

DART: A Server-side Plug-in for Resource-efficient Robust Federated Learning

Supplementary Material

7. Method Details

7.1. Data Augmentation

The DART method adopts the standard AugMix [15] augmentation strategy, which stochastically mixes diverse transformation operations such as posterize, equalize, rotate, translate, and shear. This procedure introduces multiple sources of randomness, including the selection of operations, their severity levels, the depth of transformation chains, and mixing weights.

First, three operations, $op_1, op_2,$ and $op_3,$ are sampled from the set of available operations, \mathcal{O} . The sampled operations are sequentially composed to generate transformation chains of increasing depth:

$$op_{12} = op_2 \circ op_1, \text{ and } op_{123} = op_3 \circ op_2 \circ op_1 \quad (5)$$

Then one chain is sampled uniformly at random from $\{op_1, op_{12}, op_{123}\}$. This process is repeated S times, where S is a hyperparameter denoting mixture width, resulting in S operation chains $\{\text{chain}_i\}_{i=1}^S$. S is set to 3 by default [15].

The augmented sample \mathbf{x}_{aug} is generated from a sample \mathbf{x} using:

$$\mathbf{x}_{\text{aug}} = \eta \mathbf{x} + (1 - \eta) \sum_{i=1}^S m_i \times \text{chain}_i(\mathbf{x}) \quad (6)$$

where $(m_1 \dots m_i)$ and η are random weights sampled from Dirichlet($\alpha \dots \alpha$) and from Beta(α, α) distributions, respectively. As can be seen, the augmentation procedure mixes the output of the operation chains. It also includes direct contribution from the original image \mathbf{x} to prevent the loss of semantic information.

This diverse augmentation process is designed to improve the robustness of machine learning models by exposing them to a wide range of data during training.

7.2. Overall System Algorithm

Section 4.1 presented the overall DART-enhanced FL system design. In a standard FL setting, clients perform local training while a central server aggregates their updates. Our proposed plug-in, **DART**, is integrated on the server side, operating immediately after model aggregation and before the updated global model is redistributed to the clients. Algorithm 1 details the DART-enhanced FL system.

7.3. DART Algorithm

Section 4.2 introduced our novel server-side plug-in, **DART**, designed to enable resource-efficient and robust

Algorithm 1 System Description.

K : number of clients; D_k : local private dataset of client k ; D_0 : public unlabeled dataset; G : number of global updates; E : local rounds per global update; T_{DART} : global rounds per DART update

Server executes:

```
1: Initialize  $\mathbf{w}_0$ ;  
2: for each global round  $t = 1, 2, \dots, G$  do  
3:   for each client  $k = 1, 2, \dots, K$  in parallel do  
4:      $\mathbf{w}_k^{t+1} \leftarrow \text{ClientUpdate}(k, \mathbf{w}_k^t)$ ;  
5:   end for  
6:    $\mathbf{w}_0^{t+1} \leftarrow \text{Aggregate}(\mathbf{w}_1^{t+1}, \dots, \mathbf{w}_K^{t+1})$ ;  
7:   if  $t \bmod T_{\text{DART}} = 0$  then  
8:      $\mathbf{w}_{\text{rob}}^{t+1} \leftarrow \text{DART}(\mathbf{w}_0^{t+1}, D_0)$ ;  
9:      $\mathbf{w}_0^{t+1} \leftarrow \mathbf{w}_{\text{rob}}^{t+1}$ ;  
10:  end if  
11: end for
```

ClientUpdate(k, \mathbf{w}):

```
1: for each local epoch  $i = 1, 2, \dots, E$  do  
2:   for batch  $\mathbf{b} \in D_k$  do  
3:      $\mathbf{w} \leftarrow \text{LocalTraining}(\mathbf{w}, \mathbf{b})$ ;  
4:   end for  
5: end for  
6: return  $\mathbf{w}$  to server
```

federated learning. Algorithm 2 details the DART procedure.

8. Theorem Proof

The theoretical analysis in this work is inspired by Lin et al. [29], who study ensemble distillation for model fusion in model heterogeneous FL. In their work, a generalization bound is developed that relates the loss of distilled federated models to that of a centrally trained model, explicitly characterizing the impact of distribution mismatch across clients and server. Building on this foundation, we derive a bound that relates the pre-DART (teacher) and post-DART (student) clean risks, enabling us to formally quantify the utility degradation introduced by DART while enhancing robustness.

Setup and notation. Let \mathcal{X} be the input space and \mathcal{Y} the label set. Let \mathcal{D}_{in} and \mathcal{D}_{out} be distributions on $\mathcal{X} \times \mathcal{Y}$. A classifier $h : \mathcal{X} \rightarrow \mathcal{Y}$ has risk on a distribution Q

$$R_Q(h) = \Pr_{(\mathbf{x}, \mathbf{y}) \sim Q} (h(\mathbf{x}) \neq \mathbf{y}) \quad (7)$$

Algorithm 2 DART.

