss ( $\mathcal{L}_{mv}$ ). Even in the case of higher skew ( $\alpha = 0.01$ ), the improvement can be majorly attributed to the model-variant contrastive loss. However, the maximum average test accuracy is obtained by adding both data-variant and model-variant contrastive loss terms.

## 6. Discussion and Limitations

The proposed *Cross-feature Contrastive Loss* has two potential limitations – (a) Compute overhead of the model-variant cross-features and (b) Communication overhead of the data-variant cross-features. Each agent has to compute model-variant cross-features at every iteration. This requires every agent to perform  $p$  additional forward passes where  $p$  is the number of peers/neighbors per agent. Assume that  $c_f$  is the compute required for the forward pass. Now, the compute overhead can be given as  $\mathcal{O}(p * c_f)$ . Quantitatively, we measure the compute overhead as the fraction of additional compute required for the model-variant cross-features computation (Equation. 6).

$$\text{Compute overhead} = \frac{\text{Compute for cross-features}}{\text{Total compute}} \quad (6)$$

Table. 7 presents the compute overhead for various settings. We observe that for a ring topology, the compute overhead is around 35 – 40%. This overhead shoots up to 57% for a torus graph. Note that the compute overhead depends on the number of peers per agent rather than the total graph size.

Every agent communicates the class-wise sum of data-variant cross-features and class count along with the model parameters to each of their neighbors. The overhead is from the communication of data-variant cross-features. For example, for a dataset with  $C$  classes and  $r$  is the size of a cross-feature, every agent communicates  $C$  cross-features of size  $r$  (one for each class) and a vector of size  $C$

Table 7. Compute overhead incurred per agent due to *Cross-feature Contrastive Loss*.

Dataset	Model	Peers	Compute overhead
Fashion-MNIST	LeNet-5	2	0.351
CIFAR-10	ResNet-20	2	0.397
CIFAR-10	ResNet-20	3	0.496
CIFAR-10	ResNet-20	4	0.567
CIFAR-100	ResNet-20	2	0.397
ImageNette	MobileNet-V2	2	0.397

indicating the number of samples per class in the mini-batch. Thus, the communication overhead can be given as  $\mathcal{O}(p * C * (r + 1))$ . This overhead is negligible compared to the communication of model parameters. Table. 8 compares the communication cost per iteration of *CCL* with *QG-DSGDm-N* in MegaBytes (MB) for training various datasets over a 16-agents ring topology. We observe that the communication overhead is around 0.2% for CIFAR-10, 2.3% for CIFAR-100, 1.4% for Fashion-MNIST, and 0.6% for ImageNette. This shows that the communication overhead incurred by the proposed framework is insignificant.

Table 8. Communication cost per agent per iteration over a ring graph of 16 agents.

Dataset	Model	Method	Comm. Cost (MB)
Fashion-MNIST	LeNet-5	QG-DSGDm-N	0.471 (1×)
		<i>CCL</i>	0.477 (1.013×)
CIFAR-10	ResNet-20	QG-DSGDm-N	2.079 (1×)
		<i>CCL</i>	2.084 (1.002×)
CIFAR-100	ResNet-20	QG-DSGDm-N	2.123 (1×)
		<i>CCL</i>	2.173 (1.023×)
ImageNette	MobileNet-V2	QG-DSGDm-N	17.52 (1×)
		<i>CCL</i>	17.62 (1.006×)

Further, a minor limitation of the proposed *CCL* is that it adds two additional hyper-parameters  $\lambda_m, \lambda_d$ . These hyper-parameters need to be tuned similarly to the learning rate. We used a grid search to find these hyperparameters. We observed that these hyperparameters typically take a value of 0.1 or 0.01. We consider the exploration of the compute efficient *CCL* and adaptive *CCL* as potential future research directions.

## 7. Conclusion

Supporting decentralized training on heterogeneous data is a critical factor for machine learning applications to effectively harness the vast amounts of user-generated private

data. In this paper, we propose a novel *Cross-feature Contrastive Loss* which is better suited for decentralized learning on heterogeneous data. In particular, minimizing proposed contrastive loss terms increases the similarity of local feature representations with the model-variant and data-variant cross-features. We evaluate the *CCL* with Quasi-Global Momentum through an exhaustive set of experiments on various datasets, model architectures, and graph topologies. Our experiments confirm the superior performance of the *Cross-feature Contrastive Loss* compared to existing state-of-the-art methods for decentralized learning on heterogeneous data.

## Acknowledgements

This work was supported in part by, the Center for the Co-Design of Cognitive Systems (COCOSYS), a DARPA-sponsored JUMP center, the Semiconductor Research Corporation (SRC), and the National Science Foundation.

