



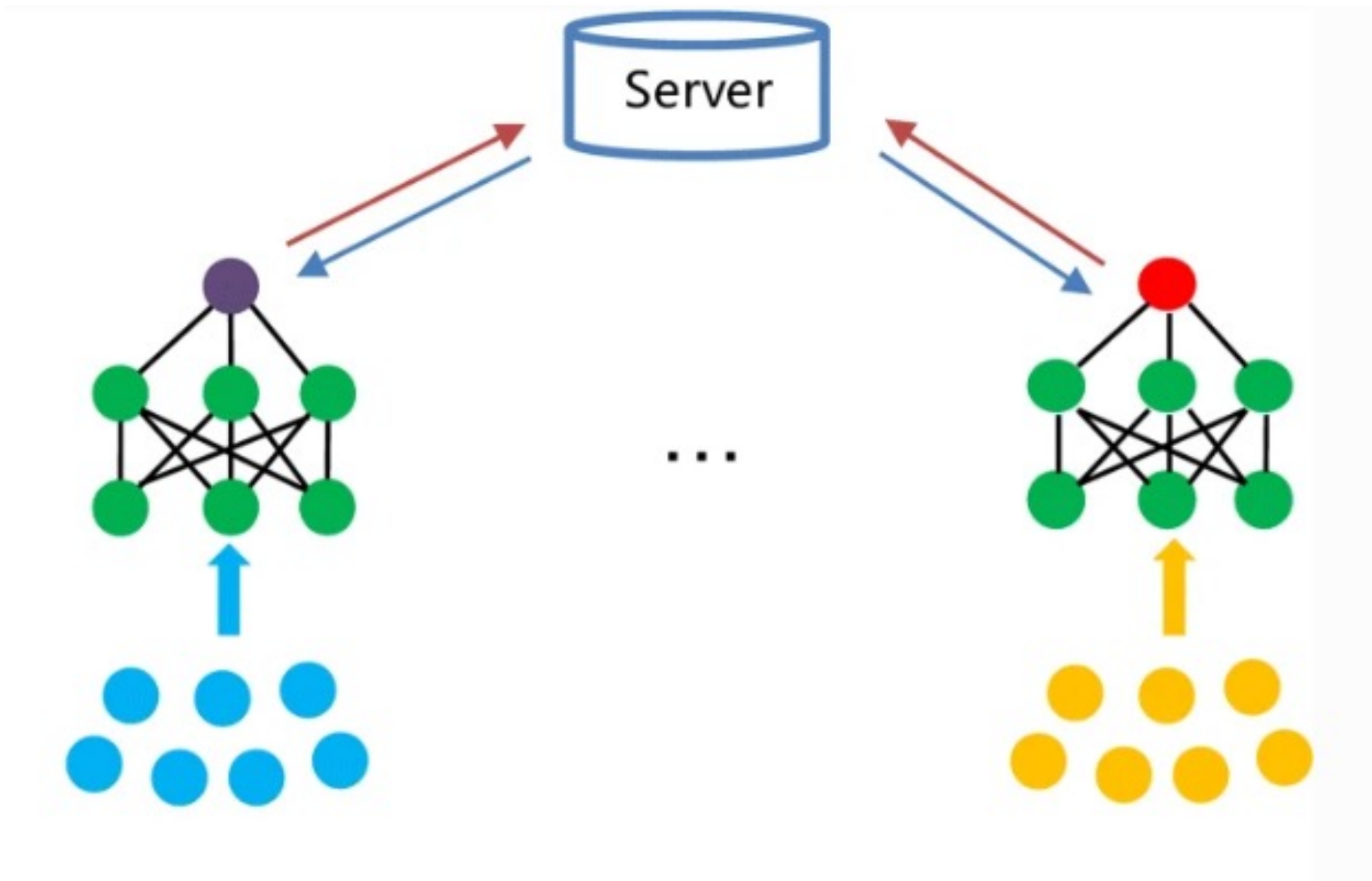
# Efficient Federated Learning

**Karolina Bogacka**

[karolina.bogacka.dokt@pw.edu.pl](mailto:karolina.bogacka.dokt@pw.edu.pl)

**21 grudnia 2022**

Szkoła Doktorska Politechniki Warszawskiej  
Wydział Matematyki i Nauk Informatycznych



Huang, X., Ding, Y., Jiang, Z.L. *et al.* DP-FL: a novel differentially private federated learning framework for the unbalanced data. *World Wide Web* **23**, 2529–2545 (2020). <https://doi.org/10.1007/s11280-020-00780-4>

