# A. Appendix

## A.1. Proof of Unbiasedness and Convergence for JointSQ

Our quantizer selection is similar to the random uniform quantizer as referenced in [1]. We first prove the unbiasedness and convergence for the general form of Co-compressor. For a gradient vector $\mathbf{g}$, the $i$-th gradient element is quantized as follows:

$$\mathcal{Q}_b[g_i] = \|\mathbf{g}\| \cdot \mathrm{sgn}(g_i) \cdot \zeta(g_i, s),$$

where $\|\mathbf{g}\|$ is the $l_2$ norm of $\mathbf{g}$; $\mathrm{sgn}(g_i) = \{+1, -1\}$ is the sign of $g_i$; $s$ is the quantization level. If we use $b$ bits to quantize $g_i$, we will use one bit to represent its sign and the other $b-1$ bits to represent $\zeta(g_i, s)$, thus resulting in a quantization level $s = 2^{b-1}-1$. And $\zeta(g_i, s)$ is an unbiased stochastic function that maps scalar $|g_i|/\|\mathbf{g}\|$ to one of the values in set $\{0, 1/s, 2/s, \ldots, s/s\}$: if $|g_i|/\|\mathbf{g}\| \in [l/s, (l+1)/s]$, we have:

$$\zeta(g_i, s) = \begin{cases} l/s, & \text{with probability } 1 - p_r, \\ (l+1)/s, & \text{with probability } p_r = s\frac{|g_i|}{\|\mathbf{g}\|} - l. \end{cases}$$

So we have:

$$\begin{aligned}\mathbb{E}[\zeta(g_i, s)] =& \frac{l}{s}\left[1 - s\frac{|g_i|}{\|\mathbf{g}\|} + l\right] \\ &+ \frac{l+1}{s}\left[s\frac{|g_i|}{\|\mathbf{g}\|} - l\right] = \frac{|g_i|}{\|\mathbf{g}\|}.\end{aligned}$$

Then:

$$\begin{aligned}\mathbb{E}\left[\zeta(g_i, s)^2\right] &= \mathbb{E}[\zeta(g_i, s)]^2 + \mathbb{V}[\zeta(g_i, s)] \\ &= \frac{|g_i|^2}{\|\mathbf{g}\|^2} + \frac{1}{s^2}p(1-p) \\ &\leq \frac{|g_i|^2}{\|\mathbf{g}\|^2} + \frac{1}{4s^2}.\end{aligned}$$

Considering that $Q_s(g_i) = \|\mathbf{g}\| \cdot \mathrm{sgn}(g_i) \cdot \zeta(g_i, s)$, we have:

$$\begin{aligned}\mathbb{E}\left[\|Q_b[\mathbf{g}]\|^2\right] &= \sum_{i=0}^{d}\mathbb{E}\left[\|\mathbf{g}\|^2\zeta(g_i, s)^2\right] \\ &\leq \sum_{i=0}^{d}\|\mathbf{g}\|^2\left(\frac{|g_i|^2}{\|\mathbf{g}\|^2} + \frac{1}{4s^2}\right) \\ &= \|\mathbf{g}\|^2 + \frac{d}{4s^2}\|\mathbf{g}\|^2.\end{aligned}$$

We can get:

$$\mathbb{E}[Q_b[\mathbf{g}]] = \mathbf{g},$$

$$\mathbb{E}\left[\|Q_b[\mathbf{g}]\|^2\right] \leq \left[1 + \frac{d}{4^b}\right]\|\mathbf{g}\|^2.$$

In the proof presented in [1], the sparsifier is set as the Rand-$k$ sparsifier with an amplification factor of $d/k$. Here, we generalize it to a general unbiased sparsifier. For the stochastic gradient vector $\mathbf{g}$, with a sparsification parameter of $k$, we have the following expression:

$$\mathbb{E}[S_k(\mathbf{g})] = \mathbf{g},$$
$$\mathbb{E}\left[\|S_k(\mathbf{g})\|^2\right] \leq \|\mathbf{g}\|^2.$$

Therefore, for the general form of the Co-compressor that utilizes uniform random quantization and unbiased sparsification, we obtain:

$$E[\hat{\mathbf{g}}] = \mathbb{E}[Q_b[S_k(\mathbf{g})]] = \mathbb{E}[S_k(\mathbf{g})] = \mathbf{g}, \qquad (1)$$

$$\begin{aligned}E\left[\|\hat{\mathbf{g}}\|^2\right] = \mathbb{E}\left[\|Q_b[S_k(\mathbf{g})]\|^2\right] &\leq \left[1 + \frac{k}{4^b}\right]\|S_k(\mathbf{g})\|^2 \\ &= \left[1 + \frac{k}{4^b}\right]\|\mathbf{g}\|^2.\end{aligned} \tag{2}$$

Eq. (17) demonstrates the unbiasedness for Co-compressor and Eq. (18) provides the convergence analysis for Co-compressor.