D_{tr} and D_{val} : public unlabeled training and validation datasets partitioned from D_0 ; T_{max} : maximum number of DART epochs; T_{val} : early stopping threshold; η_s : DART learning rate; \mathbf{w}_0 : pre-trained weights.

```

1: counter  $\leftarrow$  0;  $\mathbf{w} \leftarrow \mathbf{w}_0$ ;  $\text{loss}_{\min} \leftarrow \infty$ 
2: for each epoch  $i = 1, 2, \dots, T_{\text{max}}$  do
3:   for batch  $\mathbf{b} \in D_{\text{tr}}$  do
4:      $\mathbf{w} \leftarrow \mathbf{w} - \eta_s \nabla \mathcal{L}_{\text{DART}}(\mathbf{w}; \mathbf{b})$ 
5:   end for
6:   if  $\mathcal{L}_{\text{DART}}(\mathbf{w}; D_{\text{val}}) \leq \text{loss}_{\min}$  then
7:      $\text{loss}_{\min} \leftarrow \mathcal{L}_{\text{DART}}(\mathbf{w}; D_{\text{val}})$ 
8:      $\mathbf{w}_{\text{best}} \leftarrow \mathbf{w}$ ; counter  $\leftarrow$  0
9:   else
10:    counter ++
11:   end if
12:   if counter =  $T_{\text{val}} - 1$  then
13:     break
14:   end if
15: end for
16:  $\mathbf{w}_{\text{rob}} \leftarrow \mathbf{w}_{\text{best}}$ 
17: return  $\mathbf{w}_{\text{rob}}$ 

```

For $h, h' : \mathcal{X} \rightarrow \mathcal{Y}$, define their *disagreement probability* on a distribution Q by

$$\text{Dis}_Q(h, h') = \Pr_{\mathbf{x} \sim Q} (h(\mathbf{x}) \neq h'(\mathbf{x})) \quad (8)$$

We consider a teacher $f_{\mathbf{w}_t}$ and a student $f_{\mathbf{w}_s}$ that make predictions via posteriors $p_t(\cdot | \mathbf{x})$ and $p_s(\cdot | \mathbf{x})$, with decisions $f_{\mathbf{w}_t}(\mathbf{x}) = \arg \max_{\mathbf{y}} p_t(\mathbf{y} | \mathbf{x})$ and $f_{\mathbf{w}_s}(\mathbf{x}) = \arg \max_{\mathbf{y}} p_s(\mathbf{y} | \mathbf{x})$. The client distribution is \mathcal{D}_{in} , while \mathcal{D}_{out} is the server distribution.

Define the *teacher margin* at \mathbf{x} :

$$m_t(\mathbf{x}) = p_t(\mathbf{y}_t^* | \mathbf{x}) - \max_{\mathbf{y} \neq \mathbf{y}_t^*} p_t(\mathbf{y} | \mathbf{x}) \quad (9)$$

where $\mathbf{y}_t^* = \arg \max_{\mathbf{y}} p_t(\mathbf{y} | \mathbf{x})$ and $m_t(\mathbf{x}) \in [0, 1]$. Define the *margin tail* on \mathcal{D}_{out} by $\theta_{\mathcal{D}_{\text{out}}}(\gamma) = \Pr_{\mathbf{x} \sim \mathcal{D}_{\text{out}}} (m_t(\mathbf{x}) \leq \gamma)$ for $\gamma \in (0, 1]$. The total variation distance is

$$\text{TV}(\mathbf{x}) = \frac{1}{2} \|p_t(\cdot | \mathbf{x}) - p_s(\cdot | \mathbf{x})\|_1 \quad (10)$$

and the expected distillation KL divergence on \mathcal{D}_{out} :

$$\delta = \mathbb{E}_{\mathbf{x} \sim \mathcal{D}_{\text{out}}} [\text{KL}(p_t(\cdot | \mathbf{x}) \| p_s(\cdot | \mathbf{x}))] \quad (11)$$

Finally, the $\mathcal{H}\Delta\mathcal{H}$ -divergence [3] between \mathcal{D}_{in} and \mathcal{D}_{out} is

$$d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{\text{in}}, \mathcal{D}_{\text{out}}) = 2 \sup_{h, h' \in \mathcal{H}} |\text{Dis}_{\mathcal{D}_{\text{in}}}(h, h') - \text{Dis}_{\mathcal{D}_{\text{out}}}(h, h')|$$

for some hypothesis class \mathcal{H} .

Lemma 1. For any $\mathbf{x} \in \mathcal{X}$, if $f_{\mathbf{w}_s} \neq f_{\mathbf{w}_t}$, then

$$m_t(\mathbf{x}) \leq 2 \text{TV}(\mathbf{x})$$

Proof. Let $\mathbf{y}_t^* = f_{\mathbf{w}_t}(\mathbf{x})$ and let $\mathbf{y}' = f_{\mathbf{w}_s}(\mathbf{x}) \neq \mathbf{y}_t^*$. Since the student prefers \mathbf{y}' , $p_s(\mathbf{y}' | \mathbf{x}) \geq p_s(\mathbf{y}_t^* | \mathbf{x})$. Let $\Delta_{\mathbf{y}} := p_s(\mathbf{y} | \mathbf{x}) - p_t(\mathbf{y} | \mathbf{x})$ for each class \mathbf{y} , then

$$m_t(\mathbf{x}) \leq p_t(\mathbf{y}_t^* | \mathbf{x}) - p_t(\mathbf{y}' | \mathbf{x}) \quad (12)$$

$$= (p_t(\mathbf{y}_t^* | \mathbf{x}) - p_s(\mathbf{y}_t^* | \mathbf{x})) \quad (13)$$

$$+ (p_s(\mathbf{y}_t^* | \mathbf{x}) - p_s(\mathbf{y}' | \mathbf{x})) \\ + (p_s(\mathbf{y}' | \mathbf{x}) - p_t(\mathbf{y}' | \mathbf{x}))$$

$$\leq |\Delta_{\mathbf{y}_t^*}| + 0 + |\Delta_{\mathbf{y}'}| \quad (14)$$

$$\leq \sum_{\mathbf{y} \in \mathcal{Y}} |\Delta_{\mathbf{y}}| \quad (15)$$

$$= \|p_s(\cdot | \mathbf{x}) - p_t(\cdot | \mathbf{x})\|_1 \quad (16)$$

$$= 2 \text{TV}(\mathbf{x}) \quad (17)$$

which proves the claim. \square

Theorem (Clean Student Risk Bound under Distillation).