**WHY?**

---



Privacy



Efficiency

# Common carbon footprint benchmarks

in lbs of CO2 equivalent

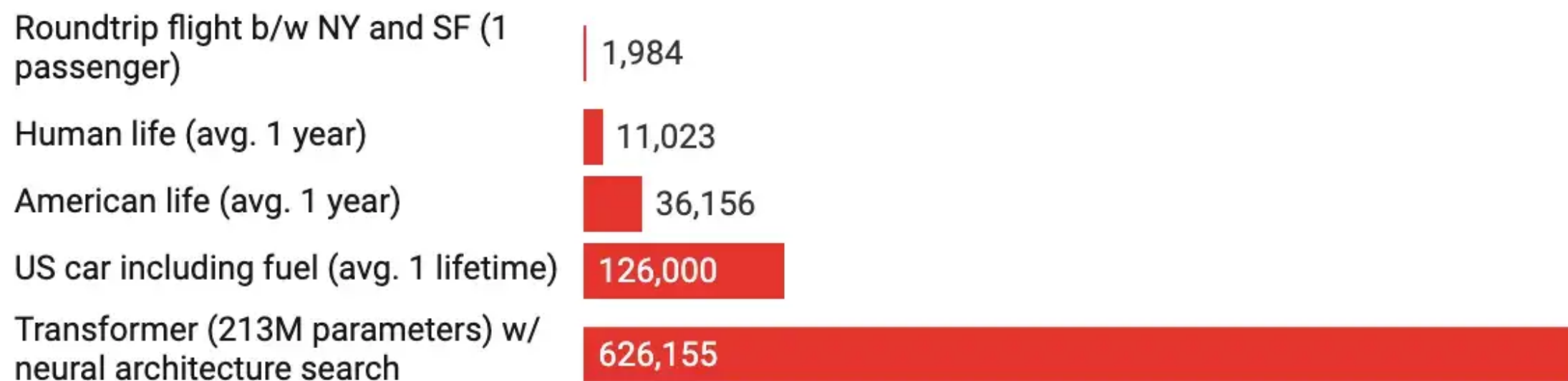
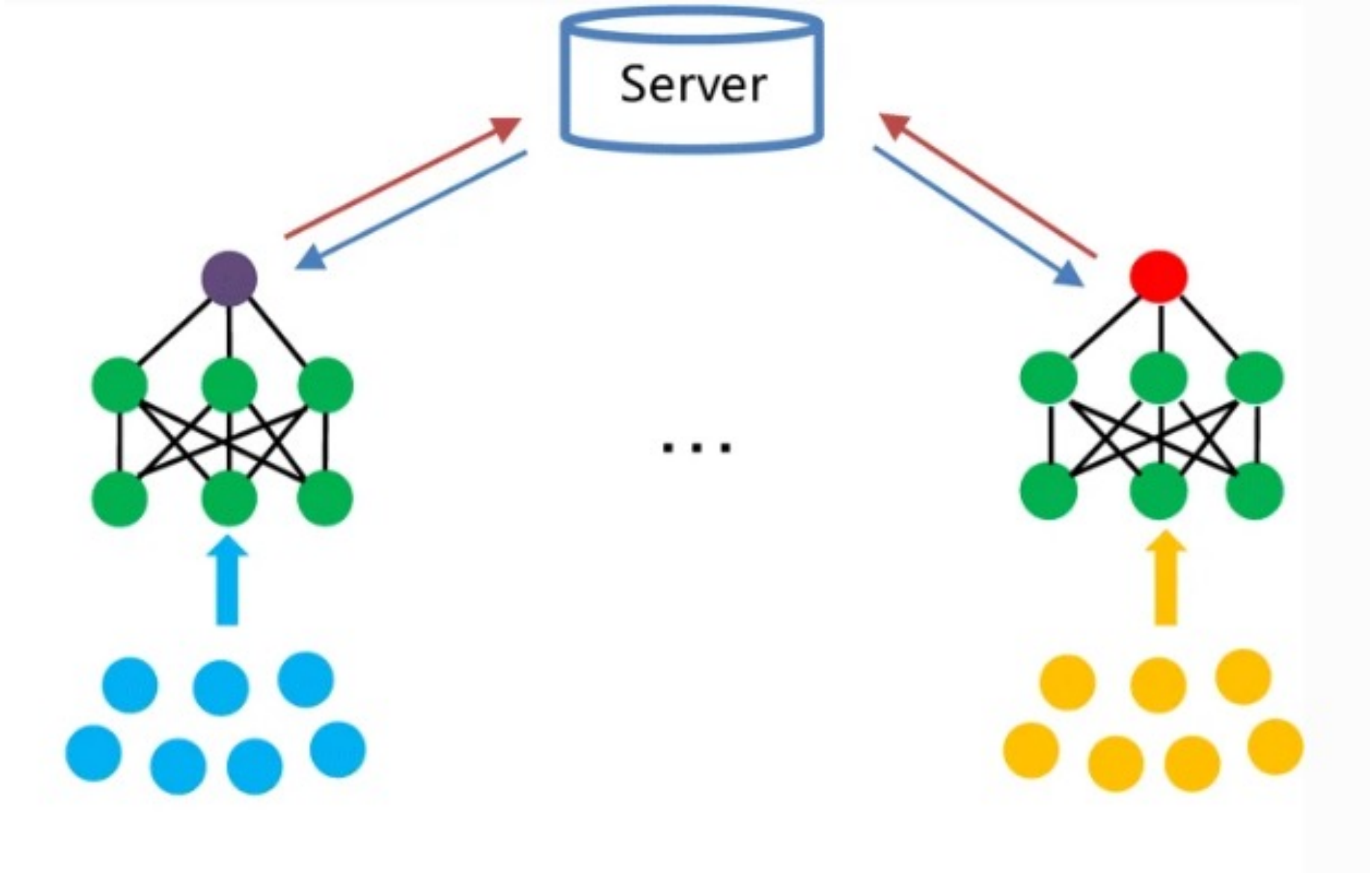


Chart: MIT Technology Review • Source: Strubell et al. • Created with Datawrapper





Number of Devices



Network  
Bandwidth



Limited Edge Node  
Computation



Statistical  
Heterogeneity



Local Updating



Client Selection



Reducing Model Updates



Decentralized Training and Different Topologies

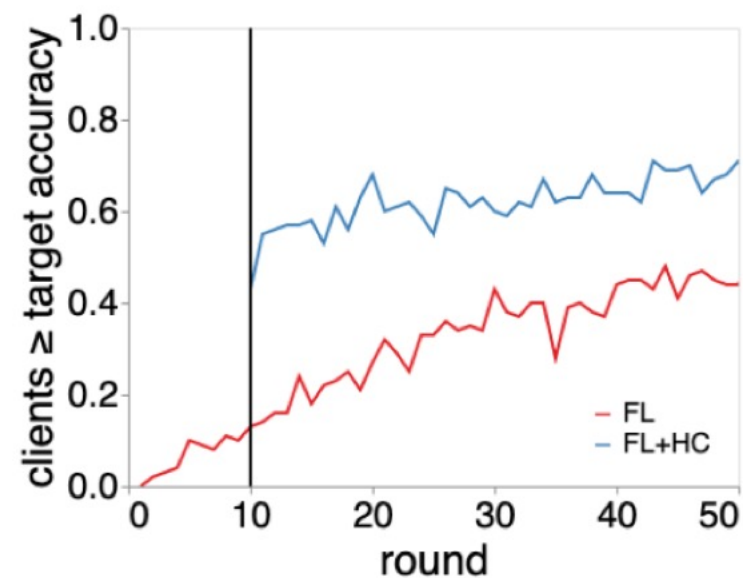
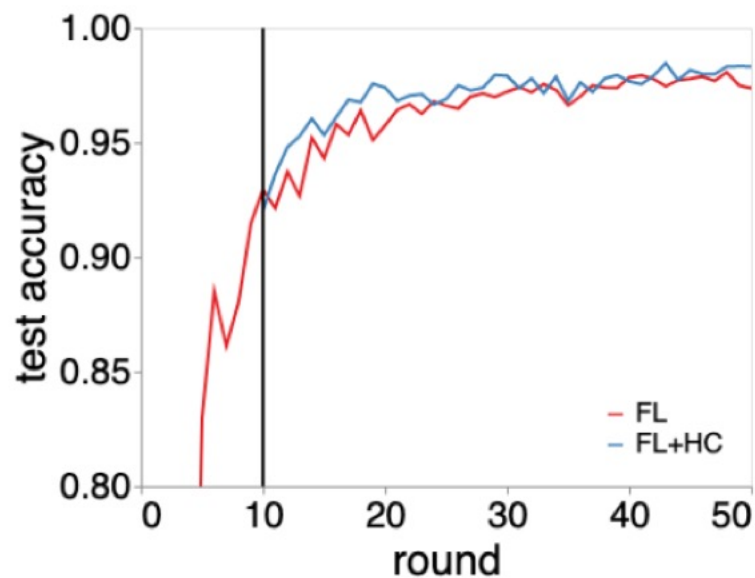


Compression Schemes



# Local Updating

Federated learning with hierarchical clustering of local updates to improve training on non-IID data



```

1: procedure FL+HC ▷ On server
2:   Initialise  $w_0$ 
3:   for each round  $t \in [1, n]$  do
4:      $w_{t+1} \leftarrow \text{FEDERATEDLEARNING}(w_t, K)$ 
5:   end for
6:    $w \leftarrow w_{t+1}$ 
7:   for each client  $k \in K$  do ▷ In parallel
8:      $\Delta w^k \leftarrow \text{CLIENTUPDATE}(k, w)$ 
9:   end for
10:   $C \leftarrow \text{HierarchicalClusteringAlgorithm}(\Delta w, P)$ 
11:  for  $c \in C$  do ▷ In parallel
12:     $w_{c,0} \leftarrow w$ 
13:    for each round  $t = 1, 2, \dots$  do
14:       $w_{c,t+1} \leftarrow \text{FEDERATEDLEARNING}(w_{c,t}, K_c)$ 
15:    end for
16:  end for
17: end procedure