Our JointSQ framework treats sparsity as 0-bit quantization and introduces the idea of mixed-precision quantization. We split the gradient vector $\mathbf{g}$ into several subgradients $\mathbf{g}_i$ of length $k_i$ and quantize them with different bit-width $b_i$. For example, the gradient vector $\{0.1, 0.2, 0.3, 0.4\}$ can be split into $\{0.1, 0.2\}$ and $\{0.3, 0.4\}$. For ease of analysis, we set the remaining positions of the subgradients to 0 to match the length of the original gradient vector. Thus, we have $\mathbf{g} = \sum_{i=1}^{n}\mathbf{g}_i$, where $n$ is the number of quantization bit levels. For JointSQ, we analyze its unbiasedness:

$$E[\hat{\mathbf{g}}] = E\left[\sum_{i=1}^{n}\hat{\mathbf{g}}_i\right] = \sum_{i=1}^{n}E(\hat{\mathbf{g}}_i).$$

Based on the unbiasedness for the general form of the Co-compressor, as shown in Eq. (17), we know that $E(\hat{\mathbf{g}}_i) = \mathbf{g}_i$. So we have:

$$E[\hat{\mathbf{g}}] = \sum_{i=1}^{n}\mathbf{g}_i = \mathbf{g}. \qquad (3)$$

Based on the definition of the Euclidean norm (L2 norm), we have $\|\mathbf{g}\|^2 = \sum_{i=1}^{n}\|\mathbf{g}_i\|^2$. Therefore:

$$\begin{aligned}E\left[\|\hat{\mathbf{g}}\|^2\right] &= E\left[\sum_{i=1}^{n}\|\mathbf{g}_i\|^2\right] \\ &= \sum_{i=1}^{n}E\left[\|\mathbf{g}_i\|^2\right].\end{aligned}$$

Based on the convergence for the general form of the Co-compressor, as shown in Eq. (18), we obtain:

$$E\left[\|\mathbf{g}_i\|^2\right] \leqslant \left[1 + \frac{k}{4^b}\right] \|\mathbf{g}_i\|^2 .$$

Therefore:

$$E\left[\|\hat{\mathbf{g}}\|^2\right] \leq \sum_{i=1}^{n} \left[1 + \frac{k}{4^b}\right] \|\mathbf{g}_i\|^2 . \tag{4}$$

Eq. (19) demonstrates the unbiasedness for JointSQ and Eq. (20) provides the convergence analysis for the JointSQ.

## A.2. Proof of Improved Convergence for JointSQ

To contrast with the general form of a Co-compressor, we assume that the gradient tensor $\mathbf{g}$ is compressed using a Co-compressor with sparsity parameter $k$ and quantization bit-width $b$. According to Eq. (3) in the main text, the compression noise in this case is obtained as follows:

$$\begin{aligned}
h(k,b) &\triangleq \frac{k}{4^b} \\
&= \frac{k-2}{4^b} \frac{\|\mathbf{g}'\|^2}{\|\mathbf{g}\|^2} + \frac{1}{4^b} \frac{\|g_1\|^2}{\|\mathbf{g}\|^2} + \frac{1}{4^b} \frac{\|g_2\|^2}{\|\mathbf{g}\|^2} .
\end{aligned}$$

In this particular case, we consider a scenario where we only change the quantization bit-width of two gradient elements in the compressed gradient. We quantize one gradient element, denoted as $g_1$, which is originally quantized to $b$ bits, to $b+x$ bits, where $x \in \mathbb{N}^*$, and we quantize another gradient element, denoted as $g_2$, to $b-x$ bits. According to Eq. (6) in the main text, the compression noise in this case can be derived as follows:

$$h'(k,b) \triangleq \frac{k-2}{4^b} \frac{\|\mathbf{g}'\|^2}{\|\mathbf{g}\|^2} + \frac{1}{4^{b+x}} \frac{\|g_1\|^2}{\|\mathbf{g}\|^2} + \frac{1}{4^{b-x}} \frac{\|g_2\|^2}{\|\mathbf{g}\|^2} .$$