For any $\gamma \in (0, 1]$,

$$\mathcal{R}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_s}) \leq \mathcal{R}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_t}) + \theta_{\mathcal{D}_{\text{out}}}(\gamma) \\ + \frac{2}{\gamma} \sqrt{\frac{\delta}{2}} + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{\text{in}}, \mathcal{D}_{\text{out}}) \quad (18)$$

Proof. **(1) Risk decomposition.** For any (\mathbf{x}, \mathbf{y}) ,

$$\mathbf{1}[f_{\mathbf{w}_s}(\mathbf{x}) \neq \mathbf{y}] \leq \mathbf{1}[f_{\mathbf{w}_t}(\mathbf{x}) \neq \mathbf{y}] + \mathbf{1}[f_{\mathbf{w}_s}(\mathbf{x}) \neq f_{\mathbf{w}_t}(\mathbf{x})] \quad (19)$$

Taking expectation over $(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_{\text{in}}$ gives

$$\mathcal{R}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_s}) \leq \mathcal{R}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_t}) + \text{Dis}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_s}, f_{\mathbf{w}_t}) \quad (20)$$

(2) Move disagreement from \mathcal{D}_{in} to \mathcal{D}_{out} . By the definition of $d_{\mathcal{H}\Delta\mathcal{H}}$,

$$\text{Dis}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_s}, f_{\mathbf{w}_t}) \leq \text{Dis}_{\mathcal{D}_{\text{out}}}(f_{\mathbf{w}_s}, f_{\mathbf{w}_t}) + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{\text{in}}, \mathcal{D}_{\text{out}}) \quad (21)$$

Combining with (20),

$$\mathcal{R}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_s}) \leq \mathcal{R}_{\mathcal{D}_{\text{in}}}(f_{\mathbf{w}_t}) + \text{Dis}_{\mathcal{D}_{\text{out}}}(f_{\mathbf{w}_s}, f_{\mathbf{w}_t}) \\ + \frac{1}{2} d_{\mathcal{H}\Delta\mathcal{H}}(\mathcal{D}_{\text{in}}, \mathcal{D}_{\text{out}}) \quad (22)$$

(3) Control disagreement on \mathcal{D}_{out} via margin and TV. By Lemma 1, for any \mathbf{x} ,

$$\{f_{\mathbf{w}_s}(\mathbf{x}) \neq f_{\mathbf{w}_t}(\mathbf{x})\} \subseteq \{m_t(\mathbf{x}) \leq 2 \text{TV}(\mathbf{x})\}$$

Fix $\gamma \in (0, 1]$. Then $\{m_t \leq 2 \text{TV}\} \subseteq \{m_t \leq \gamma\} \cup \{2 \text{TV} \geq \gamma\}$ and,

$$\text{Dis}_{\mathcal{D}_{\text{out}}}(f_{\mathbf{w}_s}, f_{\mathbf{w}_t}) \leq \theta_{\mathcal{D}_{\text{out}}}(\gamma) + \Pr_{\mathbf{x} \sim \mathcal{D}_{\text{out}}} (2 \text{TV}(\mathbf{x}) \geq \gamma)$$

Now apply Markov’s inequality (since $\text{TV} \geq 0$), Pinsker and Jensen:

$$\Pr(2 \text{TV} \geq \gamma) \leq \frac{2}{\gamma} \mathbb{E}[\text{TV}] \quad (23)$$

$$\leq \frac{2}{\gamma} \mathbb{E} \left[\sqrt{\frac{1}{2} \text{KL}(p_t \| p_s)} \right] \quad (24)$$

$$\leq \frac{2}{\gamma} \sqrt{\frac{1}{2} \mathbb{E}[\text{KL}(p_t \| p_s)]} = \frac{2}{\gamma} \sqrt{\frac{\delta}{2}} \quad (25)$$

Hence

$$\text{Dis}_{\mathcal{D}_{\text{out}}}(f_{\mathbf{w}_s}, f_{\mathbf{w}_t}) \leq \theta_{\mathcal{D}_{\text{out}}}(\gamma) + \frac{2}{\gamma} \sqrt{\frac{\delta}{2}} \quad (26)$$

(4) Combine. Insert (26) into (22) to obtain the stated bound. \square

9. Resource Setup

9.1. Profiling Client Hardware

Each client is represented by an NVIDIA Jetson Orin Nano [33], equipped with an Ampere GPU, a 6-core ARM CPU, and 8 GB of RAM, delivering up to 67 TOPS of compute. We profile batch execution time and batch energy consumption directly on the hardware. Power is measured using Jetson Stats [4], and energy is computed by integrating power over time, approximated as the product of sampled power values and their corresponding time intervals. Execution time is measured as wall-clock time. Table 6 reports the measured batch time and energy. While our experiments are executed on NVIDIA L40S GPUs, we use measurements collected on the Jetson Orin Nano to estimate the per-client training time and energy consumption for large-scale federated learning.

Table 6. Client-side batch time and batch energy, measured on NVIDIA Jetson Orin Nano using ResNet-18 with batch size 64 and input dimension $32 \times 32 \times 3$.