## References

- [1] Alekh Agarwal and John C Duchi. Distributed delayed stochastic optimization. *Advances in neural information processing systems*, 24, 2011. **1**
- [2] Sai Aparna Aketi, Abolfazl Hashemi, and Kaushik Roy. Global update tracking: A decentralized learning algorithm for heterogeneous data. *arXiv preprint arXiv:2305.04792*, 2023. **1**
- [3] Sai Aparna Aketi, Sangamesh Kodge, and Kaushik Roy. Neighborhood gradient clustering: An efficient decentralized learning method for non-iid data distributions. *arXiv preprint arXiv:2209.14390*, 2022. **1, 3**
- [4] Mahmoud Assran, Nicolas Loizou, Nicolas Ballas, and Mike Rabbat. Stochastic gradient push for distributed deep learning. In *International Conference on Machine Learning*, pages 344–353. PMLR, 2019. **3**
- [5] Aditya Balu, Zhanhong Jiang, Sin Yong Tan, Chinmay Hedge, Young M Lee, and Soumik Sarkar. Decentralized deep learning using momentum-accelerated consensus. In *ICASSP 2021-2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 3675–3679. IEEE, 2021. **3**
- [6] Frédéric Berdoz, Abhishek Singh, Martin Jaggi, and Ramesh Raskar. Scalable collaborative learning via representation sharing. *arXiv preprint arXiv:2211.10943*, 2022. **1**
- [7] Léon Bottou. Large-scale machine learning with stochastic gradient descent. In *Proceedings of COMPSTAT'2010*, pages 177–186. Springer, 2010. **3**
- [8] Jeffrey Dean, Greg Corrado, Rajat Monga, Kai Chen, Matthieu Devin, Mark Mao, Marc’aurilio Ranzato, Andrew Senior, Paul Tucker, Ke Yang, et al. Large scale distributed deep networks. *Advances in neural information processing systems*, 25, 2012. **3**
- [9] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee, 2009. **5, 11**
- [10] Yasaman Esfandiari, Sin Yong Tan, Zhanhong Jiang, Aditya Balu, Ethan Herron, Chinmay Hegde, and Soumik Sarkar. Cross-gradient aggregation for decentralized learning from non-iid data. In *International Conference on Machine Learning*, pages 3036–3046. PMLR, 2021. **1, 3**
- [11] Arya Ketabchi Haghghat, Varsha Ravichandra-Mouli, Pranamesh Chakraborty, Yasaman Esfandiari, Saeed Arabi, and Anuj Sharma. Applications of deep learning in intelligent transportation systems. *Journal of Big Data Analytics in Transportation*, 2(2):115–145, 2020. **1**
- [12] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016. **5, 11**
- [13] Kevin Hsieh, Amar Phanishayee, Onur Mutlu, and Phillip Gibbons. The non-IID data quagmire of decentralized machine learning. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 4387–4398. PMLR, 13–18 Jul 2020. **1**
- [14] Hamel Husain. Imagenette - a subset of 10 easily classified classes from the imagenet dataset. <https://github.com/fastai/imagenette>, 2018. **5, 11**
- [15] Anastasiia Koloskova, Tao Lin, and Sebastian U Stich. An improved analysis of gradient tracking for decentralized machine learning. *Advances in Neural Information Processing Systems*, 34:11422–11435, 2021. **1, 3**
- [16] Anastasia Koloskova, Nicolas Loizou, Sadra Boreiri, Martin Jaggi, and Sebastian Stich. A unified theory of decentralized sgd with changing topology and local updates. In *International Conference on Machine Learning*, pages 5381–5393. PMLR, 2020. **3**
- [17] Jakub Konečný, H Brendan McMahan, Daniel Ramage, and Peter Richtárik. Federated optimization: Distributed machine learning for on-device intelligence. 2016. **1**
- [18] Alex Krizhevsky, Vinod Nair, and Geoffrey Hinton. Cifar (canadian institute for advanced research). <http://www.cs.toronto.edu/~kriz/cifar.html>, 2014. **5, 11**
- [19] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick Haffner. Gradient-based learning applied to document recognition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998. **5, 11**
- [20] Chengxi Li, Gang Li, and Pramod K Varshney. Decentralized federated learning via mutual knowledge transfer. *IEEE Internet of Things Journal*, 9(2):1136–1147, 2021. **4**
- [21] Daliang Li and Junpu Wang. Fedmd: Heterogenous federated learning via model distillation. *arXiv preprint arXiv:1910.03581*, 2019. **1**
- [22] Xiangru Lian, Ce Zhang, Huan Zhang, Cho-Jui Hsieh, Wei Zhang, and Ji Liu. Can decentralized algorithms outperform centralized algorithms? a case study for decentralized parallel stochastic gradient descent. *Advances in Neural Information Processing Systems*, 30, 2017. **1, 2, 3, 5, 6**
- [23] Tao Lin, Sai Praneeth Karimireddy, Sebastian Stich, and Martin Jaggi. Quasi-global momentum: Accelerating decentralized deep learning on heterogeneous data. In *Proceedings*