The variation of compressed noise is:

$$\Delta h' = \left(\frac{1}{4^{b+x}} - \frac{1}{4^b}\right) \frac{g_1^2}{\|\mathbf{g}\|^2} + \left(\frac{1}{4^{b-x}} - \frac{1}{4^b}\right) \frac{g_2^2}{\|\mathbf{g}\|^2} .$$

By solving the inequality $\Delta h' < 0$, we obtain the following result:

$$|g_1| > 2^x |g_2| . \tag{5}$$

This provides a case where the convergence for JointSQ is superior to the general form of Co-compressor. In fact, this conclusion can be generalized to reducing the bit-width of multiple gradient elements to improve the bit-width of multiple gradient elements:

$$\sum_{i=1}^{n_1} \frac{4^{x_i}-1}{4^{x_i}} |g_i|^2 > \sum_{j=1}^{n_2} \left(4^{y_j} - 1\right) |g_j|^2 .$$

where $x_1 + x_2 + ... x_{n_1} = y_1 + y_2 + ... y_{n_2}$. This finding demonstrates the significant contribution of JointSQ in expanding the solution space and mitigating the occurrence of suboptimal solutions in Co-compressor.

## A.3. Core Algorithm of JointSQ

### A.3.1 Greedy Allocation

---
**Algorithm 1** Greedy Allocation

---
**Input:** Assignable bit-width $c$, gradient vector $\mathbf{g}$.
**Output:** Mixed-Precision quantization mask $x$, gradient sorting results $SP'$.
1: $Remain\_bit \leftarrow c$ // Remain backpack capacity.
2: $x_{i1} \leftarrow 1$ // Default Selection of 0-bit per Group.
3: $B \leftarrow [0, 2, 4, 8], b_j \in B$ // Available quantization bit widths.
4: $\rho_{ij} \leftarrow \frac{4^{w_{ij}}-1}{4^{w_{ij}}} \frac{g_i^2}{\|\mathbf{g}\|^2}, w_{ij} \leftarrow b_j$ // Profit and weight of items.
5: $SP \leftarrow argsort(\frac{p_{ij}-p_{i,j-1}}{w_{ij}-w_{i,j-1}})$ // Sort by Incremental Profit Density.
6: $SP' \leftarrow SP$
7: **while** $Remain\_bit > 0$ **do**
8:    $i,j \leftarrow SP[0], x_{i,j} \leftarrow 1, x_{i,j-1} \leftarrow 0$ // Select the
9:    item with the highest rank.
10:    $SP \leftarrow Update(SP)$ // Remove the $j$-th item in
11:    $i$-th group from the selection pool.
12:    $Remain\_bit \leftarrow Remain\_bit - w_{i,j}$
13: **end while**
14: return $x, SP'$

---

### A.3.2 Reallocation

---
**Algorithm 2** Reallocation

---
**Input:** Learnable Parameter $R$, assignable bit-width $c$, Mixed-Precision quantization mask $x$, gradient vector $\mathbf{g}$, sorting results in Greedy Allocation $SP$.
**Output:** Mixed-Precision quantization Mask $x$, reduction of compression noise $h$.
1: $\bar{k} \leftarrow \frac{Rc}{8}$ // Constraint Value $\bar{k}$ for Length.
2: $k \leftarrow \sum_{i=1}^{d} \sum_{j=2}^{4} x_{ij}$ // Get $k$ from the Last Reallocation.
3: $h \leftarrow 0, f \leftarrow 0$ // Initialize the compression noise reduction amount $h$ and the fine-tuning flag $f$.
4: $\Delta k = k - \bar{k}$ // Retrieve the number of fine-tuning iterations.
5: **for** $i = 1$ to $|\Delta k|$ **do**
6:    $x', f \leftarrow finetuning(x, SP, \Delta k, f)$ // Fine-tuning the mask according to the rules mentioned in the main text.
7:    $\Delta h = h(g, x') - h(g, x)$ // Calculate the difference in compression noise before and after fine-tuning.
8:    **if** $\Delta h < 0$ **then**
9:       $x \leftarrow x'$ // Keep only the fine-tuning attempts
10:       that result in a reduction of compression noise.
11:       $h \leftarrow h + \Delta h$ // Update the reduction of
12:       compression noise.
13:    **end if**
14: **end for**
15: return $x, h$

---

**Algorithm 3** Fine-tuning
***
**Input:** Mixed-Precision quantization mask $x$, sorting results in Greedy Allocation $SP$, difference between the constraint length and the current length $\Delta k$, fine-tuning flag $f$.