	Clean FL		Robust FL		
	FedAvg	FedProx	FedAFA	FedAugMix	FedPrime
Batch Time (s)	0.39	0.41	0.75	1.02	8.28
Batch Energy (J)	1.43	1.53	2.76	4.15	8.20

9.2. Time Model

Since multiple clients perform training in parallel during each global round, the training time at round i can be estimated as

$$\mathcal{T}_{\text{tr}}^i = i \times t_{\text{batch}} \times b_{\text{avg}} \quad (27)$$

where t_{batch} is the profiled batch execution time and b_{avg} is the average number of batches processed by a client. b_{avg} is computed as the average number of samples per client divided by the batch size:

$$b_{\text{avg}} = \frac{\sum_{k=1}^K |D_k|}{K \times \text{batch.size}} \quad (28)$$

where K is the number of clients and $|D_k|$ is denotes the number of samples held by client k .

9.3. Energy Model

We consider the total compute energy consumed by clients during federated training, which can be estimated for global round i as

$$\mathcal{E}_{\text{tr}}^i = i \times r \times K \times e_{\text{batch}} \times b_{\text{avg}} \quad (29)$$

where e_{batch} denotes the profiled per-batch energy consumption, r denotes the client participation rate, K is the number of clients, and b_{avg} is computed using Eq. 28.

10. CIFAR-10-C Results

10.1. CIFAR-10-C Corruptions

CIFAR-10-C [14] is a standard benchmark for evaluating model robustness to common image corruptions. Table 7 lists the corruption types included in the benchmark along with the abbreviations used throughout the paper.

10.2. Full CIFAR-10-C Results

Section 5.2 evaluates popular FL methods and their DART-enhanced variants on the CIFAR-10-C robustness benchmark. Table 8 extends the results by reporting performance across all corruption types. As discussed earlier, integrating DART consistently improves robustness on CIFAR-10-C at zero client-side overhead and at varying degrees of heterogeneity. Notably, DART improves performance across most corruption categories, with the exception of brightness and saturation, which introduce minimal visual distortion and thus exhibit smaller robustness gains.

11. Ablation Studies

11.1. Generalization Across Model Architectures

To demonstrate DART on various deep learning models, we conduct simulations for a 10 client FL setup with a fixed time budget of 1000s and we compare DART enhanced FedAvg clean, robust and average accuracy against FedAvg and FedAugMix for ResNet-18, MobileNet and VGG-16. Table 9 shows that DART enhances the robust accuracy of

Table 7. CIFAR-10-C corruption types and abbreviations used throughout the paper.

Corruption type	Abbreviation
Brightness	bright
Contrast	contr
Defocus blur	defoc
Elastic transform	elas
Fog	fog
Frost	frost
Gaussian blur	g_blur
Gaussian noise	g_noise
Glass blur	glass
Impulse noise	imp
JPEG compression	jpeg
Motion blur	motion
Pixelate	pix
Saturate	sat
Shot noise	shot
Snow	snow
Spatter	spatt
Speckle noise	speck
Zoom blur	zoom

FedAvg by 4.3% for a slight drop of 1.7% in clean accuracy on average. Additionally, DART enhanced FedAvg outperforms FedAugMix in both clean and robust accuracy by an average of 18.9% and 13.8%, respectively. Most notably, DART performs best with VGG-16 where the gains are most significant (29.2% in \mathcal{A}_{cln} , 24.3% in \mathcal{A}_{rob}). These results indicate that DART-enhanced FL generalizes consistently across model architectures while maintaining a favorable balance between clean and robust accuracy, making it well-suited for deployment in resource-constrained environments.

Fig. 4 shows that DART-enhanced FedAvg consistently outperforms both FedAvg and FedAugMix in robust accuracy (\mathcal{A}_{rob}) under constrained time and energy budgets. Across the full range of time and energy, DART-enhanced FL achieves higher robustness than FedAvg. When the per-client time budget is below 2200s or the energy budget is below 25kJ, DART also surpasses FedAugMix, highlighting its effectiveness in resource-constrained settings. Moreover, DART-enhanced FL reaches a target robust accuracy of 80% on ResNet-18 with the lowest time and energy cost among all methods. In contrast, FedAugMix requires substantially more computation to converge, limiting its practicality for deployment on resource-constrained devices. All methods are evaluated with 10 clients over 200 global rounds, with DART applied to the FedAvg model at each round. Similar trends are observed for MobileNet (Fig. 5) and VGG-16 (Fig. 6), demonstrating DART’s generalization across model architectures.

11.2. CIFAR-100 Results

To showcase the benefit of using DART with different client datasets, we provide the following results that utilize CIFAR-100 on the clients and CIFAR-10 on the server. The number of clients used is 10.

Analyzing the CIFAR-100 training curves of ResNet-18 in Fig. 7, we notice that DART enhanced FedAvg achieves the best robust accuracy when time and energy are below 2500s and 30kJ, respectively. Within these intervals, DART also significantly outperforms FedAugMix in clean accuracy, while incurring a drop relative to FedAvg. As expected, FedAugMix requires substantially more time and energy to converge. DART provides the most robust model under resource constraints, even when CIFAR-100 is the client dataset.

These experiments demonstrate that the advantages of DART are not limited to a specific client dataset. By maintaining strong clean and robust performance across diverse data distributions, DART proves to be broadly applicable and adaptable, reinforcing its potential as a versatile solution for robust federated learning in scenarios constrained by compute, memory, time, and energy.

11.3. Impact of Robustification Period

We experimentally investigate the impact of the robustness-enhancement period T_{DART} on utility and robustness. Additional DART iterations can be executed since the server is not resource-limited and this imposes no further cost on clients.

Fig. 8 shows that reducing the robustification period T_{DART} has minimal impact on clean accuracy. In general, smaller values of T_{DART} , i.e., more frequent DART updates, leads to higher robust accuracy at no additional client training time and energy. For example, ResNet-18 robust accuracy improves from 81.98% to 83.39% when T_{DART} decreases from 200 to 25. Similar trends can be observed for MobileNet and VGG-16, where robust accuracy increases by +1% and +0.72% respectively. These findings suggest that robust accuracy can be increased by performing more frequent DART updates, however, in cases where the server cannot be overloaded, applying DART once at the final global update yields competitive performance.

