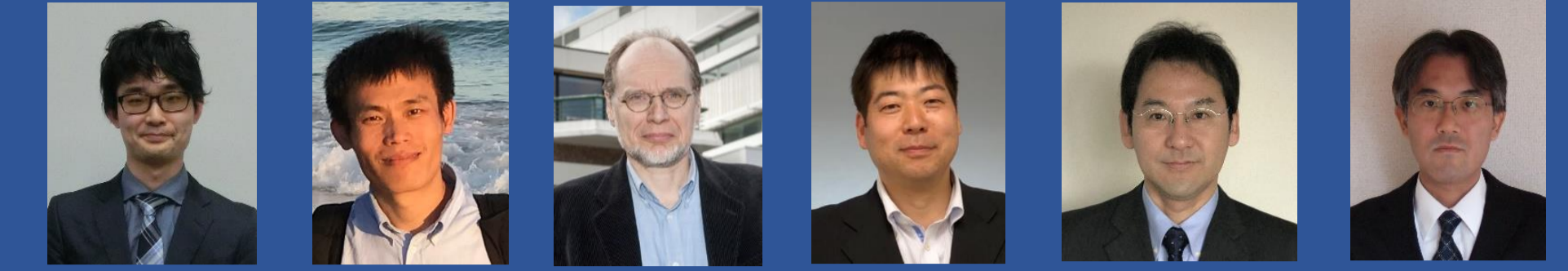


Asynchronous Decentralized Optimization with Implicit Stochastic Variance Reduction

Kenta Niwa^{1,2}, Guoqiang Zhang³, W. Bastiaan Kleijn⁴, Noboru Harada^{1,2}, Hiroshi Sawada¹, and Akinori Fujino¹

1: NTT Communication Science Laboratories, 2: NTT Media Intelligence Laboratories, 3: University Technology of Sydney, 4: Victoria University of Wellington



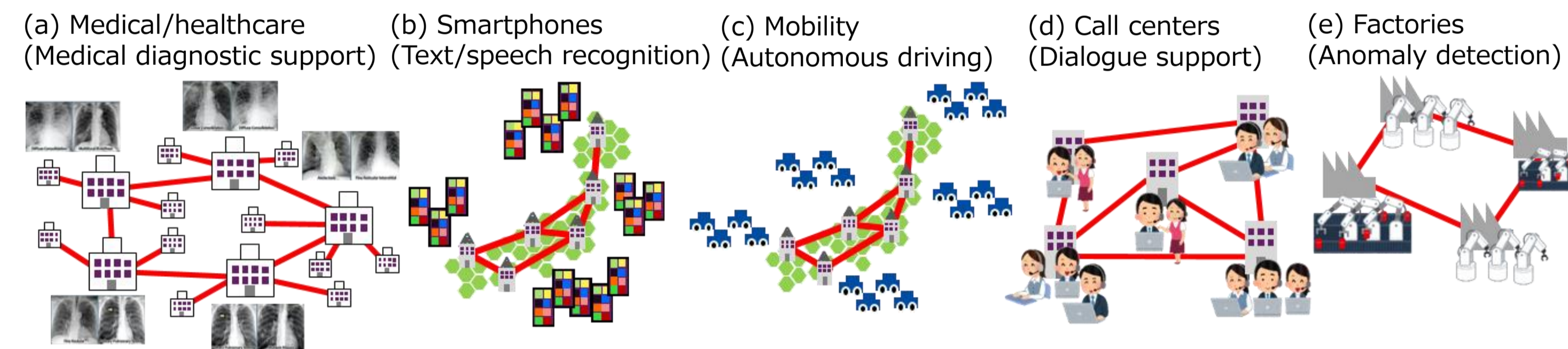
Goal

Background: We are entering an era of distributed data processing due to data volume, privacy-aware issues, and legal regulations, e.g., GDPR.

Goal: To train ML models without aggregating data to a central cloud.

- Nodes are connected by **decentralized network** for flexible scalability.
- Asynchronous communication** among nodes are allowed.
- Data subsets held on local nodes are **statistically heterogeneous (non-IID)**.

