# Supplementary Material for: Deep Learning Model Compression with Rank Reduction in Tensor Decomposition

Wei Dai, *Student Member, IEEE,* Kai Hwang, *Life Fellow, IEEE,* and Jicong Fan, *Member, IEEE*

## I. INTRODUCTION

This supplemental materials contain all detailed proofs in the original paper.

### A. Proof of Proposition 1.

*Proof.* Mathematically, the convolution using im2col can be expressed as

$$\mathcal{Y}_{(1)} = \mathcal{W}_{(1)} \cdot \text{im2col}(\mathcal{X}), \tag{1}$$

where $\text{im2col}(\mathcal{X}) \in \mathbb{R}^{chw \times h_o w_o}$ and $\mathcal{W}_{(1)} \in \mathbb{R}^{q \times chw}$ is mode-1 unfold of tensor $\mathcal{W}$.

The kernel $\mathcal{W}$ with multilinear rank of $(r_1, r_2, r_3, r_4)$ can be decomposed as $\mathcal{W} = \mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3 \times_4 U_4$. Then, we can express the mode-1 unfolding of $\mathcal{W}$ as

$$\mathcal{W}_{(1)} = U_1 \mathcal{G}_{(1)} (U_2 \otimes U_3 \otimes U_4)^{\mathsf{T}}.$$

By plugging in Eq. (1), we prove the proposition. $\mathcal{Y}_{(1)} = U_1 \mathcal{G}_{(1)} (U_2 \otimes U_3 \otimes U_4)^{\mathsf{T}} \cdot \text{im2col}(\mathcal{X})$. Q.E.D.

### B. Proof of Proposition 2.

*Proof.* Since $\mathcal{W}$ has multilinear rank of $(r_1, r_2, 1, 1)$, it can be decomposed as $\mathcal{W} = \mathcal{G} \times_1 U_1 \times_2 U_2 \times_3 U_3 \times_4 U_4$ with $\mathcal{G} \in \mathbb{R}^{r_1 \times r_2 \times 1 \times 1}$, $U_1 \in \mathbb{R}^{q \times r_1}$, $U_2 \in \mathbb{R}^{c \times r_2}$, $U_3 \in \mathbb{R}^{1 \times 1}$, and $U_4 \in \mathbb{R}^{1 \times 1}$, we can see $U_3$ and $U_4$ are essentially scalar. By setting them to 1, the decomposition can be simplified as

$$\mathcal{W} = \mathcal{G} \times_1 U_1 \times_2 U_2.$$

After tensor mode-1 unfolding and apply Eq. 2, we have $y = U_1 \mathcal{G}_{(1)} U_2^{\mathsf{T}} x$. Q.E.D.

### C. Proof of Lemma 1.

*Proof.* We summarize the update rule of the proposed scheme in the following.

$$\begin{aligned}
\mathcal{W}^t &= g(\mathcal{H}^t), \\
\hat{\mathcal{W}}^t &= \mathcal{W}^t - \eta_1 \nabla l(\mathcal{W}^t), \\
\mathcal{H}^{t+1} &= c(\hat{\mathcal{W}}^t) - \eta_2 \nabla l(c(\hat{\mathcal{W}}^t)), \\
\mathcal{W}^{t+1} &= g(\mathcal{H}^{t+1}).
\end{aligned} \tag{2}$$

Wei Dai and Kai Hwang are with The Chinese University of Hong Kong, Shenzhen. E-mail: weidai@link.cuhk.edu.cn, hwangkai@cuhk.edu.cn.

Jicong Fan is with The Chinese University of Hong Kong, Shenzhen, and also with Shenzhen Research Institute of Big data, Shenzhen, Guangdong, China. E-mail: fanjicong@cuhk.edu.cn.

By rewriting the update rule in Eq. (2) as

$$\mathcal{W}^{t+1} = \mathcal{W}^t - \eta_1 \nabla l(\mathcal{W}^t) + \mathcal{E}^t,$$

where $\mathcal{E}^t = g(\mathcal{H}^{t+1}) - \mathcal{W}^t + \eta_1 \nabla l(\mathcal{W}^t)$ denotes the low-rank update error on the $t$-th iteration.