- of the 38th International Conference on Machine Learning, volume 139 of *Proceedings of Machine Learning Research*, pages 6654–6665. PMLR, 18–24 Jul 2021. 1, 3, 4, 5, 6
- [24] Tao Lin, Lingjing Kong, Sebastian U Stich, and Martin Jaggi. Ensemble distillation for robust model fusion in federated learning. *Advances in Neural Information Processing Systems*, 33:2351–2363, 2020. 1
- [25] Hanxiao Liu, Andy Brock, Karen Simonyan, and Quoc Le. Evolving normalization-activation layers. In H. Larochelle, M. Ranzato, R. Hadsell, M. F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 13539–13550. Curran Associates, Inc., 2020. 5
- [26] Angelia Nedic. Distributed gradient methods for convex machine learning problems in networks: Distributed optimization. *IEEE Signal Processing Magazine*, 37(3):92–101, 2020. 1
- [27] Deepak Ravikumar, Gobinda Saha, Sai Aparna Aketi, and Kaushik Roy. Homogenizing non-iid datasets via in-distribution knowledge distillation for decentralized learning. *arXiv preprint arXiv:2304.04326*, 2023. 1, 4
- [28] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. Mobilenetv2: Inverted residuals and linear bottlenecks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 4510–4520, 2018. 5, 11
- [29] Yuki Takezawa, Han Bao, Kenta Niwa, Ryoma Sato, and Makoto Yamada. Momentum tracking: Momentum acceleration for decentralized deep learning on heterogeneous data. *arXiv preprint arXiv:2209.15505*, 2022. 1, 3
- [30] Thijs Vogels, Lie He, Anastasiia Koloskova, Sai Praneeth Karimireddy, Tao Lin, Sebastian U Stich, and Martin Jaggi. Relaysum for decentralized deep learning on heterogeneous data. In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 28004–28015. Curran Associates, Inc., 2021. 1, 3, 4, 5, 6
- [31] Han Xiao, Kashif Rasul, and Roland Vollgraf. Fashion-mnist: a novel image dataset for benchmarking machine learning algorithms. *arXiv preprint arXiv:1708.07747*, 2017. 5, 11
- [32] Lin Xiao and Stephen Boyd. Fast linear iterations for distributed averaging. *Systems & Control Letters*, 53(1):65–78, 2004. 3
- [33] Zhaorong Zhang, Kan Xie, Qianqian Cai, and Minyue Fu. A bp-like distributed algorithm for weighted average consensus. In *2019 12th Asian Control Conference (ASCC)*, pages 728–733. IEEE, 2019. 3
- [34] Zhuangdi Zhu, Junyuan Hong, and Jiayu Zhou. Data-free knowledge distillation for heterogeneous federated learning. In *International Conference on Machine Learning*, pages 12878–12889. PMLR, 2021. 1

## A. Appendix

### A.1. Experimental Setup Details

For the decentralized setup, we use an undirected ring, undirected Dyck graph, and undirected torus graph topologies with a uniform mixing matrix. The undirected ring topology for any graph size has 3 peers per agent including itself and each edge has a weight of  $\frac{1}{3}$ . The undirected Dyck topology with 32 agents has 4 peers per agent including itself and each edge has a weight of  $\frac{1}{4}$ . The undirected torus topology with 32 agents has 5 peers per agent including itself and each edge has a weight of  $\frac{1}{5}$ . All our experiments were conducted on a system with an NVIDIA A40 card with 4 GPUs. We report the test accuracy of the consensus model averaged over three randomly chosen seeds. The consensus model is obtained by averaging the model parameters across all agents using an all-reduce mechanism at the end of the training.

#### A.1.1 Datasets

In this section, we give a brief description of the datasets used in our experiments. We use a diverse set of datasets each originating from a different distribution of images to show the generalizability of the proposed techniques.

**CIFAR-10:** CIFAR-10 [18] is an image classification dataset with 10 classes. The image samples are colored (3 input channels) and have a resolution of  $32 \times 32$ . There are 50,000 training samples with 5000 samples per class and 10,000 test samples with 1000 samples per class.

