

TABLE I: Comparison of different quantization methods.  $\hat{w} \approx \hat{w}_q$ 

| <i>Method</i> | $w$                | $\hat{w}$   | $\hat{w}_q$   | Parameter of quantizer                          | layer-wise |
|---------------|--------------------|---|---|---|------------|
| <i>DoReFa</i> | $w \in \mathbb{R}$ | $\hat{w} = \frac{\tanh(w)}{\max( \tanh(w) )} \in [-1, 1]$   | $\hat{w}_q = \lfloor \frac{(\hat{w} + \beta)}{\alpha} \rfloor \alpha - \beta$<br>$q_{\hat{w}} \in [0, \dots, 2^k - 1]$  | $\alpha = \frac{2}{2^k - 1}; \beta = 1$         |            |
| <i>WRPN</i>   | $w \in \mathbb{R}$ | $\hat{w} = w$<br>$\hat{w} \in \mathbb{R}$   | $\hat{w}_q = \lfloor \text{clip}(\frac{\hat{w}}{\alpha}, -2^{k-1} - 1, 2^{k-1} - 1) \rfloor \alpha$<br>$q_{\hat{w}} \in [-(2^{k-1} - 1), \dots, 0, \dots, 2^{k-1} - 1]$ | $\alpha = \frac{1}{2^{k-1} - 1}$                |            |
| <i>QIL</i>    | $w \in \mathbb{R}$ | $\hat{w} = \begin{cases} 0, & \text{if }  w  < c - d \\ \text{sign}(w), & \text{if }  w  > c + d \\ (\frac{d}{2} w  + \frac{d-c}{2d})^\gamma \cdot \text{sign}(w), & \text{otherwise} \end{cases}$<br>$\hat{w} \in [-1, 1]$ | $\hat{w}_q = \lfloor \frac{\hat{w}}{\alpha} \rfloor \alpha$<br>$q_{\hat{w}} \in [-(2^{k-1} - 1), \dots, 0, \dots, 2^{k-1} - 1]$   | $\alpha = \frac{1}{2^{k-1} - 1}; \gamma; c; d;$ |            |
| <i>LSQ</i>    | $w \in \mathbb{R}$ | $\hat{w} = w$<br>$\hat{w} \in \mathbb{R}$   | $\hat{w}_q = \lfloor \text{clip}(\frac{\hat{w}}{\alpha}, -2^{k-1}, 2^{k-1} - 1) \rfloor \alpha$<br>$q_{\hat{w}} \in [-2^{k-1}, \dots, 0, \dots, 2^{k-1} - 1]$           | $\alpha$  |            |