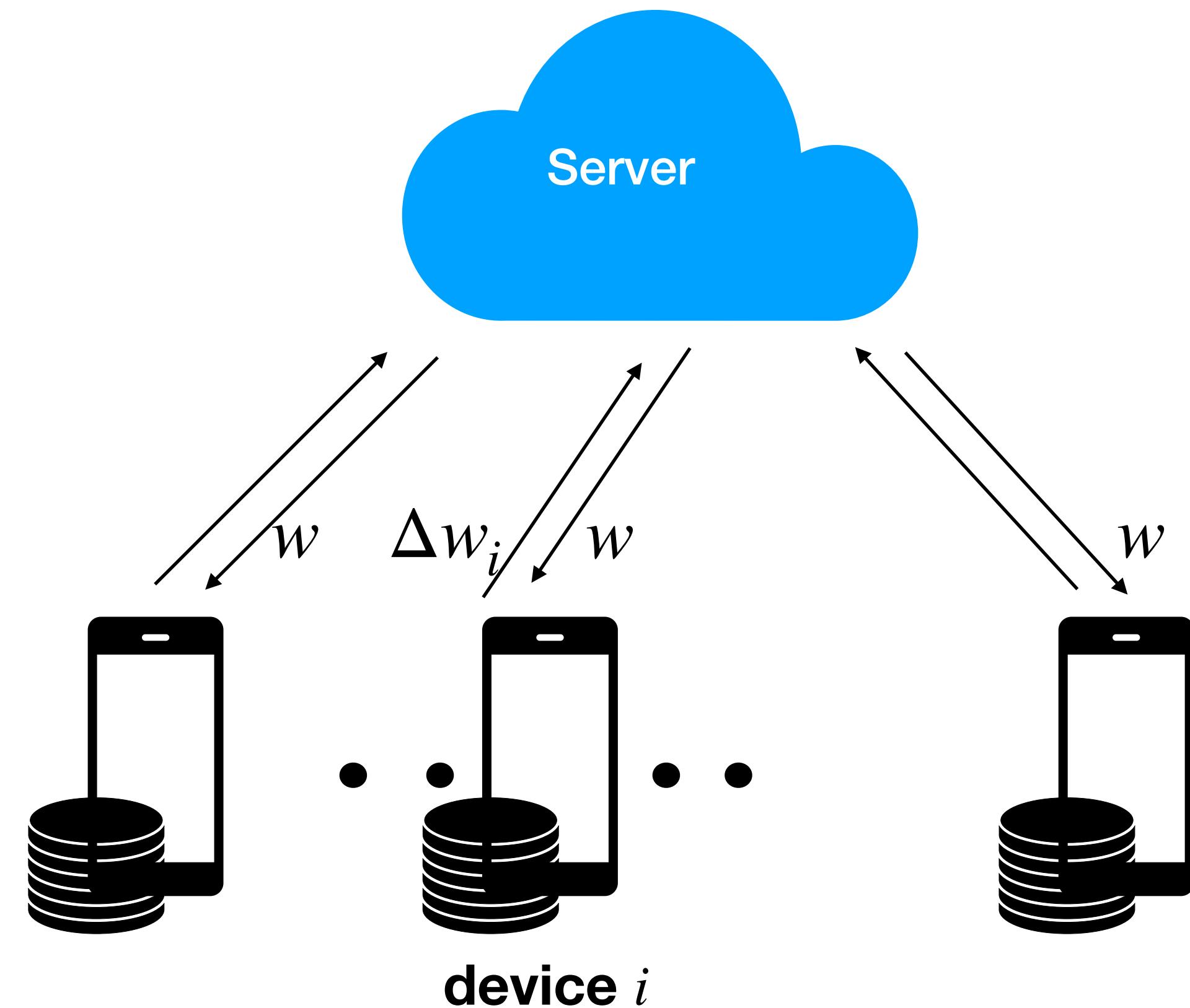


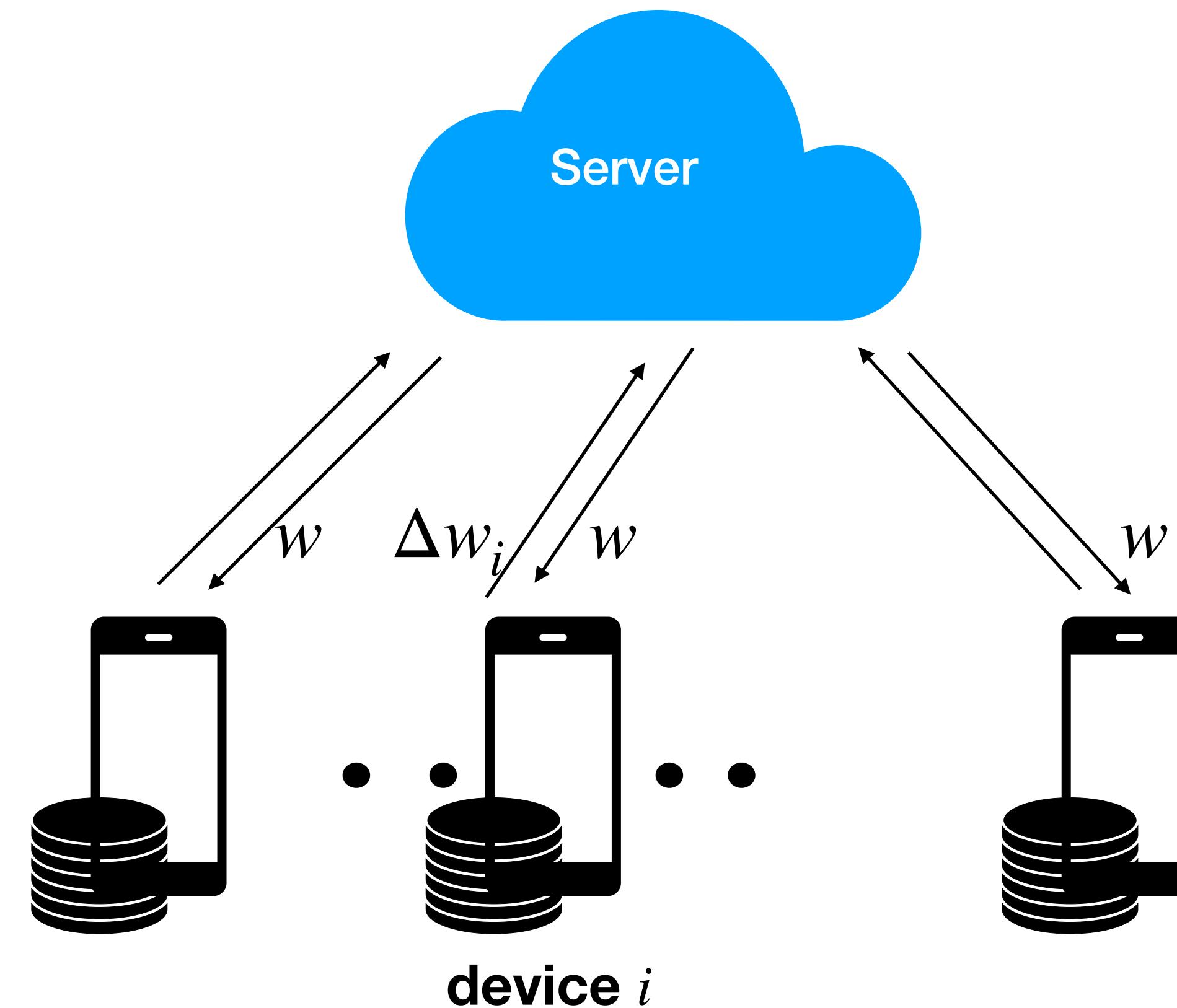
# Accelerated & Adaptive Federated Multi-task Learning

Yang Zhang, Pratik Patil, Kin Gutierrez

# Federated Learning



# Federated Learning



Why? ✓Relief on network ✓privacy ✓flexible addition of new data ✓model aggregation