Then, we can bound the low-rank update error $\mathcal{E}^t$. We first explicitly derive $\|\mathcal{E}^t\|_F$ as the following.

$$\begin{aligned}
\|\mathcal{E}^t\|_F &= \|\mathcal{E}^t_{(1)}\|_F \\
&= \|g(\mathcal{H}^{t+1})_{(1)} - \mathcal{W}^t_{(1)} + \eta_1 \nabla l(\mathcal{W}^t)_{(1)}\|_F.
\end{aligned}$$

Let

$$\begin{aligned}
&g(\mathcal{H}^{t+1})_{(1)} \\
=&g\left(c(\hat{\mathcal{W}}^t) - \eta_2 \nabla l(c(\hat{\mathcal{W}}^t))\right)_{(1)} \\
=&\left(\hat{U}_1^t - \eta_2 \nabla l(\hat{U}_1^t)\right) \left(\hat{\mathcal{G}}^t - \eta_2 \nabla l(\hat{\mathcal{G}}^t)\right)_{(1)} \\
&\left(\left(\hat{U}_2^t - \eta_2 \nabla l(\hat{U}_2^t)\right) \otimes \left(\hat{U}_3^t - \eta_2 \nabla l(\hat{U}_3^t)\right) \otimes \right. \\
&\left. \left(\hat{U}_4^t - \eta_2 \nabla l(\hat{U}_4^t)\right)\right)^{\mathsf{T}} \\
=&\hat{U}_1^t \hat{\mathcal{G}}_{(1)}^t (\hat{U}_2^t \otimes \hat{U}_3^t \otimes \hat{U}_4^t)^{\mathsf{T}} + R^t \\
=&g(c(\hat{\mathcal{W}}^t))_{(1)} + E^t - E^t + R^t \\
=&\hat{\mathcal{W}}_{(1)}^t + R^t - E^t \\
=&\mathcal{W}_{(1)}^t - \eta_1 \nabla l(\mathcal{W}^t)_{(1)} + R^t - E^t,
\end{aligned}$$

where $E^t = \hat{\mathcal{W}}_{(1)}^t - g(c(\hat{\mathcal{W}}^t))_{(1)}$, and

$$\begin{aligned}
R^t = &- \eta_2 \nabla l(\hat{U}_1^t)\hat{\mathcal{G}}_{(1)}^t (\hat{U}_2^t \otimes \hat{U}_3^t \otimes \hat{U}_4^t)^{\mathsf{T}} - ... \\
&+ \eta_2^2 \nabla l(\hat{U}_1^t) \nabla l(\hat{\mathcal{G}}^t)_{(1)} (\hat{U}_2^t \otimes \hat{U}_3^t \otimes \hat{U}_4^t)^{\mathsf{T}} + ... \\
&- \eta_2^3 \nabla l(\hat{U}_1^t) \nabla l(\hat{\mathcal{G}}^t)_{(1)} (\nabla l(\hat{U}_2^t) \otimes \hat{U}_3^t \otimes \hat{U}_4^t)^{\mathsf{T}} - ... \\
&+ \eta_2^4 \nabla l(\hat{U}_1^t) \nabla l(\hat{\mathcal{G}}^t)_{(1)} (\nabla l(\hat{U}_2^t) \otimes \nabla l(\hat{U}_3^t) \otimes \hat{U}_4^t)^{\mathsf{T}} + ... \\
&- \eta_2^5 \nabla l(\hat{U}_1^t) \nabla l(\hat{\mathcal{G}}^t)_{(1)} (\nabla l(\hat{U}_2^t) \otimes \nabla l(\hat{U}_3^t) \otimes \nabla l(\hat{U}_4^t))^{\mathsf{T}},
\end{aligned}$$

consists of 31 terms that are the permutation of low-rank weights and their gradients. Then, plugging it back, we have