Application Examples:



Conventional studies

Independently of network/communication configurations, many distributed training algorithms can be categorized into three trends.

- Average Consensus: SGD + Average**
Weak robustness to non-IID data subsets
- FedAvg [McMahan et al., 2017]
- DSGD [Chen & Sayed, 2012]
- Gossip SGD [Ormandi et al. 2013]
- FedProx [Li et al., 2019]
- Stochastic Variance Reduction (SVR):**
Stochastic gradient modification using global/local control variates.
Global control variate Local control variate
$$\mathbf{w}_i^{k+1} = \mathbf{w}_i^k - \mu(g_i(\mathbf{w}_i^k) + \bar{\mathbf{c}}_i^k - \mathbf{c}_i^k),$$

- SVRG [Johnson & Zhang, 2013]
- SAGA [Defazio et al., 2014]
- SCAFFOLD [Karimireddy et al., 2020]
- GT-SVR [Xin et al., 2020]
- Primal-dual formalism:**
Solve **model matching constraint cost-sum minimization** problem
$$\inf_{\mathbf{w}_1, \dots, \mathbf{w}_N} \frac{1}{N} \sum_{i \in \mathcal{N}} f_i(\mathbf{w}_i) \quad \text{s.t. } \mathbf{A}_{i|j} \mathbf{w}_i + \mathbf{A}_{j|i} \mathbf{w}_j = \mathbf{0}, \quad (\forall i \in \mathcal{N}, j \in \mathcal{E}_i),$$

$$\{\mathbf{A}_{i|j}, \mathbf{A}_{j|i}\} = \{\mathbf{I}, -\mathbf{I}\}$$

E.g., Update rule of Edge-Consensus Learning (ECL) [Niwa et al., 2020]
$$\mathbf{w}_i^{k+1} = (\mathbf{w}_i^k - \mu g_i(\mathbf{w}_i^k) + \mu \eta_i \sum_{j \in \mathcal{E}_i} \mathbf{A}_{i|j}^T \mathbf{z}_{i|j}^k) / (1 + \mu \eta_i E_i),$$

Stochastic gradient modification using dual variables $\mathbf{z}_{i|j}$
- Distributed ADMM [W. Shi et al., 2014]
- Primal-Dual Method of Multiplier (PDMM) [G. Zhang et al., 2017, T. Sherson et al., 2018]
- ECL [Niwa et al., 2020]
- FedSplit [Pathak & Wainwright, 2020]

Note the similarity of the adjustment of stochastic gradient descent for SVR and for the primal dual formalism (e.g., ECL).

Main contribution

Key idea: Primal-dual formalism (e.g., ECL) may have an **optimal condition where it matches SVR**.

Reformulating ECL update rule:

$$\begin{aligned} \mathbf{w}_i^{k+1} &= (\mathbf{w}_i^k - \mu g_i(\mathbf{w}_i^k) + \mu \eta_i \sum_{j \in \mathcal{E}_i} \mathbf{A}_{i|j}^T \mathbf{z}_{i|j}^k) / (1 + \mu \eta_i E_i) \\ &= \mathbf{w}_i^k - \mu [g_i(\mathbf{w}_i^k) + \frac{\eta_i}{1 + \mu \eta_i E_i} \{ \sum_{j \in \mathcal{E}_i} (\mathbf{w}_i^k - \mathbf{A}_{i|j}^T \mathbf{z}_{i|j}^k) - \mu E_i g_i(\mathbf{w}_i^k) \}] \end{aligned}$$

ECL matches with SVR $\mathbf{w}_i^{k+1} = \mathbf{w}_i^k - \mu [g_i(\mathbf{w}_i^k) + \bar{\mathbf{c}}_i^k - \mathbf{c}_i^k]$ if underlined terms are modification using global/local control variates $\bar{\mathbf{c}}_i^k - \mathbf{c}_i^k$.

Contribution: We **optimally select parameter η_i** such that ECL matches with SVR.