# Challenges of Federated Learning

Non-IID

Unbalanced

Underlying structure among nodes

*Statistical challenges*

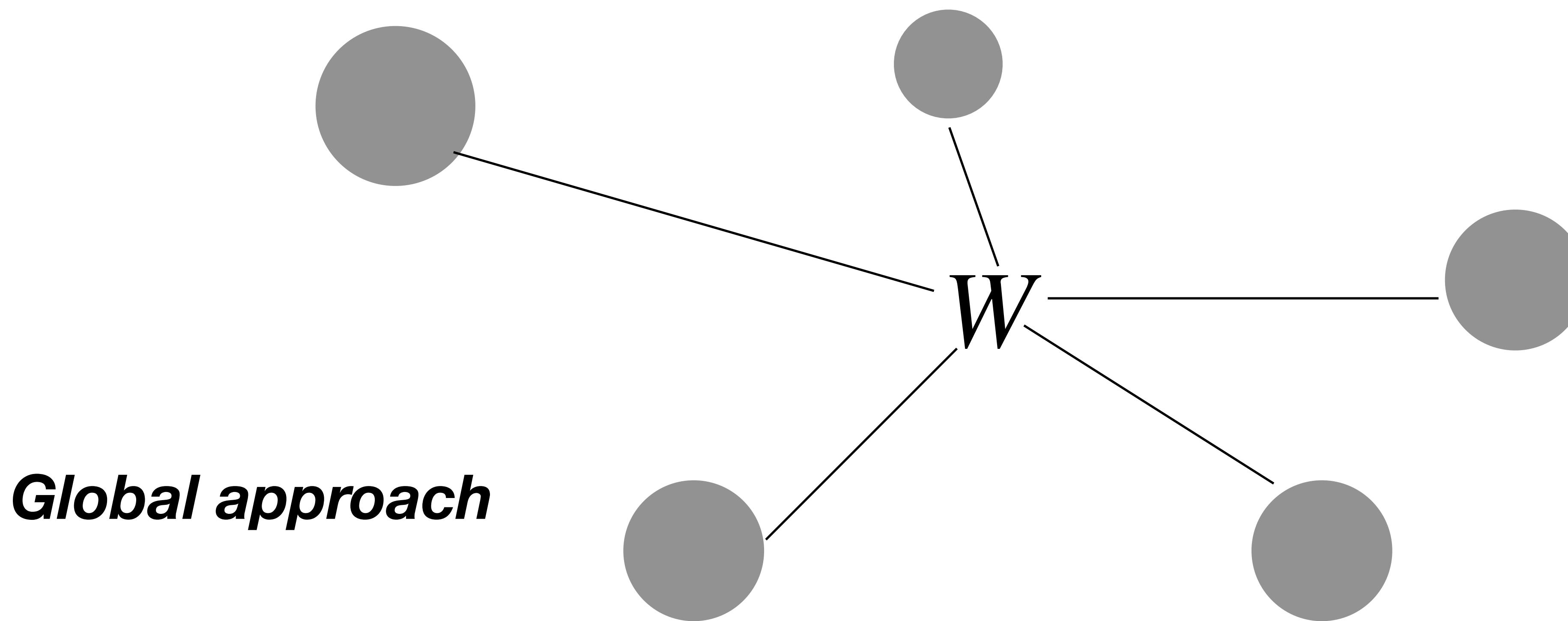
Massively distributed

Limited communication

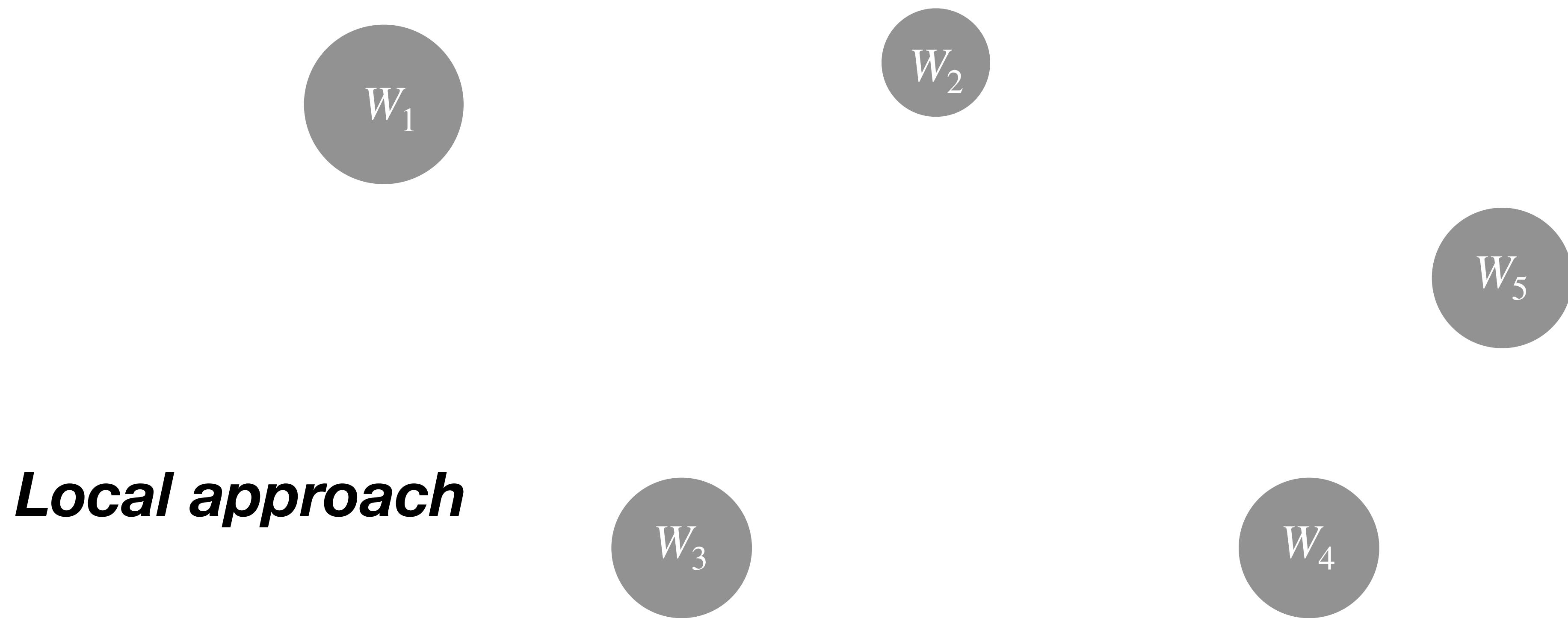
Straggler effect

*System challenges*

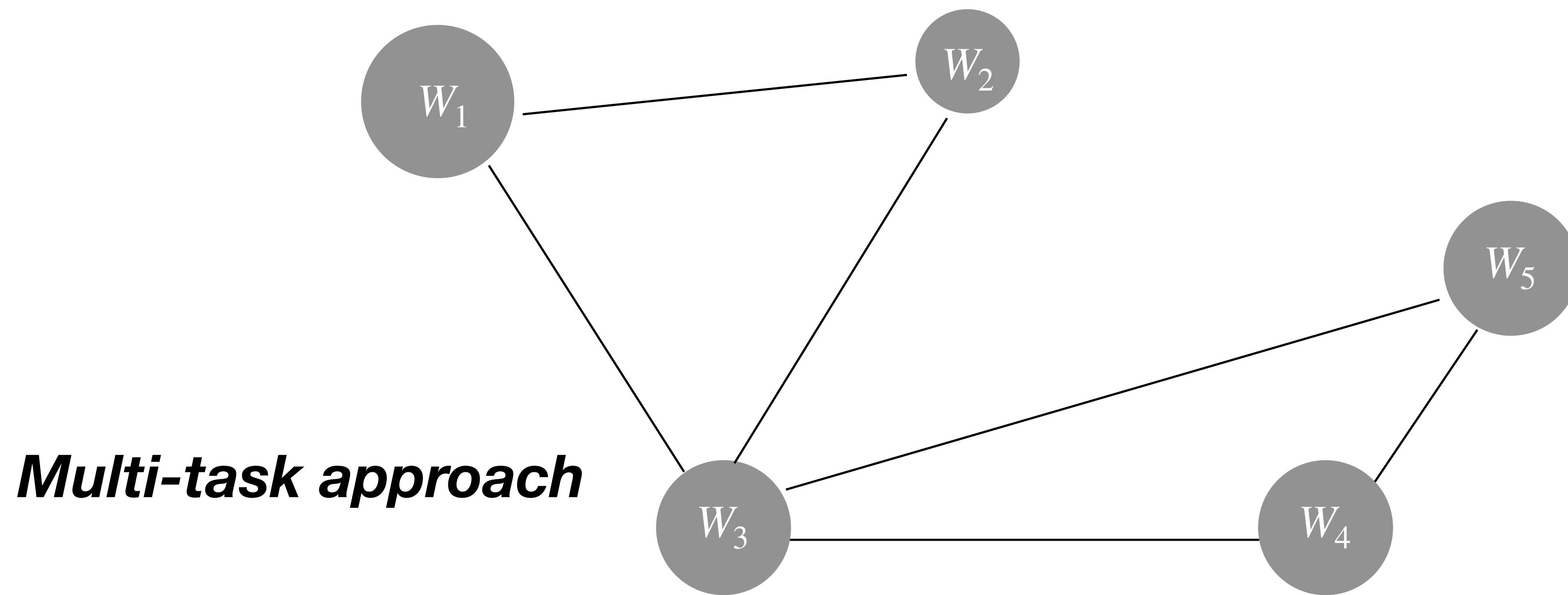
# Why Multi-task Learning?



# Why Multi-task Learning?



# Multi-task Learning Uses Underlying Structure



# Mocha Framework: Multi-task Federated Learning

**Primal**

$$\min_{\mathbf{W}, \Omega} \sum_{k=1}^K \sum_{j=1}^{n_k} \ell_k \left( \mathbf{w}_k^T \mathbf{x}_{k,j}, y_{k,j} \right) + \mathcal{R}(\mathbf{W}, \Omega)$$

# Mocha Framework: Multi-task Federated Learning

**Primal**

$$\min_{\mathbf{W}, \Omega} \sum_{k=1}^K \sum_{j=1}^{n_k} \ell_k \left( \mathbf{w}_k^T \mathbf{x}_{k,j}, y_{k,j} \right) + \mathcal{R}(\mathbf{W}, \Omega)$$

*Model Parameter*

*loss*

# Mocha Framework: Multi-task Federated Learning

**Primal**

$$\min_{\mathbf{W}, \Omega} \sum_{k=1}^K \sum_{j=1}^{n_k} \ell_k \left( \mathbf{w}_k^T \mathbf{x}_{k,j}, y_{k,j} \right) + \mathcal{R}(\mathbf{W}, \Omega)$$