**CIFAR-100:** CIFAR-100 [18] is an image classification dataset with 100 classes. The image samples are colored (3 input channels) and have a resolution of  $32 \times 32$ . There are 50,000 training samples with 500 samples per class and 10,000 test samples with 100 samples per class. CIFAR-100 classification is a harder task compared to CIFAR-10 as it has 100 classes with very few samples per class to learn from.

**Fashion MNIST:** Fashion MNIST [31] is an image classification dataset with 10 classes. The image samples are in grayscale (1 input channel) and have a resolution of  $28 \times 28$ . There are 60,000 training samples with 6000 samples per class and 10,000 test samples with 1000 samples per class.

**Imagenette:** Imagenette [14] is a 10-class subset of the ImageNet dataset. The image samples are colored (3 input channels) and have a resolution of  $224 \times 224$ . There are 9469 training samples with roughly 950 samples per class and 3925 test samples.

**ImageNet:** ImageNet dataset [9] spans 1000 object classes and contains 1,281,167 training images, 50,000 validation images, and 100,000 test images. The image samples are colored (3 input channels) and have a resolution of  $224 \times 224$ .

#### A.1.2 Network Architecture

We replace ReLU+BatchNorm layers of all the model architectures with EvoNorm-S0 as it was shown to be better suited for decentralized learning over non-IID distributions.

**ResNet-20:** For ResNet-20 [12], we use the standard architecture with  $0.27M$  trainable parameters except that BatchNorm+ReLU layers are replaced by EvoNorm-S0.

**ResNet-18:** For ResNet-18 [12], we use the standard architecture with  $11M$  trainable parameters except that BatchNorm+ReLU layers are replaced by EvoNorm-S0.

**LeNet-5:** For LeNet-5 [19], we use the standard architecture with 61,706 trainable parameters.

**MobileNet-V2:** We use the the standard MobileNet-V2 [28] architecture used for CIFAR dataset with  $2.3M$  parameters except that BatchNorm+ReLU layers are replaced by EvoNorm-S0.

#### A.1.3 Hyper-parameters

This section presents a detailed description of the hyper-parameters used in our experiments. All the experiments were run for three randomly chosen seeds. We decay the step size by 10x after 50% and 75% of the training, unless mentioned otherwise. For all the experiments, we have used a momentum of 0.9 with Nesterov, a weight decay of 0.0001, and a mini-batch size of 32 per agent.

Table 9. The value of  $\lambda_m, \lambda_v$  used for training CIFAR-10 with non-IID data using ResNet-20 architecture presented in Table 1

Agents ( $n$ )	Method	ResNet-20	
		$\alpha = 0.1$	$\alpha = 0.01$
16	CCL (ours)	0.01, 0.0	0.01, 0.01
32	CCL (ours)	0.1, 0.1	0.1, 0.1

**Hyper-parameters for experiments in Table 1:** All the experiments have the stopping criteria set to 200 epochs. The initial learning rate is set to 0.1. We decay the step size by  $10 \times$  in multiple steps at  $100^{th}$  and  $150^{th}$  epoch. Table 9 presents values of the scaling factor  $\lambda_m, \lambda_d$  used in the experiments.

Table 10. The value of  $\lambda_m, \lambda_v$  used for training various datasets with CCL (presented in Table 2).

Dataset	$\alpha = 0.1$	$\alpha = 0.01$
Fashion MNIST	0.001, 0.001	0.01, 0.01
CIFAR-100	0.1, 0.1	0.1, 0.1
Imagenette	0.001, 0.001	1.0, 1.0