$$\|\mathcal{E}_{(1)}^t\|_F = \|R^t - E^t\|_F \leq \|R^t\|_F + \|E^t\|_F.$$

According to Equation (8) in original paper, $\|E^t\|_F$ is bounded that

$$\|E^t\|_F \leq (1 - \rho)\|\hat{\mathcal{W}}^t\|_F \leq (1 - \rho)\varphi.$$

Using the assumption (2) and (3), if $0 \leq \eta_2 \leq 1$, $\mathbb{E}\left[\|\boldsymbol{R}^t\|\|_F$ is bounded that

$$
\begin{aligned}
&\mathbb{E}\left[\|\boldsymbol{R}^t\|_F\right] \\
\leq &\eta_2 \left( \|\nabla l(\hat{\boldsymbol{U}}_1^t)\|_2 \|\hat{\boldsymbol{\mathcal{G}}}_{(1)}^t\|_F \|\hat{\boldsymbol{U}}_2^t\|_2 \|\hat{\boldsymbol{U}}_3^t\|_2 \|\hat{\boldsymbol{U}}_4^t\|_2 + ... \right. \\
&+ \|\nabla l(\hat{\boldsymbol{U}}_1^t)\|_2 \|\nabla l(\hat{\boldsymbol{\mathcal{G}}}^t)_{(1)}\|_F \|\hat{\boldsymbol{U}}_2^t\|_2 \|\hat{\boldsymbol{U}}_3^t\|_2 \|\hat{\boldsymbol{U}}_4^t\|_2 + ... \\
&+ \|\nabla l(\hat{\boldsymbol{U}}_1^t)\|_2 \|\nabla l(\hat{\boldsymbol{\mathcal{G}}}^t)_{(1)}\|_F \|\nabla l(\hat{\boldsymbol{U}}_2^t)\|_2 \|\hat{\boldsymbol{U}}_3^t\|_2 \|\hat{\boldsymbol{U}}_4^t\|_2 + ... \\
&+ \|\nabla l(\hat{\boldsymbol{U}}_1^t)\|_2 \|\nabla l(\hat{\boldsymbol{\mathcal{G}}}^t)_{(1)}\|_F \|\nabla l(\hat{\boldsymbol{U}}_2^t)\|_2 \|\nabla l(\hat{\boldsymbol{U}}_3^t)\|_2 \|\hat{\boldsymbol{U}}_4^t\|_2 + ... \\
&+ \|\nabla l(\hat{\boldsymbol{U}}_1^t)\|_2 \|\nabla l(\hat{\boldsymbol{\mathcal{G}}}^t)_{(1)}\|_F \\
&\quad \left. \|\nabla l(\hat{\boldsymbol{U}}_2^t)\|_2 \|\nabla l(\hat{\boldsymbol{U}}_3^t)\|_2 \|\nabla l(\hat{\boldsymbol{U}}_4^t)\|_2 \right) \\
\leq &\eta_2 \left( \varphi \left( \sum_{i=1}^4 \binom{4}{i} 2^i \right) + G_2 \left( \sum_{i=0}^4 \binom{4}{i} 2^i \right) \right) \\
= &\eta_2 \left( 80\varphi + 81 G_2 \right),
\end{aligned}
$$

where the matrix norm inequality that $\|\boldsymbol{A}\boldsymbol{B}\|_F \leq \|\boldsymbol{A}\|_2 \|\boldsymbol{B}\|_F$ is applied. Then, by taking the expectation, it yields

$$
\mathbb{E}\left[\|\boldsymbol{\mathcal{E}}^t\|_F\right] \leq \eta_2 \left( 80\varphi + 81 G_2 \right) + (1-\rho)\varphi
$$

Q.E.D.

### D. Proof of Theorem 1

*Proof.* The proof outline is consistent with [1]. Consider the following.