*Model Parameter*      *Task Relationship Matrix*

*loss*      *Regularization*

# Mocha Framework: Multi-task Federated Learning

**Primal**

$$\min_{\mathbf{W}, \Omega} \sum_{k=1}^K \sum_{j=1}^{n_k} \ell_k \left( \mathbf{w}_k^T \mathbf{x}_{k,j}, y_{k,j} \right) + \mathcal{R}(\mathbf{W}, \Omega)$$

*Model Parameter*

*Task Relationship Matrix*

*loss*

$\lambda \operatorname{tr} (\mathbf{W} \Omega \mathbf{W}^T)$

# Mocha Framework: Multi-task Federated Learning

**Primal**

$$\min_{\mathbf{W}, \Omega} \sum_{k=1}^K \sum_{j=1}^{n_k} \ell_k \left( \mathbf{w}_k^T \mathbf{x}_{k,j}, y_{k,j} \right) + \mathcal{R}(\mathbf{W}, \Omega)$$

$\geq$

**Dual**

$$\max_{\boldsymbol{\alpha}} \sum_{k=1}^K \sum_{j=1}^{n_k} -\ell_k^* \left( -\alpha_{k,j} \right) - \mathcal{R}^*(\mathbf{X}\boldsymbol{\alpha})$$

# Mocha Framework: Multi-task Federated Learning

**Primal**

$$\min_{\mathbf{W}, \Omega} \sum_{k=1}^K \sum_{j=1}^{n_k} \ell_k \left( \mathbf{w}_k^T \mathbf{x}_{k,j}, y_{k,j} \right) + \mathcal{R}(\mathbf{W}, \Omega)$$

$\geq$

**Dual**

$$\max_{\boldsymbol{\alpha}} \sum_{k=1}^K \sum_{j=1}^{n_k} -\ell_k^* \left( -\alpha_{k,j} \right) - \mathcal{R}^*(\mathbf{X}\boldsymbol{\alpha})$$

# Mocha Framework: Multi-task Federated Learning

**Primal**

$$\min_{\mathbf{W}, \Omega} \sum_{k=1}^K \sum_{j=1}^{n_k} \frac{1}{n_k} \left[ 1 - \left( \mathbf{w}_k^T \mathbf{x}_{k,j} \right) y_{k,j} \right]_+ + \frac{\lambda}{2} \mathbf{w}^T \bar{\Omega} \mathbf{w}$$

$\geq$

**Dual**

$$\max_{\boldsymbol{\alpha}} \sum_{k=1}^K \sum_{j=1}^{n_k} y_{k,j} \alpha_{k,j} - \frac{1}{2\lambda} \boldsymbol{\alpha}^T \mathbf{X}^T \bar{\Omega}^{-1} \mathbf{X} \boldsymbol{\alpha}$$

# Dual Framework

**Global objective** 
$$\max_{\alpha} \sum_{k=1}^K \sum_{j=1}^{n_k} y_{k,j} \alpha_{k,j} - \frac{1}{2\lambda} \alpha^T \mathbf{X}^T \bar{\Omega}^{-1} \mathbf{X} \alpha$$

# Dual Framework

**Global objective** 
$$\max_{\alpha} \sum_{k=1}^K \sum_{j=1}^{n_k} y_{k,j} \alpha_{k,j} - \frac{1}{2\lambda} \alpha^T \mathbf{X}^T \bar{\Omega}^{-1} \mathbf{X} \alpha$$

# Dual Framework

**Global objective**

$$\max_{\alpha} \sum_{k=1}^K \sum_{j=1}^{n_k} y_{k,j} \alpha_{k,j} - \frac{1}{2\lambda} \alpha^T \mathbf{X}^T \bar{\Omega}^{-1} \mathbf{X} \alpha$$

**Local objective**

$$\begin{aligned} & \max_{\alpha_k^+} \sum_{j=1}^{n_k} y_{k,j} \alpha_{k,j} - \left\langle \mathbf{w} (\alpha^{(t)})_k, \mathbf{X}_k (\alpha_k^{(t+1)} - \alpha_k^{(t)}) \right\rangle \\ & - \frac{\sigma'}{2\lambda} \left\| \mathbf{X}_k (\alpha_k^{(t+1)} - \alpha_k^{(t)}) \right\|_{\bar{\Omega}_k^{-1}}^2 - \frac{1}{K} \mathcal{R}^* (\mathbf{X} \alpha) \end{aligned}$$

# Mocha: Primal-Dual Framework

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over K nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7       returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

---

# Mocha: Primal-Dual Framework

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over K nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7       returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

---

**Communication Round**

# Mocha: Primal-Dual Framework

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over K nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7       returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

---

*Communication Round*

*Local Computation*

# Mocha: Primal-Dual Framework

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

1 **for**  $i = 0, 1, \dots$  **do**  
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$   
3   **for**  $t = 0, 1, \dots T_i$  **do**  
4     Send  $\mathbf{w}_k^{(t)}$  to each node  
5     **for**  $k = 1, \dots, K$  *in parallel over K nodes* **do**  
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$   
7       returning approximate solution  $\alpha_k^{(t+1)}$   
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server  
9     **end**  
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server  
11   **end**  
12   Update  $\Omega$  centrally based on latest  $\mathbf{w}$   
13 **end**  
14 **return**  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$