```

```

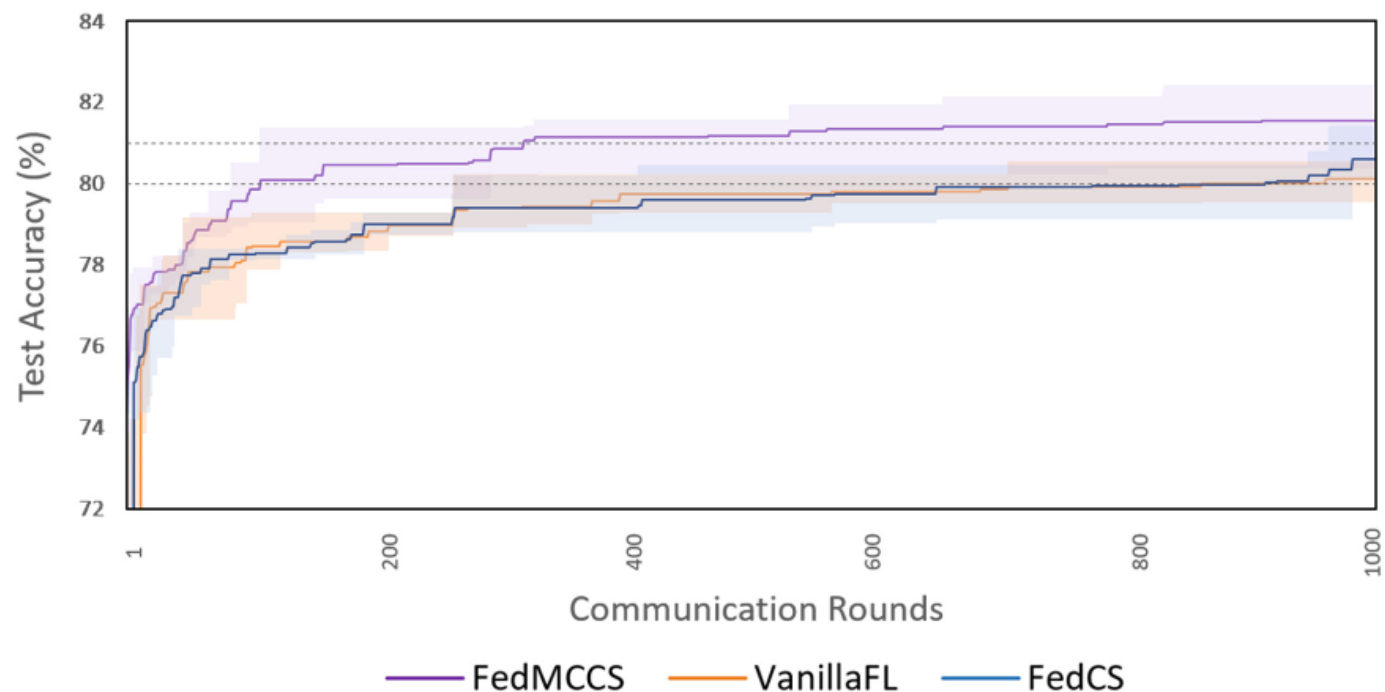
18: procedure FEDERATEDLEARNING( $w_t, K$ ) ▷ On server
19:    $m \leftarrow \max(\alpha \cdot K, 1)$ 
20:    $S_t \leftarrow$  (random set of  $m$  clients)
21:   for each client  $k \in S_t$  do ▷ In parallel
22:      $w_{t+1}^k \leftarrow \text{CLIENTUPDATE}(k, w_t)$ 
23:   end for
24:    $w_{t+1} \leftarrow \sum_{k=1}^K \frac{n_k}{n} w_{t+1}^k$ 
25: end procedure

26: procedure CLIENTUPDATE( $k, w$ ) ▷ On client  $k$ 
27:    $\mathcal{B} \leftarrow$  (Split  $\mathcal{P}_k$  into batches of size  $B$ )
28:   for each local epoch  $i$  from 1 to  $E$  do
29:     for batch  $b \in \mathcal{B}$  do
30:        $w \leftarrow w - \eta \nabla \mathcal{L}(w; b)$ 
31:     end for
32:   end for
33:   return  $w$  to server
34: end procedure

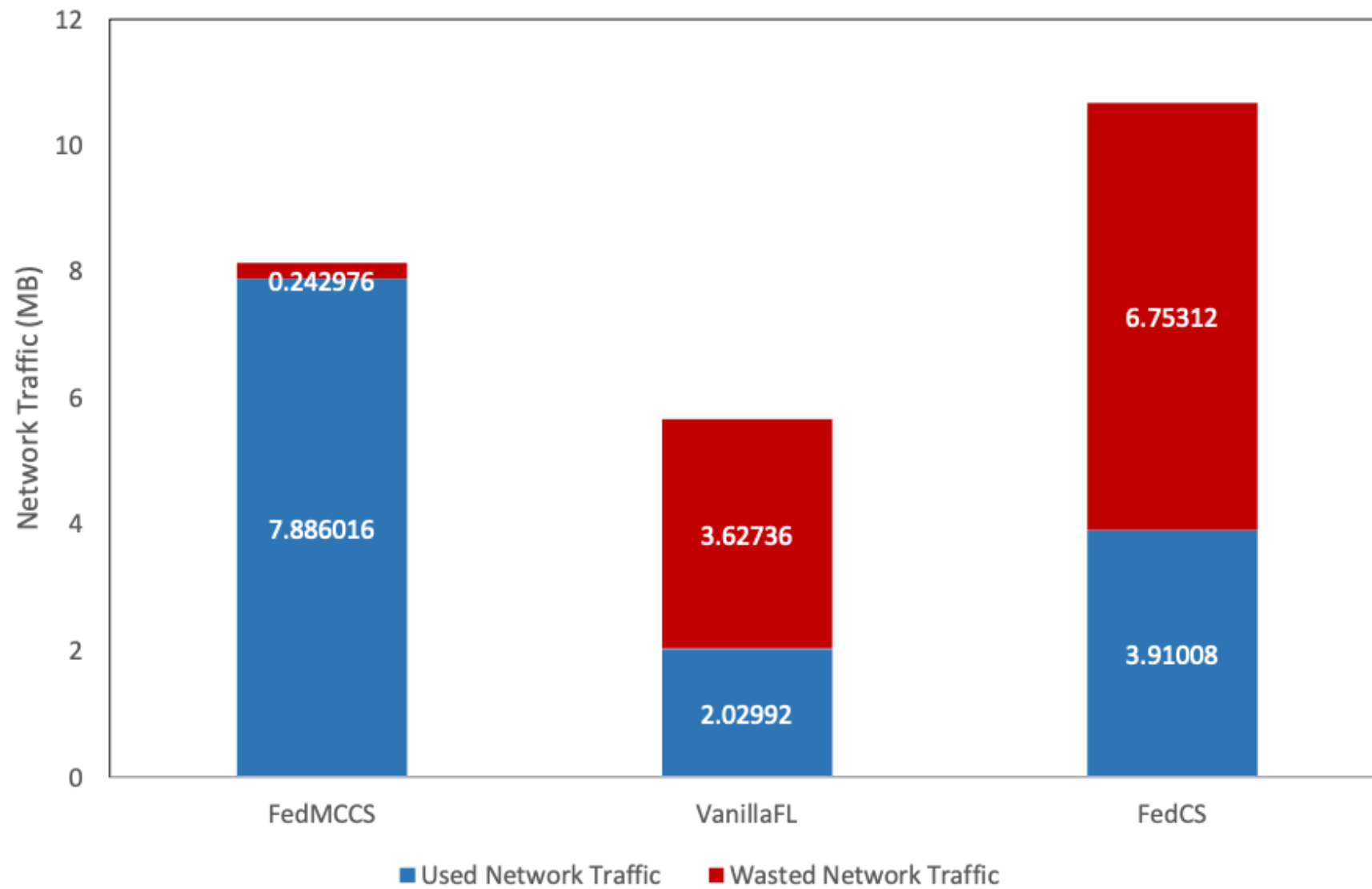
```