- Investigating physical meaning of affine dual variables $\mathbf{A}_{i|j}^T \mathbf{z}_{i|j}$, it is proportion to the sum of update difference between nodes. (-> a part of global control variate)
- By reformulating \mathbf{w} -update rule, we can optimally set η_i to follow SVR. $\eta_i = 1 / (\mu E_i (K - 1))$

Proposed method (ECL with Implicit SVR: ECL-ISVR)

- Optimal η_i is selected for the previous ECL [Niwa et al., 2020].

Cost function (with quadratic approximation)	Previous ECL
$\inf_{\mathbf{w}_1, \dots, \mathbf{w}_N} \frac{1}{N} \sum_{i \in \mathcal{N}} f_i(\mathbf{w}_i)$ s.t. $\mathbf{A}_{i j} \mathbf{w}_i + \mathbf{A}_{j i} \mathbf{w}_j = \mathbf{0}, \quad (\forall i \in \mathcal{N}, j \in \mathcal{E}_i),$	$\mathbf{w}_i^{r,k+1} = \arg \min (g_i(\mathbf{w}_i) + \sum_{j \in \mathcal{N}_i} \frac{\rho_i}{2} \ \mathbf{w}_i - \mathbf{w}_j^{r,k}\ ^2),$
$g_i(\mathbf{w}_i) = f_i(\mathbf{w}_i^{r,k}) + \langle g_i(\mathbf{w}_i^{r,k}), \mathbf{w}_i - \mathbf{w}_i^{r,k} \rangle + \frac{1}{2\mu} \ \mathbf{w}_i - \mathbf{w}_i^{r,k}\ ^2.$	$\mathbf{y}_{i j}^{r,k+1} = \mathbf{z}_{i j}^{r,k} - 2\mathbf{A}_{i j} \mathbf{w}_i^{r,k+1}$ $\mathbf{z}_{i j}^{r,k+1} = \begin{cases} \mathbf{y}_{j i}^{r,k+1} & \text{(PDMM-SGD)} \\ \frac{1}{2} \mathbf{y}_{j i}^{r,k+1} + \frac{1}{2} \mathbf{z}_{i j}^{r,k} & \text{(ADMM-SGD)} \end{cases}$

Affined dual variable is a part of global control variate. We select η_i such that ECL matches with SVR.
$$\eta_i = 1 / (\mu E_i (K - 1))$$