**Output:** Mixed-Precision quantization mask $x$.

1: **if** $\Delta k < 0$ **then**
2:  **if** $f = 0$ **then**
3:    $x_{i_1,j_4} \leftarrow 0, x_{i_1,j_3} \leftarrow 1$
4:    $x_{i_2,j_2} \leftarrow 0, x_{i_2,j_3} \leftarrow 1$
5:    $x_{i_3,j_1} \leftarrow 0, x_{i_3,j_2} \leftarrow 1$
6:    $f \leftarrow 1$
7:    // $i_1$ represents the least ranked 8-bit gradient, $i_2$ represents the highest ranked 2-bit gradient, and $i_3$ represents the highest ranked 0-bit gradient.
8:  **else**
9:    $x_{i_1,j_3} \leftarrow 0, x_{i_1,j_2} \leftarrow 1$
10:    $x_{i_2,j_1} \leftarrow 0, x_{i_2,j_2} \leftarrow 1$
11:    $f \leftarrow 0$
12:    // $i_1$ represents the least ranked 4-bit gradient, $i_2$ represents the highest ranked 0-bit gradient.
13:  **end if**
14: **else**
15:  **if** $f = 0$ **then**
16:    $x_{i_1,j_3} \leftarrow 0, x_{i_1,j_4} \leftarrow 1$
17:    $x_{i_2,j_3} \leftarrow 0, x_{i_2,j_2} \leftarrow 1$
18:    $x_{i_3,j_2} \leftarrow 0, x_{i_3,j_1} \leftarrow 1$
19:    $f \leftarrow 1$
20:    // $i_1$ represents the highest ranked 4-bit gradient, $i_2$ represents the least ranked 4-bit gradient, and $i_3$ represents the least ranked 2-bit gradient.
21:  **else**
22:    $x_{i_1,j_2} \leftarrow 0, x_{i_2,j_3} \leftarrow 1$
23:    $x_{i_2,j_2} \leftarrow 0, x_{i_2,j_1} \leftarrow 1$
24:    $f \leftarrow 0$
25:    // $i_1$ represents the highest ranked 2-bit gradient, $i_2$ represents the least ranked 2-bit gradient.
26:  **end if**
27: **end if**
28: return $x, f$
***

## A.3.3 JointSQ in Distributed Learning

**Algorithm 4** JointSQ in Distributed Learning
***
**Input:** The gradients for the current iteration of training $\mathbf{g}$, the compression ratio $C$, the number of Reallocation performed $T$, the number of nodes in distributed training $N$, the initial value $R_0$, the learning rate of $R$ $\delta_h$.

**Output:** The compressed gradients $\hat{\mathbf{g}}$.

1: On each node:
2: **for** each layer's gradient vector $\mathbf{g}_l$ in $\mathbf{g}$ **do**
3:   $c \leftarrow 32 * len(\mathbf{g}_l) * C$
4:   $x, SP \leftarrow GreedyAllocation(c, \mathbf{g}_l)$
5:   **for** $i = 1$ to $T$ **do**
6:     $x_i, h_i \leftarrow Reallocation(R_{i-1}, c, x_{i-1}, \mathbf{g}_l, SP)$
7:     $R_i \leftarrow R_{i-1} + \delta_h(h_i - h_{i-1})$.
8:   **end for**
9:   $\hat{\mathbf{g}}_l \leftarrow Quantize(\mathbf{g}_l, x)$
10: **end for**
11: All-reduce: $\hat{\mathbf{g}} \leftarrow \sum_{i=1}^{N} \hat{\mathbf{g}}$
12: return $\hat{\mathbf{g}}$
***

## References

[1] Guangfeng Yan, Tan Li, Shao-Lun Huang, Tian Lan, and Linqi Song. AC-SGD: Adaptively compressed SGD for communication-efficient distributed learning. *IEEE Journal on Selected Areas in Communications*, 40(9):2678–2693, 2022. 1