# Permutation Compressors

Our team:
- Vyacheslav Naumov
- Anastasia Yagnych
- Mmesomachi Nwachukwu
- Alina Nurysheva

For $d \geq n$ assume $d = qn$. Let $(\pi_1, \pi_2, \ldots, \pi_d)$ be a random permutation of $(1, 2, \ldots, d)$ then for $i \in \{1, \ldots, n\}$

$$C_i(x) := n \cdot \sum_{j=q(i-1)+1}^{qi} x_{\pi_j} e_{\pi_j}$$

for $d < n$ and $n = qd$ let $(\pi_1, \pi_2, \ldots, \pi_n)$ be a random permutation of $(1, \ldots, 1, 2, \ldots, 2, \ldots, d, \ldots, d)$ where each appears $q$ times, then

$$C_i(x) := dx_{\pi_j} e_{\pi_j}$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \eta C_k(\nabla f(\mathbf{x}_k)),$$

where:

- $\mathbf{x}_k$ — parameter vector at iteration $k$,

- $\eta > 0$ — step size,

- $C_k$ — compression operator,

- $\nabla f(\mathbf{x}_k)$ — gradient of the objective function $f$.

1. **Class $B_1(\alpha, \beta)$:**

$$\alpha\|\mathbf{x}\|^2 \leq \mathbb{E}[\|C(\mathbf{x})\|^2] \leq \beta\langle\mathbb{E}[C(\mathbf{x})], \mathbf{x}\rangle.$$

2. **Class $B_2(\gamma, \beta)$:**

$$\max\left\{\gamma\|\mathbf{x}\|^2, \frac{1}{\beta}\mathbb{E}[\|C(\mathbf{x})\|^2]\right\} \leq \langle\mathbb{E}[C(\mathbf{x})], \mathbf{x}\rangle.$$

3. **Class $B_3(\delta)$:**

$$\mathbb{E}[\|C(\mathbf{x}) - \mathbf{x}\|^2] \leq \left(1 - \frac{1}{\delta}\right)\|\mathbf{x}\|^2.$$

1. For $C \in B_1(\alpha, \beta)$:

$$E_k \leq \left(1 - \frac{\alpha\mu}{\beta^2 L}\right)^k E_0,$$

where $E_k = \mathbb{E}[f(\mathbf{x}_k)] - f^*$, $L$ is the smoothness constant, $\mu$ is the strong convexity constant.

2. For $C \in B_2(\gamma, \beta)$:

$$E_k \leq \left(1 - \frac{\gamma\mu}{\beta L}\right)^k E_0.$$

3. For $C \in B_3(\delta)$:

$$E_k \leq \left(1 - \frac{\mu}{\delta L}\right)^k E_0.$$

$$\min_{x \in \mathbb{R}^d} \left[ f(x) := \frac{1}{n} \sum_{i=1}^{n} f_i(x) \right]$$

where:

- $n$ is the number of workers/nodes,

- $f_i(x)$ is the loss function for the data on worker $i$,

- $x \in \mathbb{R}^d$ represents the model parameters.

# Marina algorithm

$$x^{k+1} = x^k - \gamma g^k, \qquad g^k = \frac{1}{n}\sum_{i=1}^{n} g_i^k,$$

$$g_i^{k+1} = \begin{cases} \nabla f_i(x^{k+1}) & \text{if} \quad \theta_k = 1 \\ g^k + \mathcal{C}_i^k(\nabla f_i(x^{k+1}) - \nabla f_i(x^k)) & \text{if} \quad \theta_k = 0 \end{cases}$$

1: **Input:** starting point $x^0$, stepsize $\gamma$, probability $p \in (0, 1]$, number of iterations $T$
2: Initialize $g^0 = \nabla f(x^0)$
3: **for** $k = 0, 1, \ldots, T - 1$ **do**
4:     Sample $\theta_t \sim \mathrm{Be}(p)$
5:     Broadcast $g^t$ to all workers
6:     **for** $i = 1, \ldots, n$ in parallel **do**
7:         $x^{t+1} = x^t - \gamma g^t$
8:         Set $g_i^{t+1} = \nabla f_i(x^{t+1})$ if $\theta_t = 1$, and $g_i^{t+1} = g^t + \mathcal{C}_i \left( \nabla f_i(x^{t+1}) - \nabla f_i(x^t) \right)$ otherwise
9:     **end for**
10:     $g^{t+1} = \frac{1}{n} \sum_{i=1}^n g_i^{t+1}$
11: **end for**
12: **Output:** $\hat{x}^T$ chosen uniformly at random from $\{x^t\}_{k=0}^{T-1}$