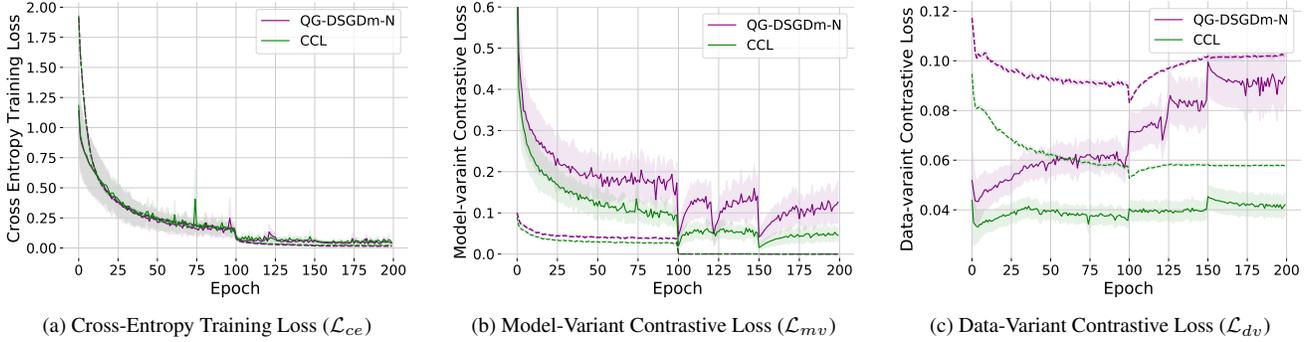


Figure 6. Comparing various training loss terms for IID (dashed lines) and non-IID (solid lines) partitions of CIFAR-10 trained on ResNet-20 over a ring topology of 16 agents. We use  $\alpha = 10$  for IID data and  $\alpha = 0.01$  for non-IID data.

**Hyper-parameters for experiments in Table 2:** All the experiments for CIFAR-100 and ImageNette have the stopping criteria set to 100 epochs and Fashion MNIST experiments have a stopping criteria of 50 epochs. The initial learning rate is set to 0.1 for CIFAR-100 and 0.01 for Fashion MNIST and Imagenette. Table 10 presents values of the scaling factor  $\lambda_m, \lambda_d$  used in the experiments.

**Hyper-parameters for experiments in Table 3:** All the experiments have the stopping criteria set to 200 epochs. The initial learning rate is set to 0.1. We decay the step size by  $10\times$  in multiple steps at  $100^{th}$  and  $150^{th}$  epoch. Table 11 presents values of the scaling factor  $\lambda_m, \lambda_d$  used in the experiments. All the experiments on the Dyck and Torus graph use an averaging rate of 0.9 (instead of the default value of 1.0).

Table 11. The value of  $\lambda_m, \lambda_v$  used for training CIFAR-10 datasets with CCL on ResNet-20 over various graph topologies (presented in Table 3).

Graph	$\alpha = 0.1$	$\alpha = 0.01$
Dyck (32 agents)	0.1, 0.1	0.1, 0.1
Torus (32 agents)	0.1, 0.1	0.1, 0.1

## A.2. Additional Results

Figure 6 measures the different training losses for both IID and non-IID distribution with  $\alpha = 0.01$  of the CIFAR-10 dataset trained on ResNet-20 architecture. We observe that the training cross-entropy loss (Fig. 6a) for IID and non-IID data converges to zero even though there is a huge gap in the validation loss. However, Fig. 6b shows that the model-variation contrastive loss for the baseline is much higher in non-IID settings compared to IID and hence is a good measure of data-heterogeneity. On the other hand, data-variant contrastive loss measures the variation in class representations across agents. Fig. 6c shows

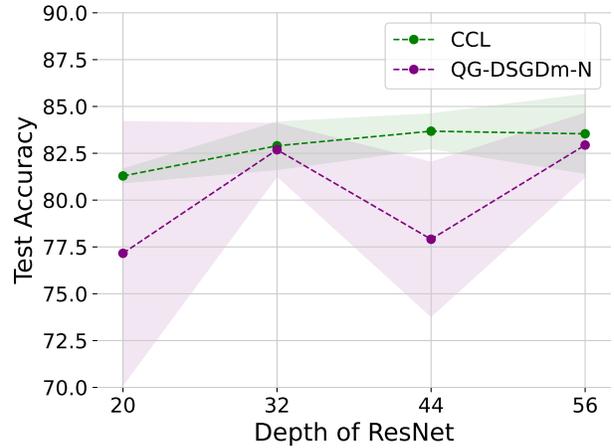


Figure 7. Test accuracy for the CIFAR-10 dataset trained on ResNet architecture with varying depth over 16-agent ring topology with a skew of  $\alpha = 0.01$ .

that this variation is relatively stable throughout the training process for QG-DSGDm-N (baseline) with IID Data. However, for QG-DSGDm-N with the non-IID setting, a significant increase in the variation of class representations across agents is evident. The proposed CCL framework explicitly minimizes the model-variant and data-variant contrastive loss. Fig. 6b shows that the CCL helps in reducing the model variance compared to QG-DSGDm-N. Fig. 6c shows that CCL has a stable variation in class representations across agents compared to QG-DSGDm-N. This results in better performance of the proposed *Cross-feature Contrastive Loss* for decentralized learning on heterogeneous data. Further, we evaluate the proposed CCL on the varying depth of ResNet architecture with ring topology of 16 agents as shown in Figure 7. We observe that the proposed CCL framework consistently outperforms the QG-DSGDm-N baseline over varying graph sizes by an average improvement of 2.68%.