# Client Selection

FedMCCS:  
Multicriteria Client  
Selection Model for  
Optimal IoT  
Federated Learning



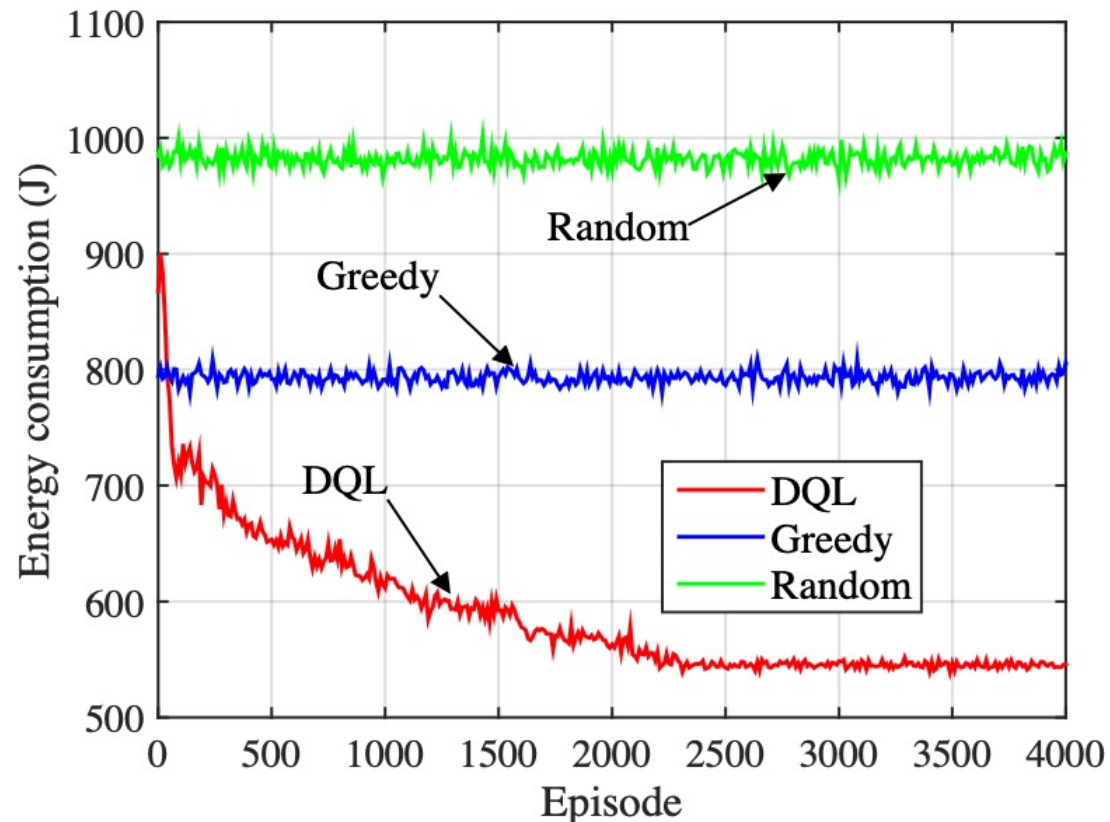
## Network Traffic between Used and Wasted



- 1: Initialization in Protocol 1.
- 2: Client Filtering : The server applies Stratified-based filtering to select clients according to their metadata, avoiding communications with irrelevant clients.
- 3: Resource Request : The server requests resource information from the filtered clients.
  
- 4: Multi-criteria Client Selection : Based on the clients responses, the server uses Multi-Criteria selection approach to determine a maximum of  $[K \times C]$  clients to participate in the remaining steps.
  
- 5: Distribution : The server disseminates the global model parameters to the selected clients.
- 6: Update and Upload in Protocol 1.
- 7: Aggregation : The server averages the parameters, when more than 70% of the requested updates are received.
- 8: All steps but Initialization are iterated as in Protocol 2.

# Client Selection

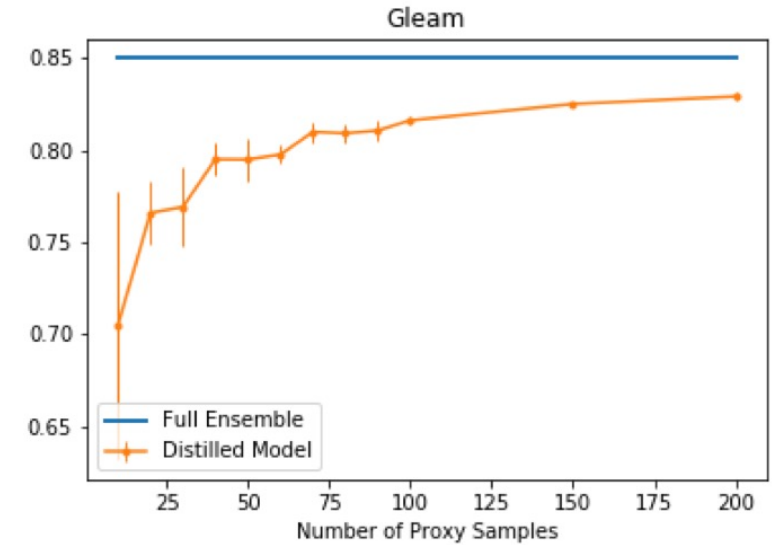
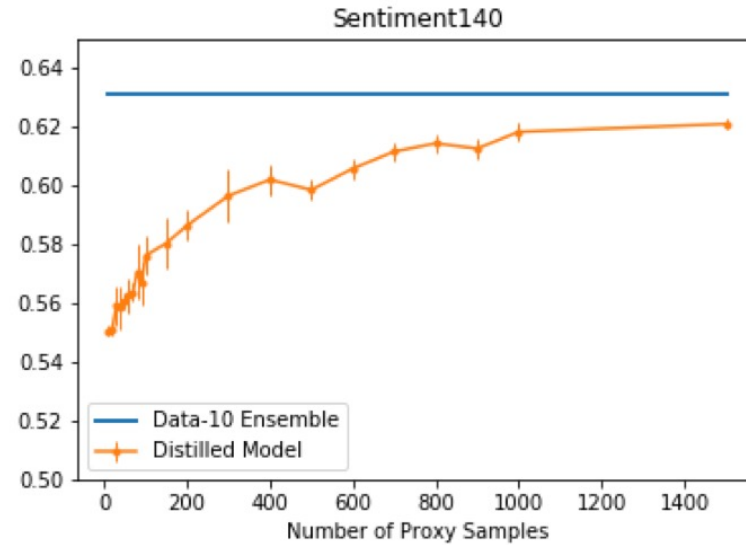
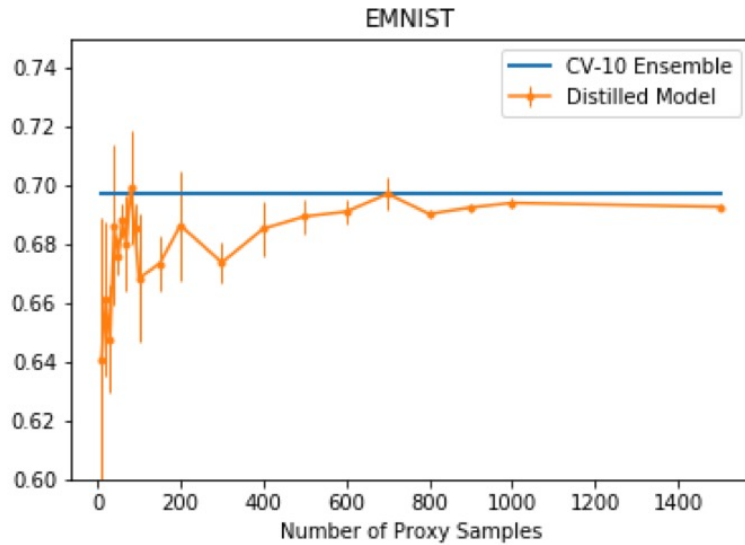
Efficient Training  
Management for Mobile  
Crowd-Machine  
Learning: A Deep  
Reinforcement Learning  
Approach