## Communication Complexity

$$T = \mathcal{O}\left(\frac{\Delta_0}{\epsilon}\left(L_- + L_+\sqrt{\frac{1-p}{p} \cdot \frac{\omega}{n}}\right)\right),$$
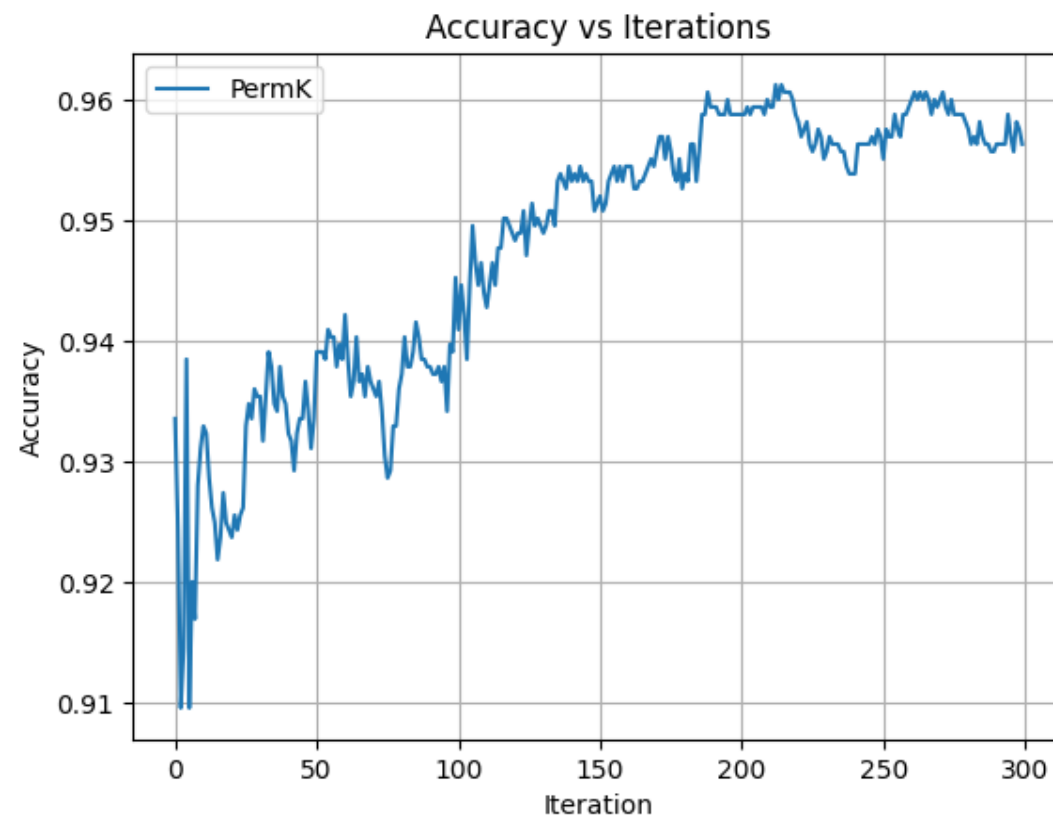
where:

- $\Delta_0 = f(x_0) - f^*$,

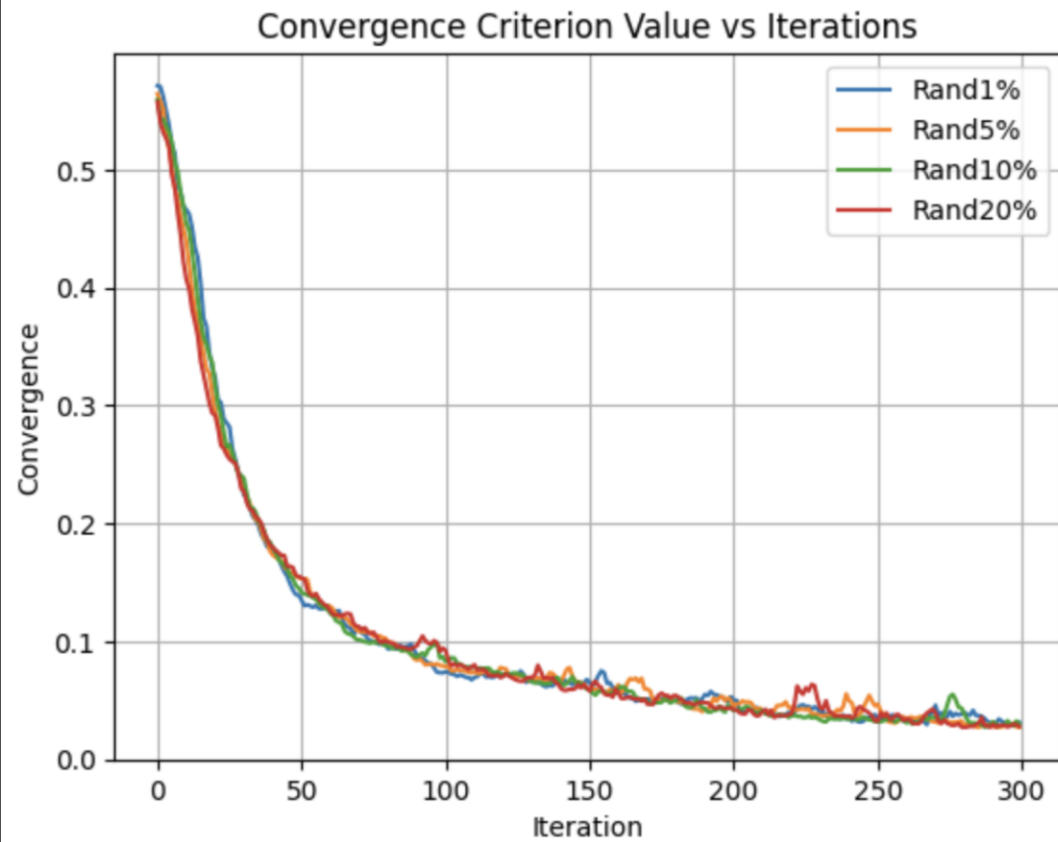- $L_-$ and $L_+$ are gradient