11.4. DART Parameter Sweeps

In DART, performance is sensitive to the loss-weighting parameter α and the mixture width S . Fig 9 reports the effect of varying α and S on robust accuracy. To maximize robustness, we observe that α values in the range of 8–10 yield the strongest performance. Increasing S introduces more complex data augmentation but also increases server-side computation. Based on this trade-off, we select $\alpha = 10$ and $S = 3$ to achieve strong robustness while limiting server overhead.

Table 8. Impact of DART in enhancing the robustness of diverse standard FL methods. Clean accuracy is measured on CIFAR-10 (\mathcal{A}_{cln}), and robust accuracy on CIFAR-10-C ($\mathcal{A}_{\text{rob}}^{\text{C}}$) for ResNet-18. Accuracies on all CIFAR-10-C corruptions are shown. DART significantly enhances the robustness (avg. 4.3%) of standard FL methods while incurring a small drop (avg. 1.6%) in clean accuracy at no additional computational cost for the clients.

Method	\mathcal{A}_{cln} (%)	$\mathcal{A}_{\text{rob}}^{\text{C}}$ (%)	bright (%)	contr (%)	defoc (%)	elas (%)	fog (%)	frost (%)	g_blur (%)	g_noise (%)	glass (%)	imp (%)	jpeg (%)	motion (%)	pix (%)	sat (%)	shot (%)	snow (%)	spatt (%)	speck (%)	zoom (%)
$\eta_{\text{FD}} = 100$																					
FedAvg	90.3	70.5	88.9	68.2	75.8	79.2	82.5	72.2	67.6	47.8	46.0	55.6	77.7	71.8	73.2	87.1	58.3	75.2	79.6	60.6	71.6
FedDyn	85.0	67.7	83.1	63.7	73.7	74.6	76.9	68.1	66.9	53.5	43.3	56.9	75.8	65.9	66.0	81.8	61.3	68.4	75.5	62.5	68.3
FedNova	86.5	69.5	84.5	61.9	75.0	76.6	76.4	69.1	68.4	56.4	45.6	58.9	77.9	68.4	71.3	83.5	63.1	70.9	78.6	63.7	70.2
FedProx	91.4	72.7	89.9	73.5	79.1	81.0	84.7	73.7	70.1	51.7	45.5	60.7	79.4	73.8	71.5	88.5	61.8	77.2	82.4	64.0	73.5
FedAvg+DART	89.2	77.4	87.4	73.8	84.8	81.2	82.5	75.2	82.6	64.0	57.7	69.2	79.2	80.1	77.8	86.2	71.7	77.2	82.8	73.8	82.5
FedDyn+DART	83.2	71.2	81.2	67.4	78.6	74.8	75.9	67.9	76.2	59.2	51.5	65.2	74.6	72.8	70.3	80.2	66.2	69.9	76.8	68.0	75.7
FedNova+DART	85.1	72.6	82.8	66.7	81.0	77.5	76.1	68.1	78.7	60.0	51.0	64.5	77.2	74.3	76.2	82.3	67.0	70.9	78.8	68.6	78.1
FedProx+DART	90.2	79.0	88.6	77.7	86.1	82.4	84.2	77.3	83.9	65.6	57.9	72.6	80.9	81.8	77.7	87.3	73.7	79.2	84.8	76.0	83.5
$\eta_{\text{FD}} = 10$																					
FedAvg	89.5	69.6	88.1	68.6	76.0	78.3	82.4	69.7	68.0	46.6	43.6	54.3	77.2	72.6	72.2	86.4	56.6	74.1	78.8	58.4	71.1
FedDyn	83.9	66.7	82.0	60.9	72.6	72.6	75.1	65.9	66.5	52.5	43.5	57.2	73.9	66.6	67.1	80.9	59.7	67.4	74.9	60.6	67.1
FedNova	85.7	67.9	83.6	60.1	73.6	75.1	74.7	67.1	66.5	54.2	44.6	57.3	76.8	65.6	70.0	82.7	66.1	69.2	77.7	61.7	68.4
FedProx	90.6	71.9	89.2	71.9	79.4	80.6	83.6	71.7	71.2	50.0	45.6	59.1	77.8	73.9	71.2	87.9	59.7	75.6	82.0	61.6	74.9
FedAvg+DART	88.3	75.8	86.4	71.8	83.4	81.2	81.6	73.2	81.1	62.1	54.7	67.4	78.0	78.9	77.6	85.7	70.0	75.2	81.2	72.1	81.1
FedDyn+DART	82.7	69.5	80.0	64.1	77.9	73.8	74.1	64.7	75.6	56.8	47.8	63.4	73.6	71.4	70.8	80.1	63.7	67.6	76.0	65.7	74.6
FedNova+DART	84.7	71.8	82.3	66.7	80.4	76.5	76.0	67.3	78.1	58.2	48.9	64.0	76.3	74.2	75.1	82.0	65.5	69.8	78.1	67.3	77.2
FedProx+DART	88.7	77.2	87.0	74.7	84.6	80.5	82.1	74.5	82.4	64.5	56.8	70.7	78.0	80.0	76.6	86.3	72.0	76.7	83.1	74.1	82.1
$\eta_{\text{FD}} = 1$																					
FedAvg	87.8	65.6	85.7	63.8	73.0	73.5	77.3	65.4	65.2	42.1	39.2	44.3	50.1	65.4	70.7	84.4	52.3	69.6	75.1	54.5	66.4
FedDyn	80.7	64.3	78.4	58.9	69.9	70.7	71.8	62.6	63.9	52.7	44.1	53.9	72.7	62.6	63.0	77.0	58.8	64.4	71.8	59.2	65.2
FedNova	82.7	64.9	80.3	55.3	69.3	70.4	69.0	64.3	62.8	54.1	45.6	57.3	74.2	60.1	68.3	79.5	60.2	66.4	74.9	60.4	63.4
FedProx	88.9	70.2	86.8	68.1	77.1	77.4	80.2	69.4	69.3	51.7	42.9	58.7	76.6	70.0	71.4	85.9	60.8	72.9	79.1	62.7	72.2
FedAvg+DART	83.2	68.4	80.1	63.9	76.1	70.5	72.7	64.8	72.8	56.2	46.6	60.8	71.6	67.8	73.0	79.7	63.0	67.6	75.1	64.8	71.8
FedDyn+DART	79.1	67.2	77.1	63.1	74.4	71.1	70.6	63.5	72.2	56.2	48.6	59.5	71.8	67.9	69.7	75.9	62.2	65.8	72.3	63.2	72.4
FedNova+DART	80.3	67.6	77.4	59.4	74.6	71.1	69.1	63.6	71.9	58.2	50.1	61.6	72.7	67.1	72.1	77.4	63.8	65.0	73.8	64.7	71.3
FedProx+DART	87.1	73.9	84.4	71.0	82.0	77.4	80.2	70.6	79.1	60.8	49.7	67.4	77.2	74.8	75.9	83.8	68.5	72.9	80.7	74.5	78.6
$\eta_{\text{FD}} = 0.1$																					
FedAvg	49.0	37.1	48.1	34.7	36.4	35.3	40.4	38.6	31.9	32.7	26.6	36.4	37.8	33.8	40.0	47.5	37.2	36.9	42.3	38.6	29.7
FedDyn	34.7	31.6	30.7	19.4	33.0	32.6	24.4	25.9	32.3	34.6	33.4	37.0	34.1	32.6	34.2	31.6	34.6	31.1	34.3	34.6	32.0
FedNova	54.3	45.0	50.5	26.1	44.0	44.3	34.7	43.4	40.6	48.4	43.9	45.2	52.5	37.9	52.1	51.1	50.2	48.8	51.6	50.3	40.1
FedProx	68.0	56.7	64.6	43.9	56.4	56.7	54.6	56.7	52.4	55.3	50.6	52.9	63.7	51.4	63.9	63.9	59.1	58.0	63.2	59.3	51.4
FedAvg+DART	43.0	32.9	41.2	29.9	36.9	33.3	35.2	30.1	34.8	26.7	20.1	29.6	34.4	34.3	36.2	41.6	30.2	30.3	36.2	31.3	33.0
FedDyn+DART	23.8	21.9	21.1	16.0	22.9	22.6	18.8	18.2	22.5	23.1	22.6	24.1	22.3	22.2	23.3	23.0	22.0	21.4	23.2	23.3	22.6
FedNova+DART	48.0	40.5	43.4	26.2	45.9	42.0	37.3	35.7	42.7	47.0	37.3	40.3	46.0	37.7	45.6	44.0	44.0	40.5	44.9	44.0	40.6
FedProx+DART	66.9	58.7	63.9	51.0	61.4	59.3	58.3	57.1	59.6	55.9	49.8	53.5	62.7	59.1	63.3	62.8	59.4	56.7	61.4	60.1	59.4