# Client Selection

```
1: Initialize:  $\theta, \theta^-$ ;
2: for episode  $i = 1$  to  $N$  do
3:   for iteration  $t = 1$  to  $T$  do
4:     Select action  $a$  according to the  $\epsilon$ -greedy policy;
5:     Execute action  $a$  and observe reward  $r$  and next state  $s'$ ;
6:     Store experience  $e = \langle s, a, r, s' \rangle$  in  $\mathcal{M}$ ;
7:     Select  $N_b$  experiences  $e_k = \langle s_k, a_k, r_k, s'_k \rangle$  from  $\mathcal{M}$ ;
8:     for  $k = 1$  to  $N_b$  do
9:       Determine  $a^{\max} = \arg \max_{a' \in \mathcal{A}} Q(s'_k, a'; \theta)$ ;
10:      Calculate  $y_k = r_k + \gamma Q(s'_k, a^{\max}; \theta^-)$ ;
11:    end for
12:    Define  $\bar{L}(\theta) = \frac{1}{N_b} \sum_{k=1}^{N_b} \left( y_k - Q(s_k, a_k; \theta) \right)^2$ ;
13:    Perform a gradient descent step on  $\bar{L}(\theta)$  to update  $\theta$ ;
14:    Reset  $\theta^- = \theta$  in every  $N^-$  iteration;
15:  end for
16: end for
```

# Reducing Model Updates

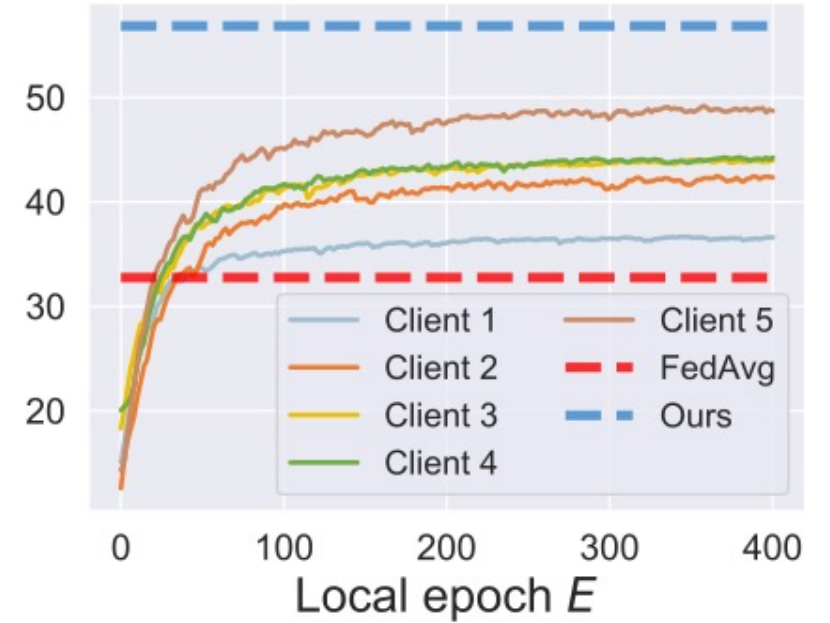
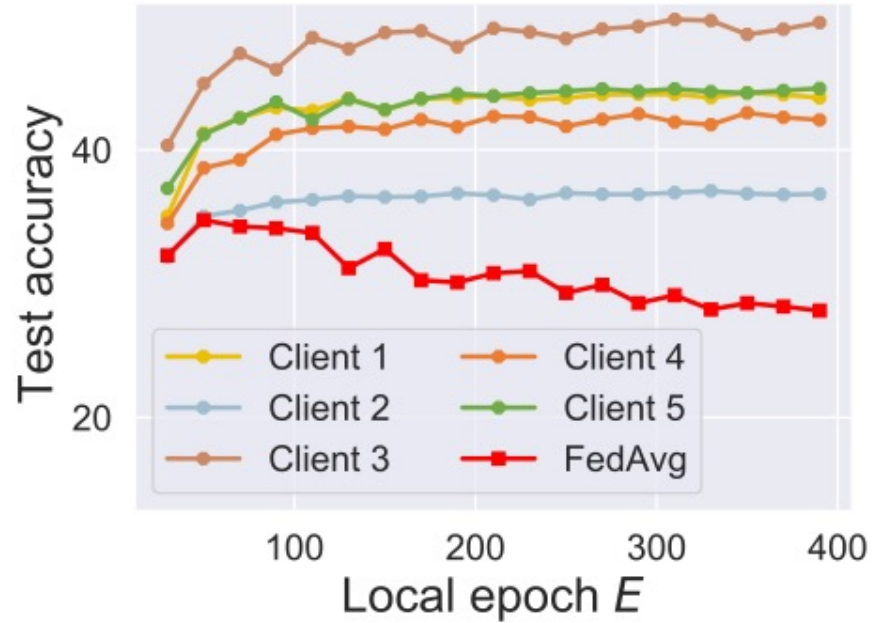


