# Supplementary Materials of
# Double-Win Quant: Aggressively Winning Robustness of Quantized Deep Neural Networks via Random Precision Training and Inference

Yonggan Fu [1]   Qixuan Yu [1]   Meng Li [2]   Vikas Chandra [2]   Yingyan Lin [1]
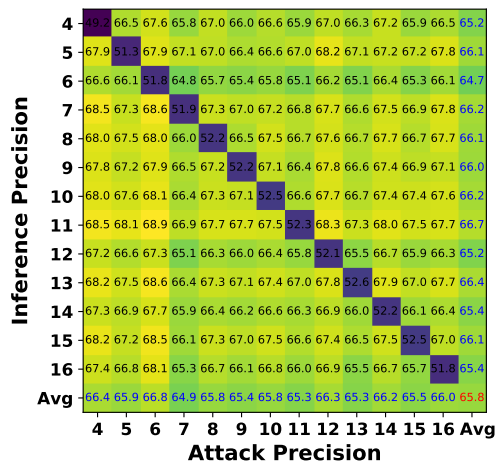
*Figure 1.* Visualizing the transferability of adversarial attacks between different precisions of WideResNet32 trained using the enhanced DWQ, where we annotate the robust accuracy under the PGD-20 attack, and provide both the average accuracy of each row and column (annotated in blue) and the overall average accuracy (annotated in red), the latter of which is also the robust accuracy if RPI is adopted.

## 1. Visualizing the transferability heatmap for WideResNet32 trained by enhanced DWQ

**Experiment settings.** Fig. 1 visualizes the transferability of PGD-20 attacks between different precisions on WideResNet32 trained using the enhanced DWQ with a training/inference precision set of 4~16-bit, where we also annotate both the overall average accuracy (in red) and the average precision of each row and column (in blue) for measuring the robustness of each precision choice.

**Results and analysis.** In addition to the consistently poor transferability of generic quantized DNNs as identified in Fig. 1 of the main content, we also observe that (1) the

enhanced DWQ equipped with RPT and switchable BNs notably enlarges the robust accuracy gaps between different precisions especially under higher precisions, i.e., **the enhanced DWQ further increases the difficulty of transferring adversarial attacks between different precisions** as compared with vanilla DWQ shown in Fig. 1 of the main content, and (2) the average robust accuracies under each attack precision (i.e., each column) are close to each other (within 2%), i.e., **there's no notable bottleneck precision that can be utilized by the attackers**. This set of experiment indicates that although both attackers and defenders can adjust the sampling strategies of their attack/inference precisions for better attack/defense in their competition, DWQ is generally effective and applicable, and our assumption of the random precision strategy in Sec. 4.1 of the main content can fairly evaluate the effectiveness of DWQ.

---

*Proceedings of the 38<sup>th</sup> International Conference on Machine Learning*, PMLR 139, 2021. Copyright 2021 by the author(s).