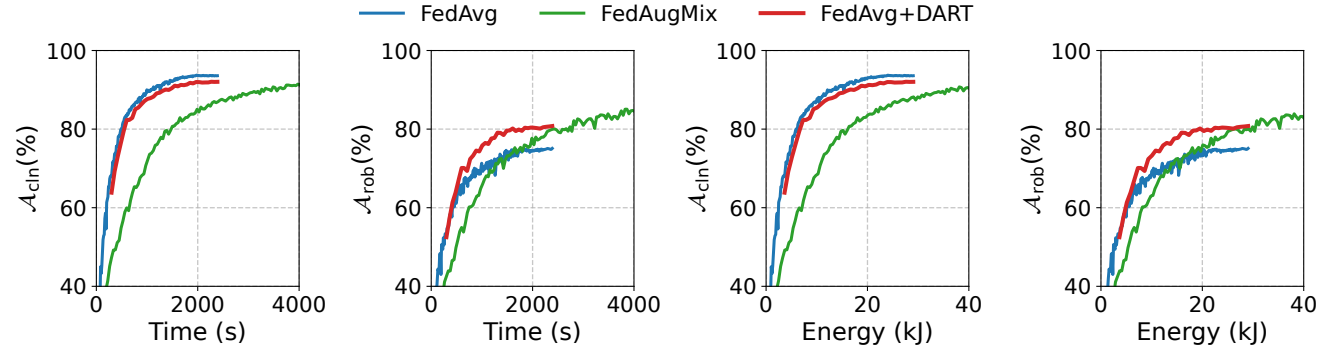


Figure 4. CIFAR-10 clean and CIFAR-10-C robust accuracy, \mathcal{A}_{cln} and \mathcal{A}_{rob} , under FedAvg, FedAugMix, and FedAvg+DART as a function of time and energy. The model used is ResNet-18, the server dataset is CIFAR-100, and 10 clients are deployed. FedAvg+DART takes the least energy and time to reach $\mathcal{A}_{\text{rob}} = 80\%$.