$$\rho_i = 0$$

Proposed ECL-ISVR

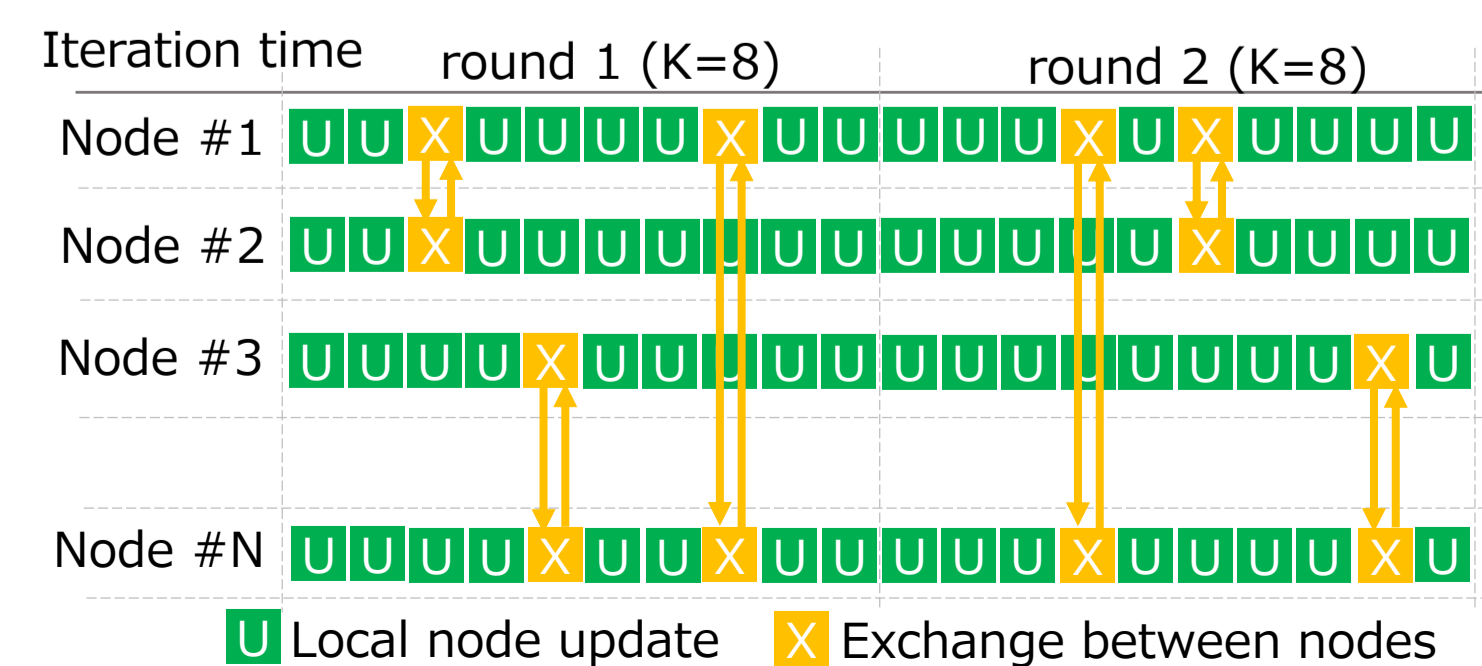
$$\begin{aligned} \mathbf{w}_i^{r,k+1} &= \mathbf{w}_i^{r,k} - \mu \{g_i(\mathbf{w}_i^{r,k}) + \bar{\mathbf{c}}_i^{r,k} - \mathbf{c}_i^{r,k}\} \\ \bar{\mathbf{c}}_i^{r,k} &= -\frac{1}{\mu K E_i} \sum_{j \in \mathcal{E}_i} \{ \mathbf{A}_{i|j}^T \mathbf{z}_{i|j}^{r,k} - \frac{1}{K} (\mathbf{A}_{i|j}^T \mathbf{z}_{i|j}^{r,k-1} + (1 - \frac{1}{K}) \mathbf{A}_{i|j}^T \mathbf{z}_{i|j}^{r,k-2} + \dots + (1 - \frac{1}{K})^{(r-1)K+k-1} \mathbf{A}_{i|j}^T \mathbf{z}_{i|j}^{1,0}) \}, \\ \mathbf{c}_i^{r,k} &= \frac{1}{K} \{g_i(\mathbf{w}_i^{r,k}) + (1 - \frac{1}{K}) g_i(\mathbf{w}_i^{r,k-1}) + (1 - \frac{1}{K})^2 g_i(\mathbf{w}_i^{r,k-2}) + \dots + (1 - \frac{1}{K})^{(r-1)K+k} g_i(\mathbf{w}_i^{1,0}) \}. \\ \mathbf{y}_{i|j}^{r,k+1} &= \mathbf{z}_{i|j}^{r,k} - 2\mathbf{A}_{i|j} \mathbf{w}_i^{r,k+1} \\ \mathbf{z}_{i|j}^{r,k+1} &= \begin{cases} \mathbf{y}_{j|i}^{r,k+1} & \text{(PDMM-ISVR)} \\ \frac{1}{2} \mathbf{y}_{j|i}^{r,k+1} + \frac{1}{2} \mathbf{z}_{i|j}^{r,k} & \text{(ADMM-ISVR)} \end{cases} \end{aligned}$$