smoothness constants.
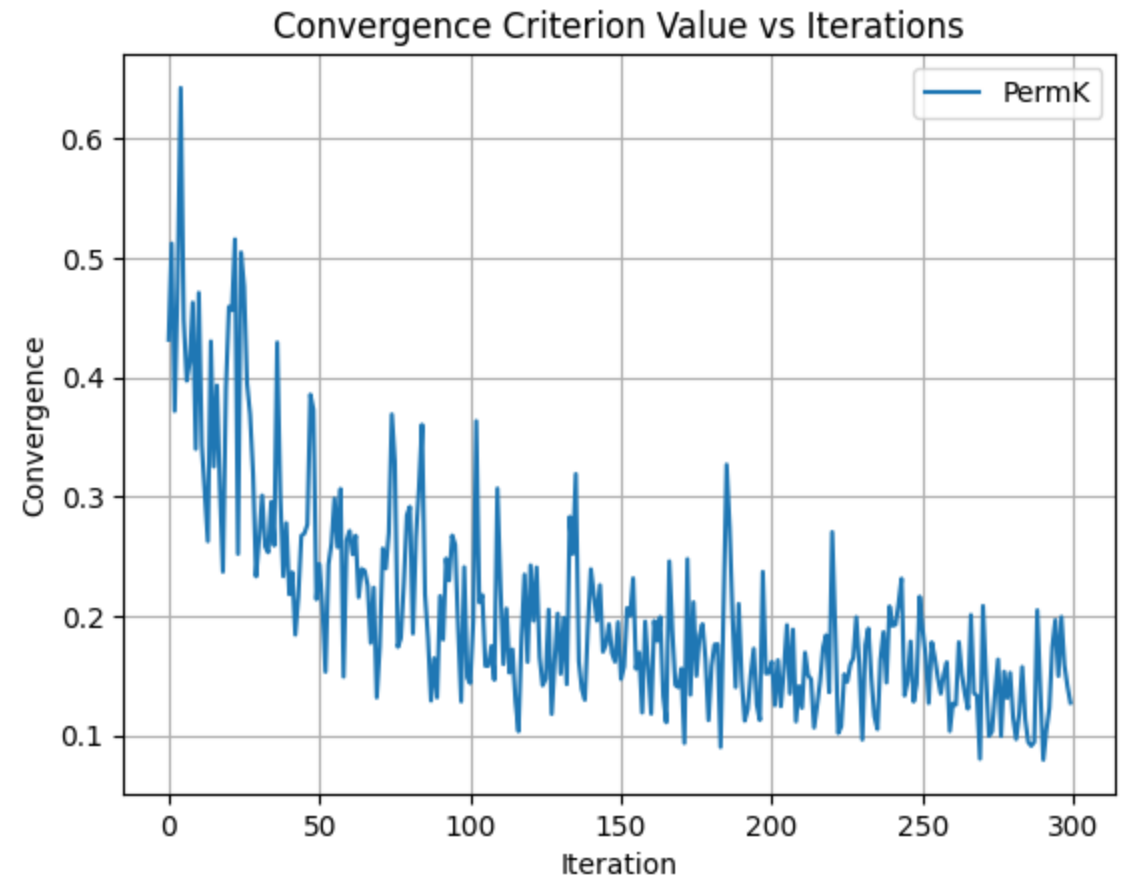
# MARINA logistic regression