$$
\begin{aligned}
&\mathbb{E}\left[\|\boldsymbol{\mathcal{W}}^{t+1} - \boldsymbol{\mathcal{W}}^*\|_F^2\right] \\
= &\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 - 2\mathbb{E}\langle \boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*, \eta_1 \nabla l(\boldsymbol{\mathcal{W}}^t) - \boldsymbol{\mathcal{E}}^t\rangle \\
&+ \mathbb{E}\|\eta_1 \nabla l(\boldsymbol{\mathcal{W}}^t) - \boldsymbol{\mathcal{E}}^t\|_F^2 \\
= &\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 - 2\eta_1 \langle \boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*, \nabla \mathcal{L}(\boldsymbol{\mathcal{W}}^t)\rangle \\
&+ 2\mathbb{E}\langle \boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*, \boldsymbol{\mathcal{E}}^t\rangle + \eta_1^2 \mathbb{E}\|\nabla l(\boldsymbol{\mathcal{W}}^t)\|_F^2 \\
&+ \mathbb{E}\|\boldsymbol{\mathcal{E}}^t\|_F^2 - 2\mathbb{E}\langle \nabla l(\boldsymbol{\mathcal{W}}^t), \boldsymbol{\mathcal{E}}^t\rangle \\
\leq &\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 - 2\eta_1 \langle \boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*, \nabla \mathcal{L}(\boldsymbol{\mathcal{W}}^t)\rangle \\
&+ 2\mathbb{E}\|\boldsymbol{\mathcal{E}}^t\|_F \|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F + 2\mathbb{E}\|\boldsymbol{\mathcal{E}}^t\|_F \|\nabla l(\boldsymbol{\mathcal{W}}^t)\|_F \\
&+ \eta_1^2 \mathbb{E}\|\nabla l(\boldsymbol{\mathcal{W}}^t)\|_F^2 + \mathbb{E}\|\boldsymbol{\mathcal{E}}^t\|_F^2 \\
\leq &2\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 - 2\eta_1 \langle \boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*, \nabla \mathcal{L}(\boldsymbol{\mathcal{W}}^t)\rangle \\
&+ \|\nabla l(\boldsymbol{\mathcal{W}}^t)\|_F^2 + \eta_1^2 \mathbb{E}\|\nabla l(\boldsymbol{\mathcal{W}}^t)\|_F^2 + 3\mathbb{E}\|\boldsymbol{\mathcal{E}}^t\|_F^2 \\
\leq &2\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 - 2\eta_1 \langle \boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*, \nabla \mathcal{L}(\boldsymbol{\mathcal{W}}^t)\rangle \\
&+ \eta_1^2 G_1^2 + 3\eta_2^2 a^2 + G_1^2,
\end{aligned}
$$

where Lemma 1 is applied and $a = 81 G_2 + 80\varphi + \frac{1}{\eta_2}(1-\rho)\varphi$. Using the assumption that $\mathcal{L}$ is $\mu$-strongly convex, we have

$$
\begin{aligned}
&\mathbb{E}\left[\|\boldsymbol{\mathcal{W}}^{t+1} - \boldsymbol{\mathcal{W}}^*\|_F^2\right] \\
\leq &(2 - \eta_1 L)\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 - 2\eta_1 \left( \mathcal{L}(\boldsymbol{\mathcal{W}}^t) - \mathcal{L}(\boldsymbol{\mathcal{W}}^*) \right) \\
&+ \eta_1^2 G_1^2 + 3\eta_2^2 a^2 + G_1^2, \\
\Rightarrow \quad &2\eta_1 \left( \mathcal{L}(\boldsymbol{\mathcal{W}}^t) - \mathcal{L}(\boldsymbol{\mathcal{W}}^*) \right) \\
\leq &(2 - \eta_1 L)\mathbb{E}\left[\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2\right] - \mathbb{E}\left[\|\boldsymbol{\mathcal{W}}^{t+1} - \boldsymbol{\mathcal{W}}^*\|_F^2\right] \\
&+ \eta_1^2 G_1^2 + 3\eta_2^2 a^2 + G_1^2, \\
\Rightarrow \quad &\mathbb{E}\left( \mathcal{L}(\boldsymbol{\mathcal{W}}^t) - \mathcal{L}(\boldsymbol{\mathcal{W}}^*) \right) \\
\leq &(\frac{2}{2\eta_1} - \frac{\mu}{2})\mathbb{E}\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 - \frac{1}{2\eta_1}\mathbb{E}\left[\|\boldsymbol{\mathcal{W}}^{t+1} - \boldsymbol{\mathcal{W}}^*\|_F^2\right]
\end{aligned}
$$

$$
+ \frac{\eta_1}{2} G_1^2 + \frac{3\eta_2^2}{2\eta_1} a^2 + \frac{1}{2\eta_1} G_1^2,
$$