11.5. Loss Component Ablations

Table 10 presents an ablation study on the components \mathcal{L}_c and \mathcal{L}_d of the DART loss (1) as compared to standard clean training. Applying DART without the consistency loss \mathcal{L}_c results in a model with clean and robust accuracy comparable to the original pre-trained model. In contrast, removing the distillation loss \mathcal{L}_d leads to significant drops in both utility and robustness. Only when both loss components are present does DART maintain clean accuracy (with a minor

1.5% drop) while improving robustness by 5.7%.

These findings are consistent with those in Section 4: \mathcal{L}_d preserves clean accuracy, while \mathcal{L}_c promotes consistency, which translates to improved robust accuracy with minimal loss in clean accuracy, when both loss components are jointly optimized.

12. Using Synthetic Data on the Server

When high-quality server data is unavailable, generative models can be used to supply the server dataset for DART.

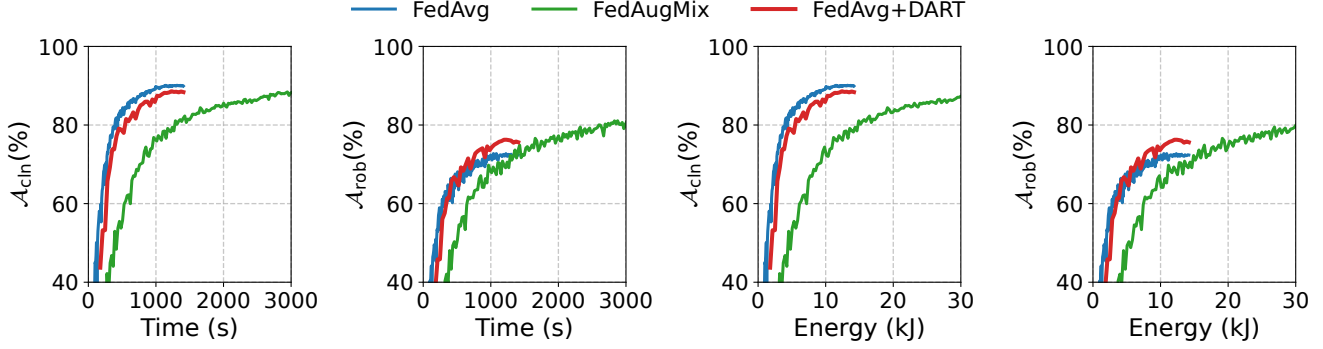


Figure 5. CIFAR-10 clean and CIFAR-10-C robust accuracy, \mathcal{A}_{cln} and \mathcal{A}_{rob} , under FedAvg, FedAugMix, and FedAvg+DART as a function of time and energy. The model used is MobileNet, the server dataset is CIFAR-100, and 10 clients are deployed.

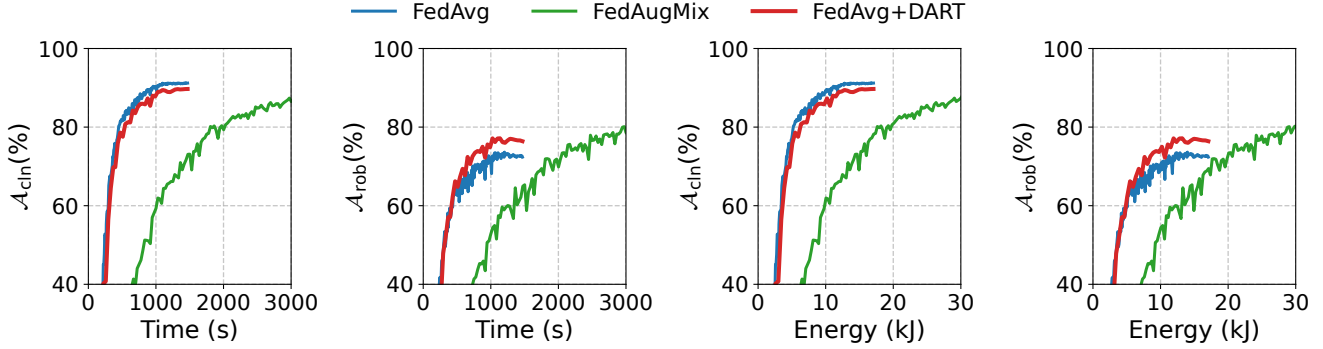


Figure 6. CIFAR-10 clean and CIFAR-10-C robust accuracy, \mathcal{A}_{cln} and \mathcal{A}_{rob} , under FedAvg, FedAugMix, and FedAvg+DART as a function of time and energy. The model used is VGG-16, the server dataset is CIFAR-100, and 10 clients are deployed.

Table 9. Accuracy at iso-time of $\approx 1000\text{s}$ for different models. The value of G is set to meet the time budget. FedAvg+DART achieves the highest robustness at zero client overhead.

Model	Method	G	\mathcal{A}_{cln} (%)	\mathcal{A}_{rob} (%)	\mathcal{A}_{avg} (%)
ResNet-18	FedAvg	83	90.00	71.32	80.66
	FedAugMix	27	72.22	65.97	69.10
	FedAvg+DART	83	87.62	75.85	81.74
MobileNet	FedAvg	141	89.83	72.38	81.11
	FedAugMix	43	76.30	68.41	72.36
	FedAvg+DART	141	88.07	75.73	81.90
VGG-16	FedAvg	135	89.90	72.08	80.99
	FedAugMix	36	59.31	52.83	56.07
	FedAvg+DART	135	88.81	77.14	82.98

Table 11 shows that using BigGAN [5]-generated data with FedAvg+DART still yields substantial robustness gains, improving \mathcal{A}_{rob} by 4.8

Table 10. CIFAR-10 clean and CIFAR-10-C robust accuracy results for FedAvg and FedAvg+DART with custom DART. The model used is ResNet-18. To improve robustness, DART must include both the consistency loss \mathcal{L}_c and the distillation loss \mathcal{L}_d .