Rand

Perm

# Rand

Convergence Criterion Value vs Iterations
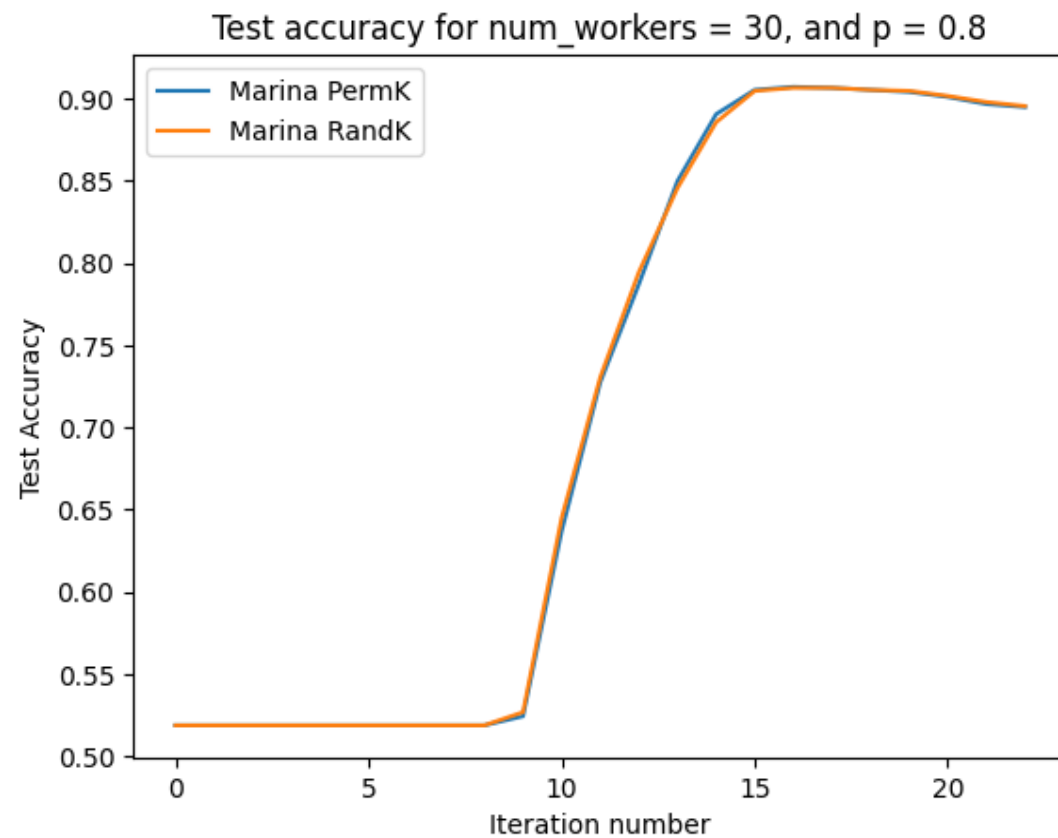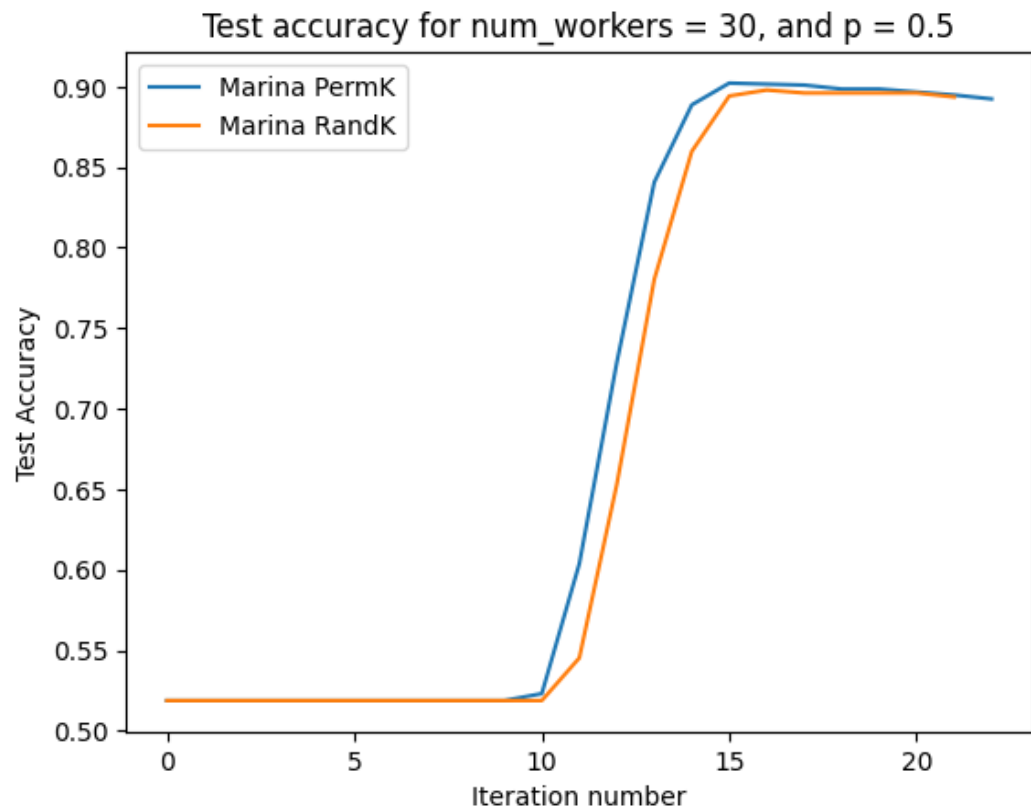
# Perm