## One-Shot Federated Learning

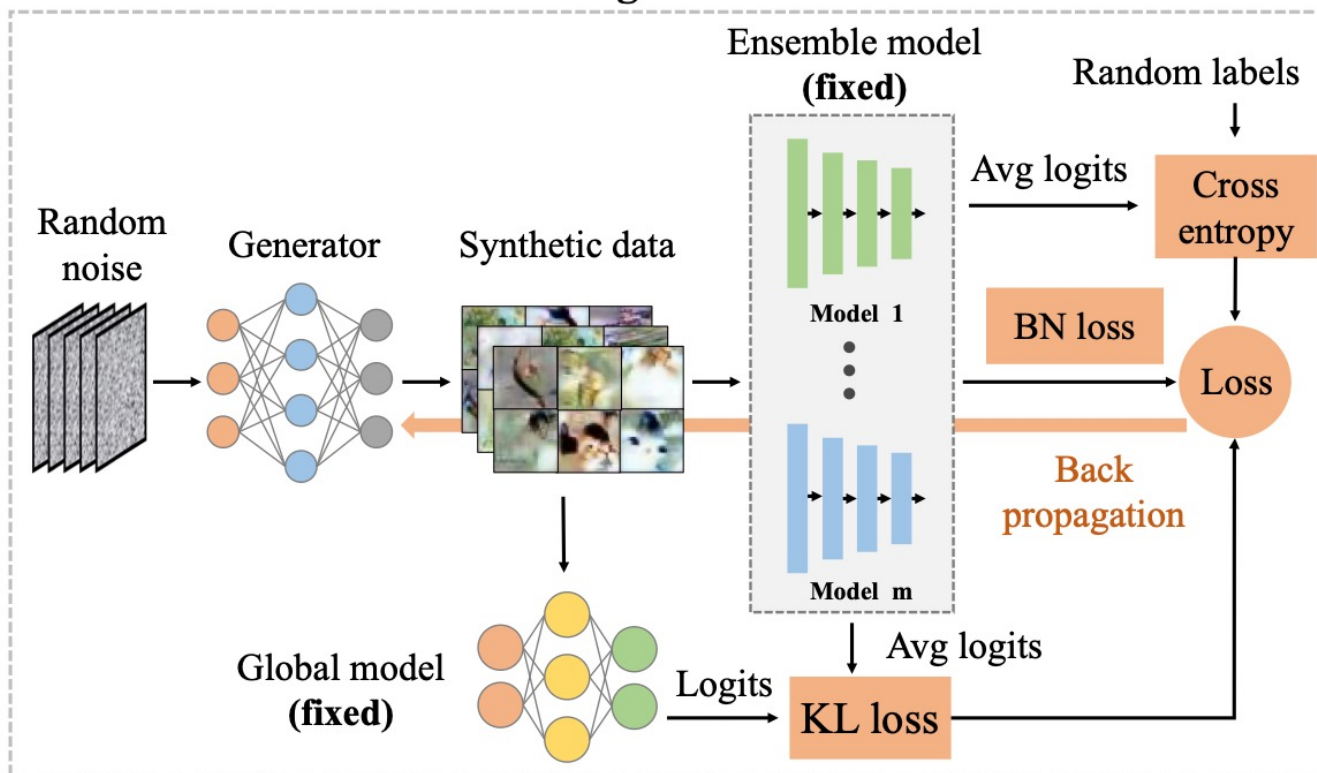


# Reducing Model Updates

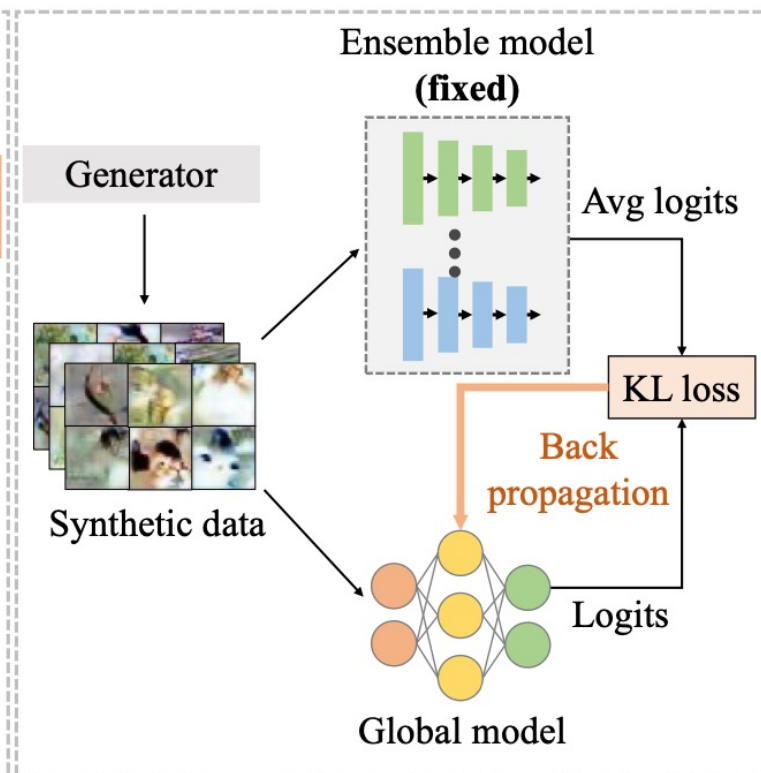
DENSE:  
Data-Free  
One-Shot  
Federated  
Learning



## Data generation



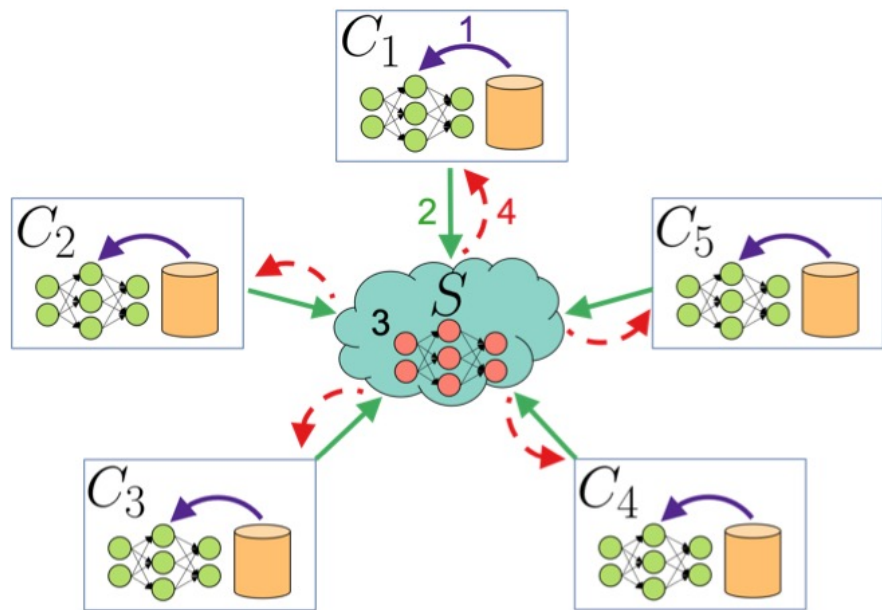
## Model distillation



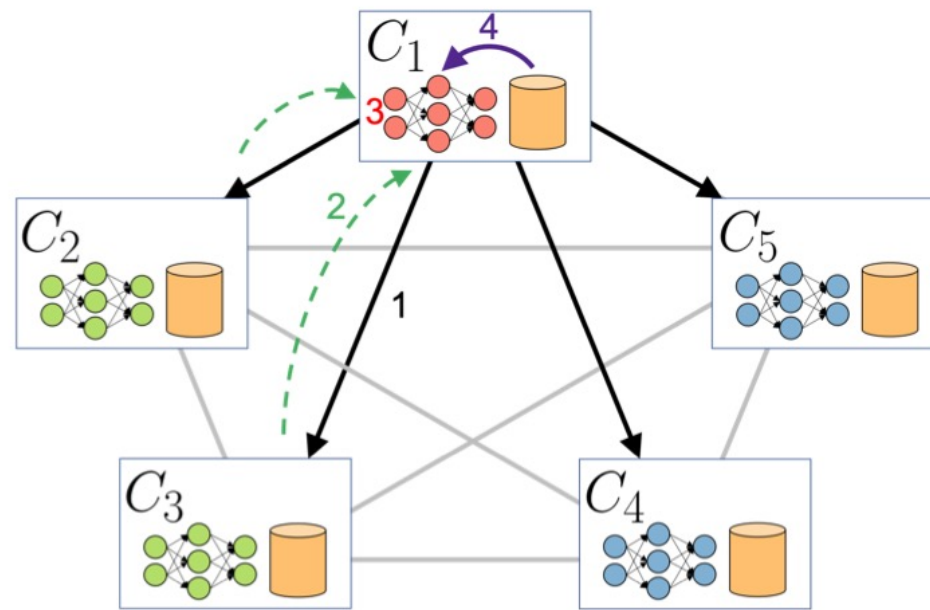
# Decentralized Training

BrainTorrent:  
A Peer-to-Peer  
Environment  
for  
Decentralized  
Federated  
Learning