If $\eta_1 = \frac{1}{1/2^t + \mu}$, and $\rho = 1 - \eta_2$, we have

$$
\begin{aligned}
&\mathbb{E}\left( \mathcal{L}(\boldsymbol{\mathcal{W}}^t) - \mathcal{L}(\boldsymbol{\mathcal{W}}^*) \right) \leq \left( \frac{1}{(2)^t} + \frac{\mu}{2} \right) \mathbb{E}\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 \\
&- \left( \frac{1}{2^{(t+1)}} + \frac{\mu}{2} \right) \mathbb{E}\|\boldsymbol{\mathcal{W}}^{t+1} - \boldsymbol{\mathcal{W}}^*\|_F^2 + \frac{1}{4/2^{(t+1)} + 2\mu} G_1^2 \\
&+ \left( \frac{1}{(2)^t} + \mu \right) (\frac{3}{2}\eta_2^2 b^2 + G_1^2) \\
\leq &\left( \frac{1}{(2)^t} + \frac{\mu}{2} \right) \mathbb{E}\|\boldsymbol{\mathcal{W}}^t - \boldsymbol{\mathcal{W}}^*\|_F^2 \\
&- \left( \frac{1}{2^{(t+1)}} + \frac{\mu}{2} \right) \mathbb{E}\|\boldsymbol{\mathcal{W}}^{t+1} - \boldsymbol{\mathcal{W}}^*\|_F^2 + \frac{1}{2\mu} G_1^2 \\
&+ \left( (\frac{1}{2})^t + \mu \right) (\frac{3}{2}\eta_2^2 b^2 + \frac{1}{2} G_1^2).
\end{aligned}
$$

where $b = 81(G_2 + \varphi)$. Then, setting $\eta_2 = (\frac{1}{\sqrt{C}})^t$ and applying the telescope sum from $t = 0$ to $T$, we have

$$
\begin{aligned}
&\frac{1}{T} \sum_{t=0}^T \mathbb{E}\left( \mathcal{L}(\boldsymbol{\mathcal{W}}^t) - \mathcal{L}(\boldsymbol{\mathcal{W}}^*) \right) \\
\leq &\frac{1}{T} \sum_{t=0}^T (\frac{3}{2}\eta_2^2 b^2 + \frac{1}{2} G_1^2) \left( \left( \frac{1}{2} \right)^t + \mu \right) + \frac{c}{T} \\
&- \left( \frac{1}{2^{(T+1)}} + \frac{\mu}{2} \right) \mathbb{E}\|\boldsymbol{\mathcal{W}}^{T+1} - \boldsymbol{\mathcal{W}}^*\|_F^2 + \frac{1}{2\mu} G_1^2 \\
\leq &\frac{1}{2T} \sum_{t=0}^T \left( \frac{3}{2}\eta_2^2 b^2 + \frac{1}{2} G_1^2 \right)^2 + \frac{1}{2T} \sum_{t=0}^T \left( \left( \frac{1}{2} \right)^t + \mu \right)^2 \\
&+ \frac{c}{T} + \frac{1}{2\mu} G_1^2 \\
= &\frac{9 b^4}{8T} \sum_{t=0}^T \left( \frac{1}{C^2} \right)^t + \frac{3 b^2 G_1^2}{4T} \left( \frac{1}{C} \right)^t \\
&+ \frac{1}{2T} \left( \sum_{t=0}^T (\frac{1}{4})^t + 2\mu \sum_{t=0}^T (\frac{1}{2})^t + 2c \right) + \frac{G_1^4}{8} + \frac{\mu^2}{2} + \frac{1}{2\mu} G_1^2 \\
\leq &\frac{1}{2T} \left( 4\mu + \frac{3 G_1^2 C b^2}{(C-1)} + \frac{9 C^2 b^4}{4(C^2 - 1)} + 2c + \frac{4}{3} \right) \\
&+ \frac{1}{2}(\frac{G_1^2}{\mu} + \mu^2 + \frac{G_1^4}{4}),
\end{aligned}
$$