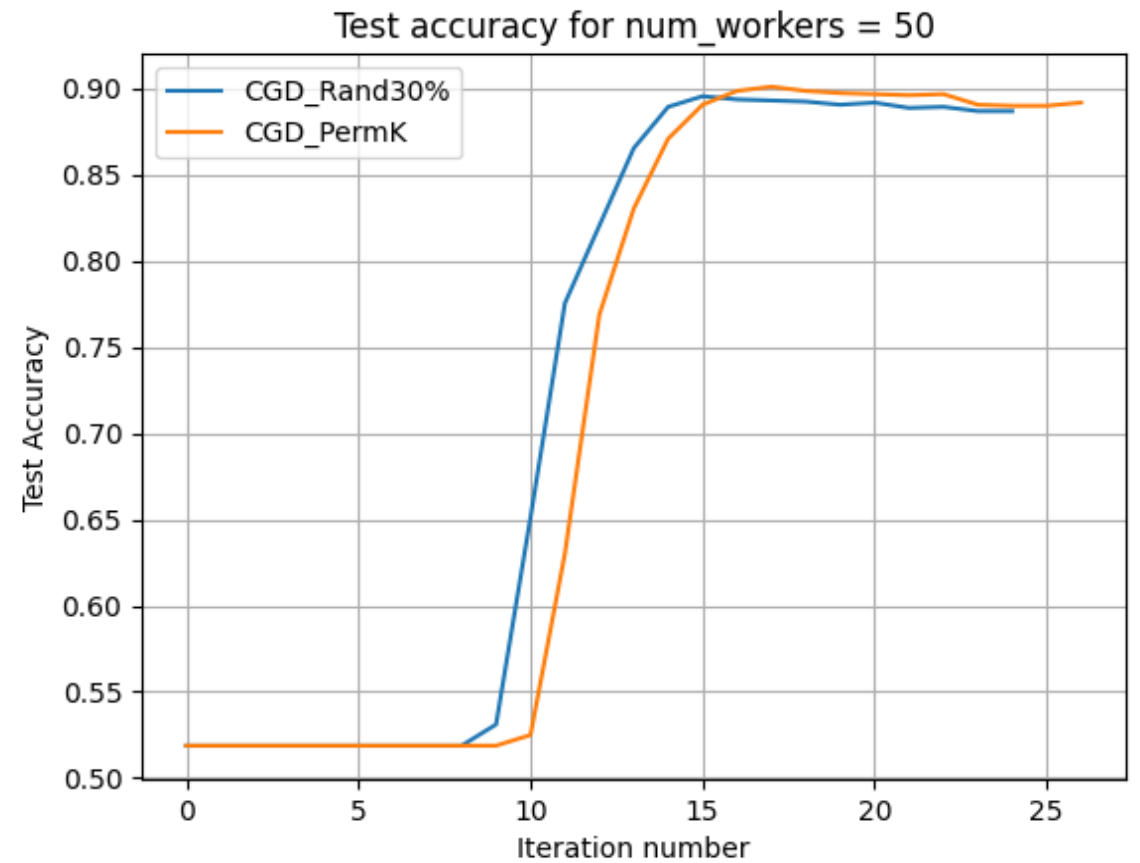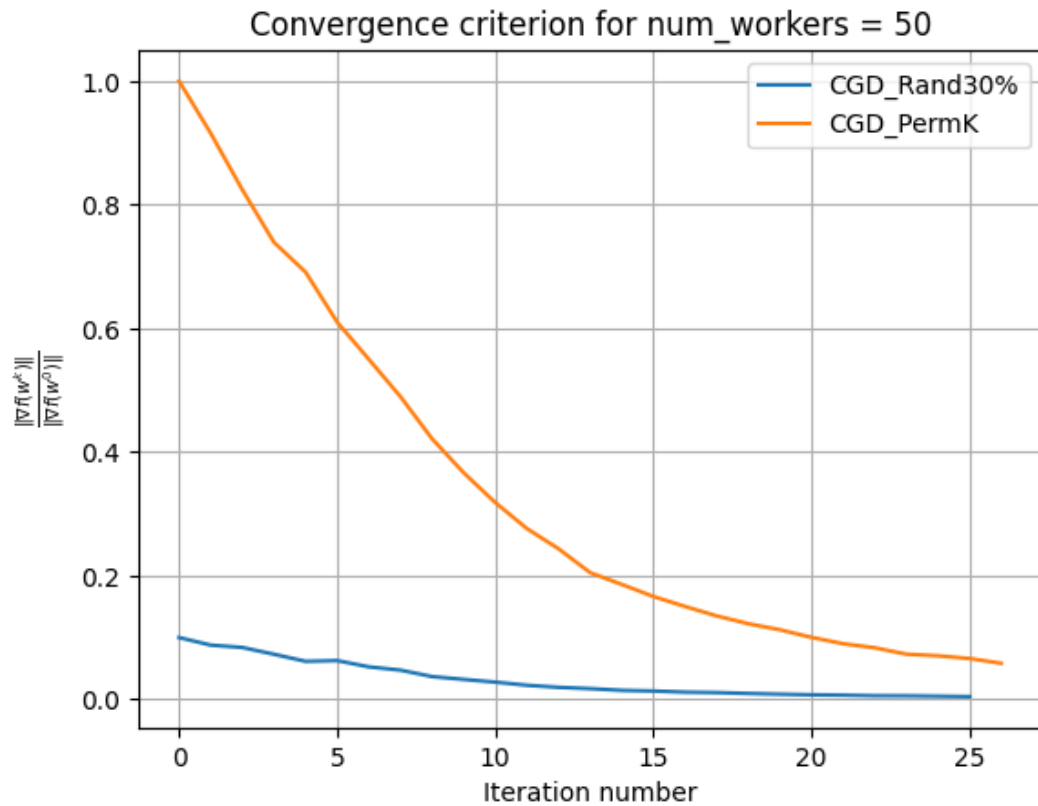Convergence Criterion Value vs Iterations

# Tuning the parameter p, in Marina

# PermK vs RandK in usual CGD

Thank you!