# Clients	Scans/client	Avg. Dice over Clients		Aggregated Model	
		FLS	BrainTorrent	FLS	BrainTorrent
5	4	0.812	0.851	0.845	<b>0.863</b>
7	3	0.753	0.837	0.843	0.861
10	2	0.792	0.807	0.842	0.850
20	1	0.570	0.578	0.687	0.728
Pooled Model			<b>0.866</b>		



(a) Federated Learning with server



(b) BrainTorrent: P2P serverless Federated Learning



Initialize  $N$  clients models,  $\mathbf{C} = \{C_1, \dots, C_N\}$  with random weights;  
Initialize  $N$  version vectors  $\mathbf{V} = [\mathbf{v}^1, \dots, \mathbf{v}^N]$  with all zero entries;  
**for** *round*  $r$  *in*  $1, 2, \dots$  **do**  
    Randomly select a client  $i$  from  $\{1, \dots, N\}$ ;  
     $\mathbf{v}_{\text{old}} \leftarrow \mathbf{v}^i$ ;  
     $\mathbf{v}_{\text{new}} \leftarrow \text{ping\_request}(C_i \rightarrow \mathbf{C})$ ;  
     $\mathbf{W} \leftarrow \frac{a_i}{a} \mathbf{W}^i$  ;  
    **for**  $j \in \{1, \dots, i - 1, i + 1, \dots, N\}$  **do**  
        **if**  $v_{\text{new}}^j > v_{\text{old}}^j$  **then**  
            Receive updated  $\mathbf{W}^j$  and  $a_j$  from  $C_j$  ;  
        **end**  
         $\mathbf{W} \leftarrow \mathbf{W} + \frac{a_j}{a} \mathbf{W}^j$  ;  
    **end**  
     $\mathbf{W}^i \leftarrow \text{FineTune}(\mathbf{W}, \mathcal{D}_i)$  ;  
    Increment  $\mathbf{v}^i(i)$ ;  
**end**

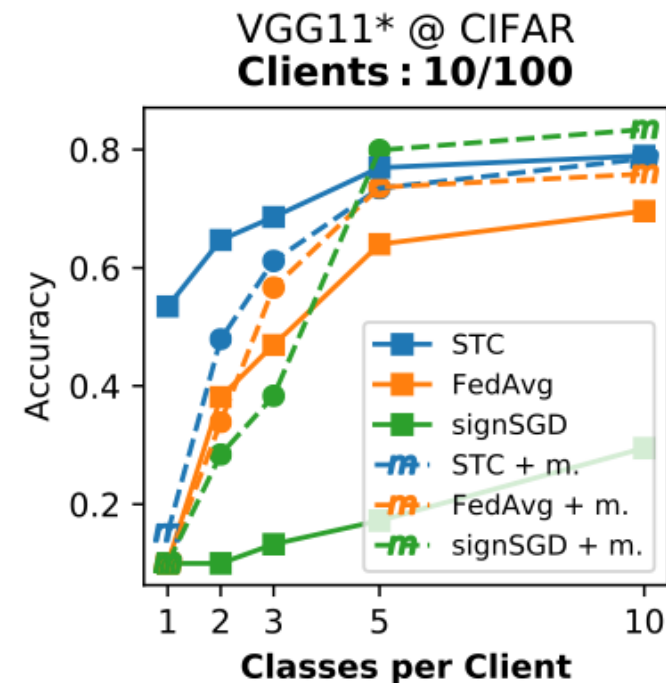
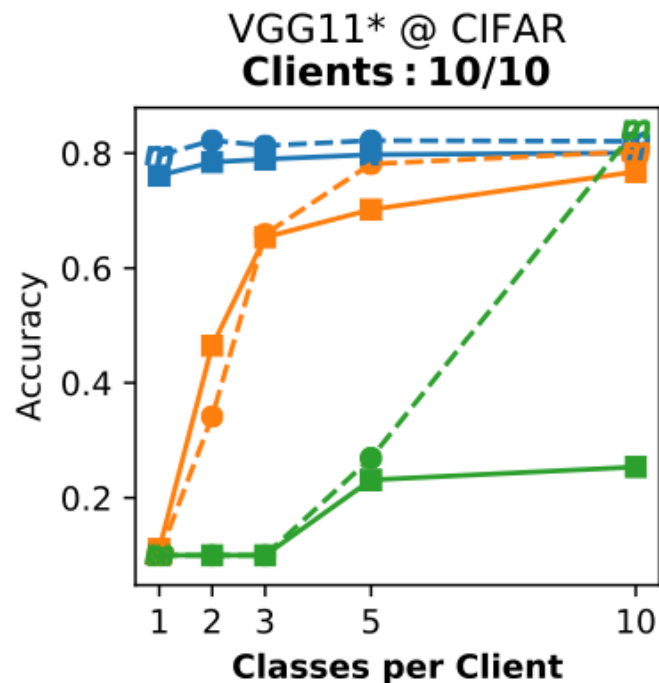
# Compression Schemes

Sparsification

Quantization

# Compression Schemes

Robust and  
Communication  
-Efficient  
Federated  
Learning from  
Non-IID Data



1 **input:** initial parameters  $\mathcal{W}$   
2 **output:** improved parameters  $\mathcal{W}$   
3 **init:** all clients  $C_i, i = 1, \dots, [\text{Number of Clients}]$  are initialized with the same parameters  $\mathcal{W}_i \leftarrow \mathcal{W}$ . Every Client holds a different dataset  $D_i$ , with  $|\{y : (x, y) \in D_i\}| = [\text{Classes per Client}]$  of size  $|D_i| = \varphi_i |\cup_j D_j|$ . The residuals are initialized to zero  $\Delta\mathcal{W}, \mathcal{R}_i, \mathcal{R} \leftarrow 0$ .

4 **for**  $t = 1, \dots, T$  **do**

5     **for**  $i \in I_t \subseteq \{1, \dots, [\text{Number of Clients}]\}$  *in parallel*  
   **do**

6         Client  $C_i$  does:

- 7         •  $\text{msg} \leftarrow \text{download}_{S \rightarrow C_i}(\text{msg})$   
8         •  $\Delta\mathcal{W} \leftarrow \text{decode}(\text{msg})$   
9         •  $\mathcal{W}_i \leftarrow \mathcal{W}_i + \Delta\mathcal{W}$   
10        •  $\Delta\mathcal{W}_i \leftarrow \mathcal{R}_i + \text{SGD}(\mathcal{W}_i, D_i, b) - \mathcal{W}_i$   
11        •  $\Delta\tilde{\mathcal{W}}_i \leftarrow \text{STC}_{p_{up}}(\Delta\mathcal{W}_i)$   
12        •  $\mathcal{R}_i \leftarrow \Delta\mathcal{W}_i - \Delta\tilde{\mathcal{W}}_i$

13        •  $\text{msg}_i \leftarrow \text{encode}(\Delta\tilde{\mathcal{W}}_i)$