---

*Communication Round*

*Local Computation*

# Mocha Drawbacks

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over K nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7         returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

---

**Intensive Local Computation**

# Mocha Drawbacks

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over K nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7       returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

---

**Intensive Local Computation**  
-> **naive acceleration**

# Mocha Drawbacks

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over K nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7       returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

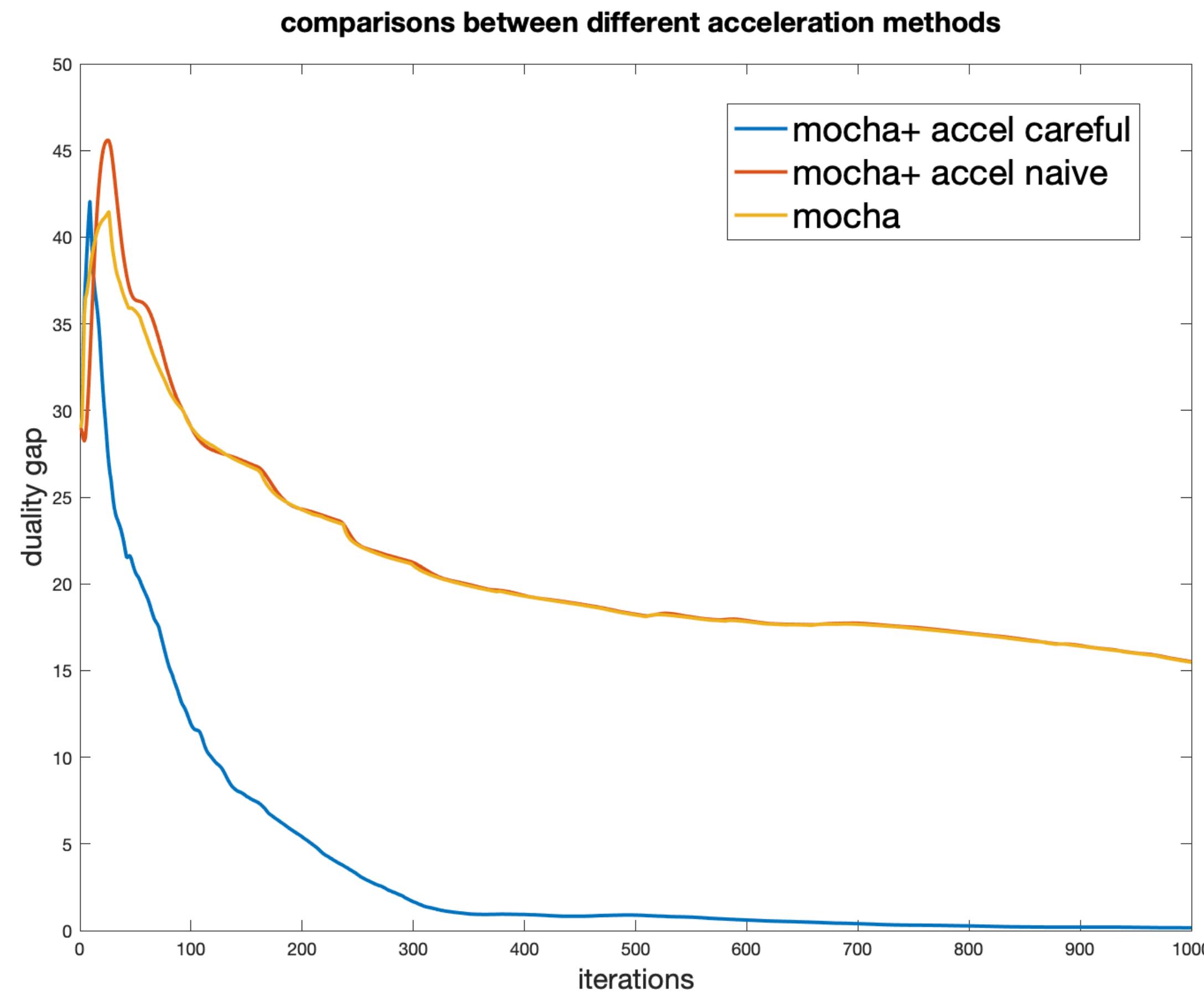
---

**Intensive Local Computation**

-> careful acceleration

# Why Acceleration?

Faster local computation



# Mocha Drawbacks

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over  $K$  nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7       returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