Method	\mathcal{A}_{cln} (%)	\mathcal{A}_{rob} (%)	\mathcal{A}_{avg} (%)
Clean training	93.58	75.09	84.34
DART w/o \mathcal{L}_c	93.46	75.36	84.41
DART w/o \mathcal{L}_d	19.24	17.53	18.39
DART	92.08	80.79	86.44

13. Accounting for Server-side Costs in DART-enhanced FL

DART offloads the resource intensive robust training to the server. Fig. 10 shows that the large server side resources and parallelism reduce overall training time when DART is applied to FedAvg compared to other robust FL works. Crucially, client-side time and energy remain comparable

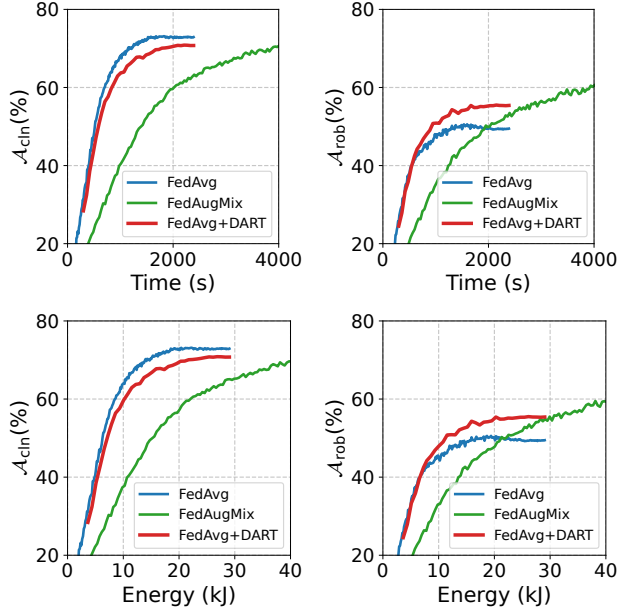


Figure 7. CIFAR-100 clean and CIFAR-100-C robust accuracy, \mathcal{A}_{cIn} and \mathcal{A}_{rob} , under FedAvg, FedAugMix, and FedAvg+DART as a function of time and energy. The model used is ResNet-18, the server dataset is CIFAR-10, and 10 clients are deployed.

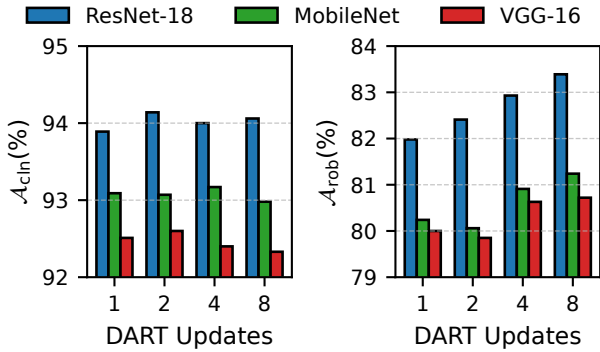


Figure 8. FedAvg+DART CIFAR-10 clean and CIFAR-10-C robust accuracy for different numbers of DART updates. The number of global updates, E , is set to 200. The number of clients is 4 and the server dataset is CIFAR-100.

to non-robust methods as shown in 4.

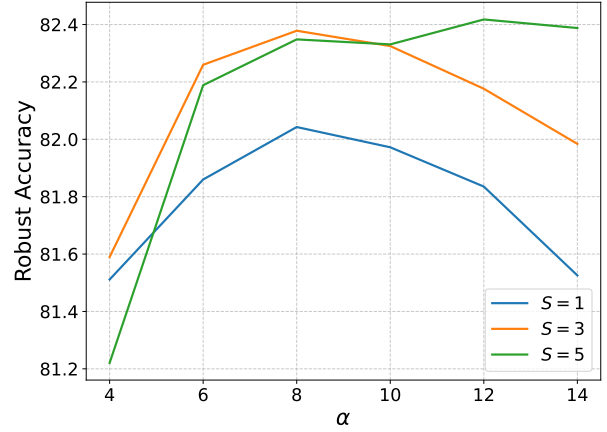


Figure 9. Impact of DART loss weighing parameter α and data augmentation mixture width S on CIFAR-10-C robust accuracy. We choose $\alpha = 10$ and $S = 3$ to maximize robustness and limit the server workload.

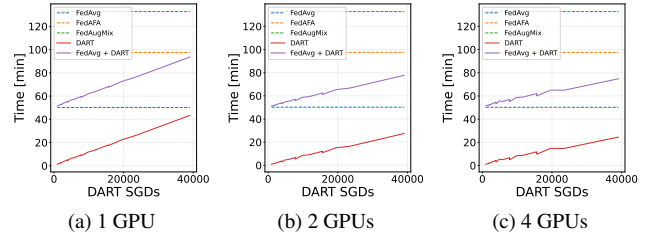


Figure 10. Total time versus number of DART SGD updates. DART-enhanced FedAvg maintains lower training time compared to robust FL methods.

Table 11. Impact of DART on CIFAR-10 clean and robust accuracy with BigGAN-generated server data.

Method	α_{iid}	\mathcal{A}_{cIn} (%)	\mathcal{A}_{rob} (%)
FedAvg	100	89.8	70.8
FedAvg	10	89.5	70.1
DART+FedAvg	100	86.4	73.7
DART+FedAvg	10	86.3	74.9