14        •  $\text{upload}_{C_i \rightarrow S}(\text{msg}_i)$

15     **end**

16     Server  $S$  does:

17     •  $\text{gather}_{C_i \rightarrow S}(\Delta\tilde{\mathcal{W}}_i), i \in I_t$

18     •  $\Delta\mathcal{W} \leftarrow \mathcal{R} + \frac{1}{|I_t|} \sum_{i \in I_t} \Delta\tilde{\mathcal{W}}_i$

19     •  $\Delta\tilde{\mathcal{W}} \leftarrow \text{STC}_{p_{down}}(\Delta\mathcal{W})$

20     •  $\mathcal{R} \leftarrow \Delta\mathcal{W} - \Delta\tilde{\mathcal{W}}$

21     •  $\mathcal{W} \leftarrow \mathcal{W} + \Delta\tilde{\mathcal{W}}$

22     •  $\text{msg} \leftarrow \text{encode}(\Delta\tilde{\mathcal{W}})$

23     •  $\text{broadcast}_{S \rightarrow C_i}(\text{msg}), i = 1, \dots, M$

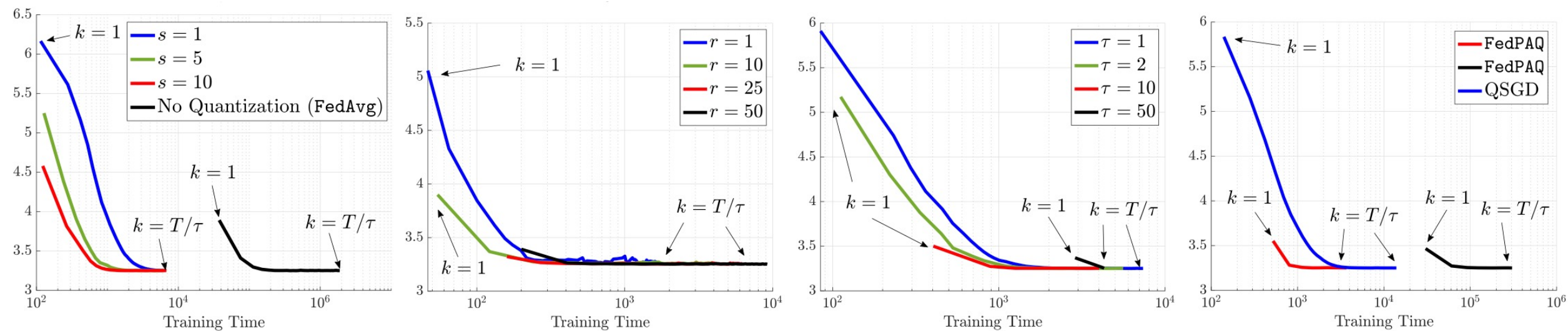
24 **end**

25 **return**  $\mathcal{W}$



- 1 **input:** flattened tensor  $T \in \mathbb{R}^n$ , sparsity  $p$
- 2 **output:** sparse ternary tensor  $T^* \in \{-\mu, 0, \mu\}^n$
- 3 •  $k \leftarrow \max(np, 1)$
- 4 •  $v \leftarrow \text{top}_k(|T|)$
- 5 •  $\text{mask} \leftarrow (|T| \geq v) \in \{0, 1\}^n$
- 6 •  $T^{\text{masked}} \leftarrow \text{mask} \odot T$
- 7 •  $\mu \leftarrow \frac{1}{k} \sum_{i=1}^n |T_i^{\text{masked}}|$
- 8 **return**  $T^* \leftarrow \mu \times \text{sign}(T^{\text{masked}})$

# Compression Schemes



**FedPAQ: A Communication-Efficient Federated Learning Method with Periodic Averaging and Quantization**

```

1: for  $k = 0, 1, \dots, K - 1$  do
2:   server picks  $r$  nodes  $\mathcal{S}_k$  uniformly at random
3:   server sends  $\mathbf{x}_k$  to nodes in  $\mathcal{S}_k$ 
4:   for node  $i \in \mathcal{S}_k$  do
5:      $\mathbf{x}_{k,0}^{(i)} \leftarrow \mathbf{x}_k$ 
6:     for  $t = 0, 1, \dots, \tau - 1$  do
7:       compute stochastic gradient
8:        $\tilde{\nabla} f_i(\mathbf{x}) = \nabla \ell(\mathbf{x}, \xi)$  for a  $\xi \in \mathcal{P}^i$ 
9:       set  $\mathbf{x}_{k,t+1}^{(i)} \leftarrow \mathbf{x}_{k,t}^{(i)} - \eta_{k,t} \tilde{\nabla} f_i(\mathbf{x}_{k,t}^{(i)})$ 
10:    end for
11:    send  $Q(\mathbf{x}_{k,\tau}^{(i)} - \mathbf{x}_k)$  to the server
12:  end for
13:  server finds  $\mathbf{x}_{k+1} \leftarrow \mathbf{x}_k + \frac{1}{r} \sum_{i \in \mathcal{S}_k} Q(\mathbf{x}_{k,\tau}^{(i)} - \mathbf{x}_k)$ 
14: end for

```

# Bibliography

---

- Huang, X., Ding, Y., Jiang, Z.L. *et al.* DP-FL: a novel differentially private federated learning framework for the unbalanced data. *World Wide Web* **23**, 2529–2545 (2020). <https://doi.org/10.1007/s11280-020-00780-4>
- Strubell, Emma & Ganesh, Ananya & McCallum, Andrew. (2019). Energy and Policy Considerations for Deep Learning in NLP. 3645-3650. 10.18653/v1/P19-1355.
- Shahid, Osama, et al. "Communication efficiency in federated learning: Achievements and challenges." *arXiv preprint arXiv:2107.10996* (2021).
- Briggs, Christopher, Zhong Fan, and Peter Andras. "Federated learning with hierarchical clustering of local updates to improve training on non-IID data." *2020 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 2020.
- Anh, Tran The, et al. "Efficient training management for mobile crowd-machine learning: A deep reinforcement learning approach." *IEEE Wireless Communications Letters* 8.5 (2019): 1345-1348.

# Bibliography

---

- S. Abdulrahman, H. Tout, A. Mourad and C. Talhi, "FedMCCS: Multicriteria Client Selection Model for Optimal IoT Federated Learning," in *IEEE Internet of Things Journal*, vol. 8, no. 6, pp. 4723-4735, 15 March 2021, doi: 10.1109/JIOT.2020.3028742.
- Guha, Neel, Ameet Talwalkar, and Virginia Smith. "One-shot federated learning." *arXiv preprint arXiv:1902.11175* (2019).
- Zhang, Jie, et al. "A Practical Data-Free Approach to One-shot Federated Learning with Heterogeneity." *arXiv preprint arXiv:2112.12371* (2021).
- Roy, Abhijit Guha et al. "BrainTorrent: A Peer-to-Peer Environment for Decentralized Federated Learning." *ArXiv abs/1905.06731* (2019): n. pag.
- F. Sattler, S. Wiedemann, K. -R. Müller and W. Samek, "Robust and Communication-Efficient Federated Learning From Non-i.i.d. Data," in *IEEE Transactions on Neural Networks and Learning Systems*, vol. 31, no. 9, pp. 3400-3413, Sept. 2020, doi: 10.1109/TNNLS.2019.2944481.
- Reisizadeh, Amirhossein, et al. "Fedpaq: A communication-efficient federated learning method with periodic averaging and quantization." *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020.