---

**High Communication**

# Mocha Drawbacks

---

**Algorithm 1:** Vanilla Mocha: Primal Dual Framework

---

**Input:** Data  $\mathbf{X}_k$  from  $k = 1, \dots, K$  tasks, initial matrix  $\Omega_0$   
Starting point  $\alpha^{(0)} := \mathbf{0} \in \mathbb{R}^n$ ,  $\mathbf{w}^{(0)} := \mathbf{0} \in \mathbb{R}^{K \times d}$

```
1 for  $i = 0, 1, \dots$  do
2   Set subproblem  $\sigma'$  and number of federated iterations  $T_i$ 
3   for  $t = 0, 1, \dots T_i$  do
4     Send  $\mathbf{w}_k^{(t)}$  to each node
5     for  $k = 1, \dots, K$  in parallel over K nodes do
6       Solve local dual problem  $\mathcal{D}_k(\alpha_k^{(t+1)}; \mathbf{X}_k, \alpha_k^{(t)}, \mathbf{w}_k^{(t)}, \Omega_{(k,k)}^{-1}, \sigma')$ 
7       returning approximate solution  $\alpha_k^{(t+1)}$ 
8       return  $\mathbf{X}_k \alpha_k^{(t+1)}$  to server
9     end
10    Compute  $\mathbf{w}^{(t+1)} = \nabla \mathcal{R}^\star(\mathbf{X} \alpha^{(t+1)})$  on server
11  end
12  Update  $\Omega$  centrally based on latest  $\mathbf{w}$ 
13 end
14 return  $\mathbf{W} := [\mathbf{w}_1, \dots, \mathbf{w}_K]$ 
```

---

**High Communication**

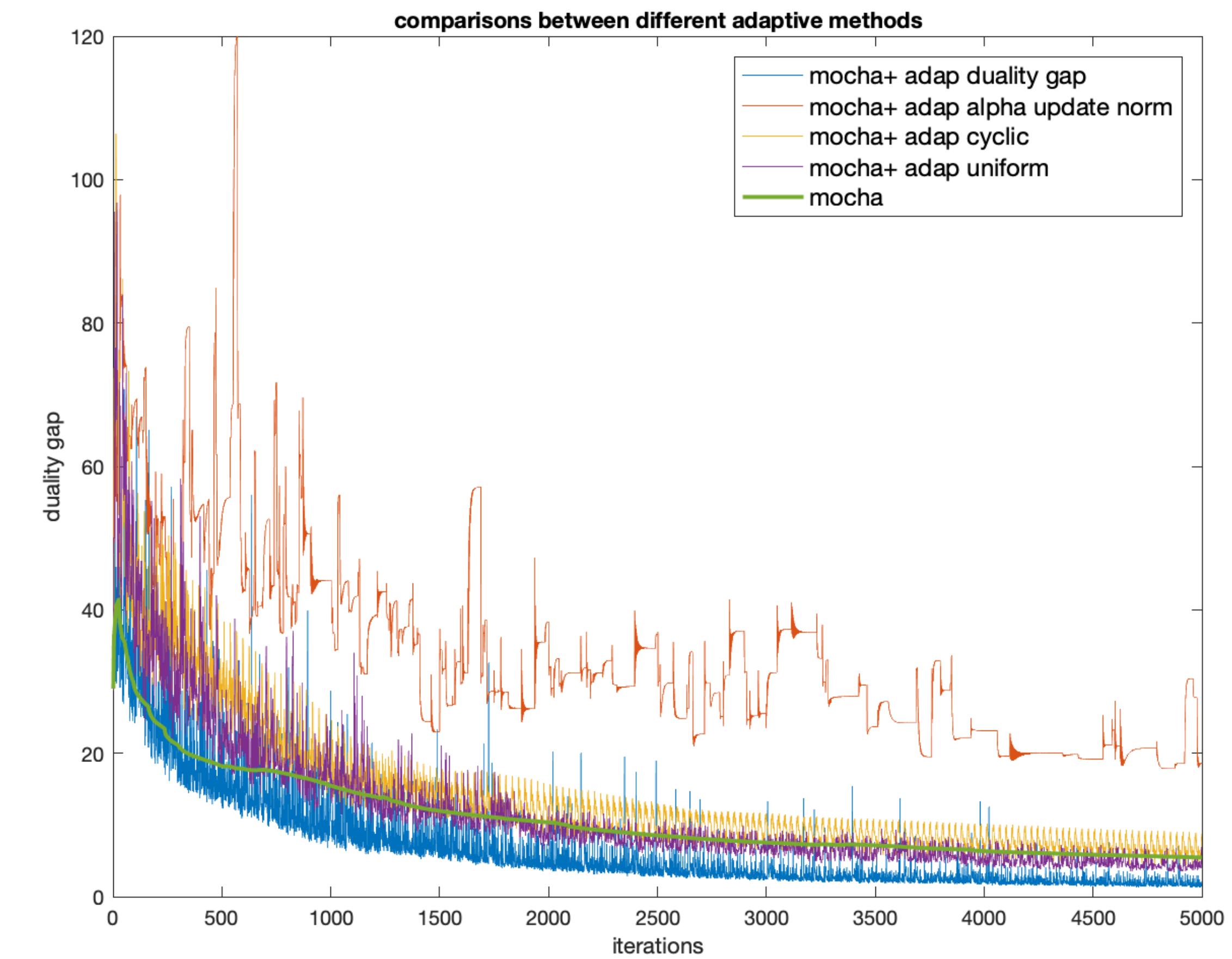
-> **adaptive node selection  
using duality gap**