where $c = \frac{\mu + 2}{2}\|\boldsymbol{\mathcal{W}}^0 - \boldsymbol{\mathcal{W}}^*\|_F^2$. Using Jensen's inequality, we have

$$
\mathbb{E}\left( \mathcal{L}(\bar{\boldsymbol{\mathcal{W}}}^T) - \mathcal{L}(\boldsymbol{\mathcal{W}}^*) \right) \leq \frac{1}{T} \sum_{t=0}^T \mathbb{E}\left( \mathcal{L}(\boldsymbol{\mathcal{W}}^t) - \mathcal{L}(\boldsymbol{\mathcal{W}}^*) \right),
$$

where $\bar{\boldsymbol{\mathcal{W}}}^T = \frac{1}{T} \sum_{t=1}^T = \boldsymbol{\mathcal{W}}^t$. Hence,

$$
\begin{aligned}
\mathbb{E}[\bar{\boldsymbol{\mathcal{W}}}^T - \mathcal{L}(\boldsymbol{\mathcal{W}}^*)] \leq &\frac{1}{2T} \left( 4L + \frac{3 G_1^2 U b^2}{(U-1)} + \frac{9 U^2 b^4}{4(U^2 - 1)} + \frac{4}{3} \right) \\
&+ \frac{1}{2}(\frac{G_1^2}{L} + \mu^2 + \frac{G_1^4}{4}).
\end{aligned}
$$

As $T$ converges to $\infty$, it converges to $\frac{1}{2}(\frac{G_1^2}{L} + \mu^2 + \frac{G_1^4}{4})$.

Q.E.D.

### E. Proof of Corollary 1.

*Proof.* $M \geq \frac{n^d}{r^d + dnr} \geq \frac{n^d}{r^d + dnr^d}$. If $r \leq \left(\frac{n^d}{\phi(1+dn)}\right)^{1/d}$, $\frac{n^d}{r^d + dnr^d} \geq \phi$. Therefore, $M \geq \phi$.      Q.E.D.

### F. Proof of Theorem 3

*Proof.* For the uncompressed case, the computation is essentially matrix-matrix multiplication between matrices with size $n \times n^3$ and $n^3 \times m$ according to Eq. (1). So, the computation complexity is $\Theta(mn^4)$.

In the case of low-rank convolutional layer, there is an efficient way to implement the Eq. (11). Denoting $\text{im2col}(\boldsymbol{\mathcal{X}}) = \left[\text{vec}(\mathcal{P}^{(1)}), \text{vec}(\mathcal{P}^{(2)}), ..., \text{vec}(\mathcal{P}^{(m)})\right]$ where $\text{vec}(\cdot)$ is the vectorization operation and $\mathcal{P}^{(m)} \in \mathbb{R}^{n \times n \times n}$ is the patch of image to be convolved, then by applying the matrix equations in Kronecker product [2], for each $k = 1, 2, ..., m$ we have

$$(\boldsymbol{U}_2 \otimes \boldsymbol{U}_3 \otimes \boldsymbol{U}_4)^\mathsf{T} \cdot \text{vec}(\mathcal{P}^{(k)}) = (\boldsymbol{U}_3 \otimes \boldsymbol{U}_4)^\mathsf{T}(\mathcal{P}^{(k)}_{(1)})^\mathsf{T}\boldsymbol{U}_2. \tag{3}$$

Furthermore, applying the matrix equations again, we have

$$(\boldsymbol{U}_3 \otimes \boldsymbol{U}_4)^\mathsf{T}(\mathcal{P}^{(k)}_{(1)})^\mathsf{T}_{:,i} = \boldsymbol{U}_4^\mathsf{T}(\mathcal{P}^{(k)}_{i,:,:})^\mathsf{T}\boldsymbol{U}_3, \ \forall i = 1, 2, ..., n. \tag{4}$$