Algorithm 1 Proposed ECL-ISVR

```

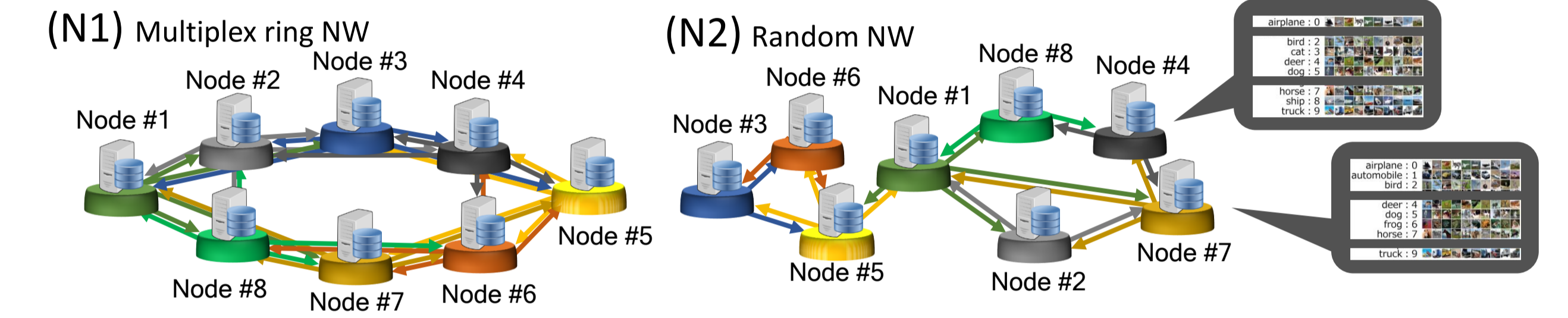
1: Set  $\mathbf{w}_i = \mathbf{w}_j \sim \text{Norm}$ ,  $\mathbf{z}_{i|j} = \mathbf{0}$ ,  $\mu$ ,  $\mathbf{x}_i$ ,  $\mathbf{A}_{i|j}$ 
2: for  $r \in \{1, \dots, R\}$  (Outer loop round) do
3:   for  $i \in \mathcal{N}$  do
4:     for  $k \in \{1, \dots, K\}$  (Inner loop iteration) do
5:       Stochastic gradient calculation
6:        $g_i(\mathbf{w}_i) \leftarrow \nabla f_i(\mathbf{w}_i, \mathbf{x}_i)$ 
7:       Update local primal and lifted dual variables
8:        $\mathbf{w}_i \leftarrow \mathbf{w}_i - \mu [g_i(\mathbf{w}_i) + \frac{1}{\mu K E_i} \{ \sum_{j \in \mathcal{E}_i} (\mathbf{w}_i - \mathbf{A}_{i|j}^T \mathbf{z}_{i|j}) - \mu E_i g_i(\mathbf{w}_i) \}]$ 
9:       for  $j \in \mathcal{E}_i$  do
10:         $\mathbf{y}_{j|i} \leftarrow \mathbf{z}_{i|j} - 2\mathbf{A}_{i|j} \mathbf{w}_i$ 
11:      end for
12:      Procedure when communicated with  $j$ -th node
13:      for  $j \in \mathcal{E}_i^c$  (at random time) do
14:        communicate  $\mathbf{e}_{j \rightarrow i}(\mathbf{y}_{j|i})$ 
15:         $\mathbf{z}_{i|j} \leftarrow \begin{cases} \mathbf{y}_{j|i} & \text{(PDMM-ISVR)} \\ \frac{1}{2} \mathbf{y}_{j|i} + \frac{1}{2} \mathbf{z}_{i|j} & \text{(ADMM-ISVR)} \end{cases}$ 
16:      end for
17:    end for
18:  end for
19: end for
    
```

- 2 algorithm flavors are existed.
1: PDMM-ISVR: Peaceman-Rachford Splitting is applied.
2: ADMM-ISVR: Douglas-Rachford Splitting is applied.
- Procedure is composed of alternatingly repeating **U** local node updates and **X** asynchronously exchange dual variables.



Numerical experiments

- Aim of experiments is to identify algorithms that nearly reach performance of **reference** case where all data are available on a single node.
- **Decentralized networks**: (N1) multiplex ring and (N2) random topologies
- **Asynchronous communication**: Once per $K=8$ inner iterations on average
- **Heterogeneous data**: (T1) fashion MNIST and (T2) CIFAR-10 is divided to $N=8$ nodes where each node has 8 classes out of a total of 10 classes.



- Proposed methods (**PDMM-ISVR** and **ADMM-ISVR**) performed closest to the single-node **reference** scores with fast processing time.
- Previous ECL (**PDMM-SGD**) was next best. However, **ADMM-SGD** was unstable with long processing time due to doubled communication requirement.