# Why Adaptiveness?

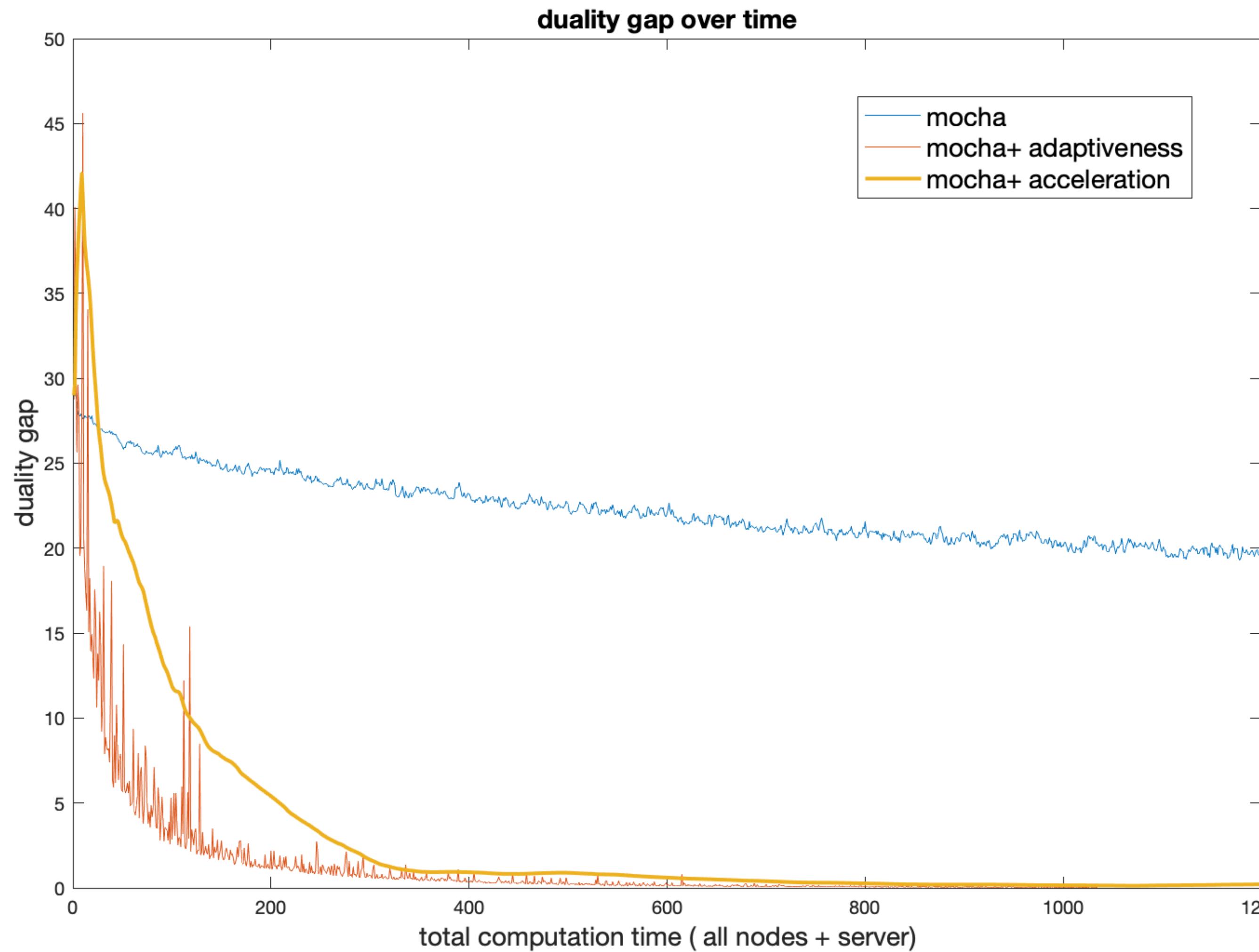
Less communication

-> Faster iterations

-> more tolerant to stragglers



# Total Computation Comparison



# Summary & Future Work

Federated multi-task learning framework

Extension:

- ✓ Adaptive node selection -> lower communication cost
- ✓ Accelerated local computation -> lower computation cost
- Acceleration on top of adaptiveness