Hence, to compute $(\boldsymbol{U}_2 \otimes \boldsymbol{U}_3 \otimes \boldsymbol{U}_4)^\mathsf{T} \cdot \text{im2col}(\boldsymbol{\mathcal{X}})$, we need to compute $m$ times Eq. (3) and $mn$ times Eq. (4). The computation complexity is $\Theta\left(m\left(n(n^2r + nr^2) + nr^3\right)\right)$. Then, by performing matrix multiplication like $\boldsymbol{U}_1\left(\boldsymbol{\mathcal{G}}_{(1)}\left((\boldsymbol{U}_2 \otimes \boldsymbol{U}_3 \otimes \boldsymbol{U}_4)^\mathsf{T} \cdot \text{im2col}(\boldsymbol{\mathcal{X}})\right)\right)$, the total complexity is $\Theta\left(m(n^3r + n^2r^2 + nr^3 + r^4 + nr)\right)$. Hence, the speed-up ratio $E$ is lower bounded by $\Omega\left(\frac{n^4}{r^4 + n^3r + n^2r^2 + nr^3 + nr}\right)$.      Q.E.D.

### G. Proof of Corollary 2.

*Proof.* $E \geq \frac{n^4}{r^4 + n^3r + n^2r^2 + nr^3 + nr} \geq \frac{n^4}{r^4(n^3 + n^2 + 2n + 1)}$. If $r \leq \left(\frac{n^4}{\tau(n^3 + n^2 + 2n + 1)}\right)^{\frac{1}{4}}$, then $\frac{n^4}{r^4(n^3 + n^2 + 2n + 1)} \geq \tau$. Therefore, $E \geq \tau$.      Q.E.D.

### REFERENCES

[1] H. Li, S. De, Z. Xu, C. Studer, H. Samet, and T. Goldstein, "Training quantized nets: A deeper understanding," in *Proceedings of the 31st International Conference on Neural Information Processing Systems*, 2017, pp. 5813–5823.

[2] T. G. Kolda and B. W. Bader, "Tensor decompositions and applications," *SIAM review*, vol. 51, no. 3, pp. 455–500, 2009.
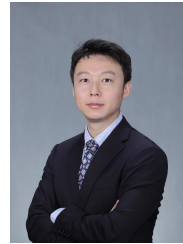
**Wei Dai** received B.S. degree from the Chinese University of Hong Kong, Shenzhen in 2019 and the Master degree from University of Minnesota, Twin Cities in 2020. He is currently pursuing Ph.D. degree in Computer and Information Engineering at the Chinese University of Hong Kong, Shenzhen. His research interests include distributed cloud/edge computing, federated learning, and artificial intelligence.

**Kai Hwang** is a Presidential Chair Professor at the Chinese University of Hong Kong (CUHK), Shenzhen, China. He received the Ph.D. in EECS from the University of California at Berkeley. He has worked at Purdue University and University of Southern California for many years prior joining the CUHK in 2018.

Dr. Hwang has published 10 scientific books and over 280 scientific papers. An IEEE Life Fellow, He has received the Outstanding Achievement Award in 2005 from China Computer Federation and the Lifetime Achievement Award from IEEE CloudCom 2012. In 2020, he received the Tenth Wu Wenjun Artificial Intelligence Natural Science Award from China's Artificial Intelligence Association for his recent work on AI-oriented clouds/datacenters.

**Jicong Fan** received his B.E and M.E degrees in Automation and Control Science & Engineering, from Beijing University of Chemical Technology, Beijing, P.R., China, in 2010 and 2013, respectively. From 2013 to 2015, he was a research assistant at the University of Hong Kong. He received his Ph.D. degree in Electronic Engineering, from City University of Hong Kong, Hong Kong S.A.R. in 2018. From 2018.01 to 2018.06, he was a visiting scholar at the Department of Electrical and Computer Engineering, University of Wisconsin-Madison, USA. From 2018.10 to 2020.07, he was a Postdoc Associate at the School of Operations Research and Information Engineering, Cornell University, Ithaca, USA. Currently, he is a Research Assistant Professor at the School of Data Science, The Chinese University of Hong Kong (Shenzhen) and Shenzhen Research Institute of Big Data, Shenzhen, China. His research interests include statistical process control, signal processing, computer vision, optimization, and machine learning.