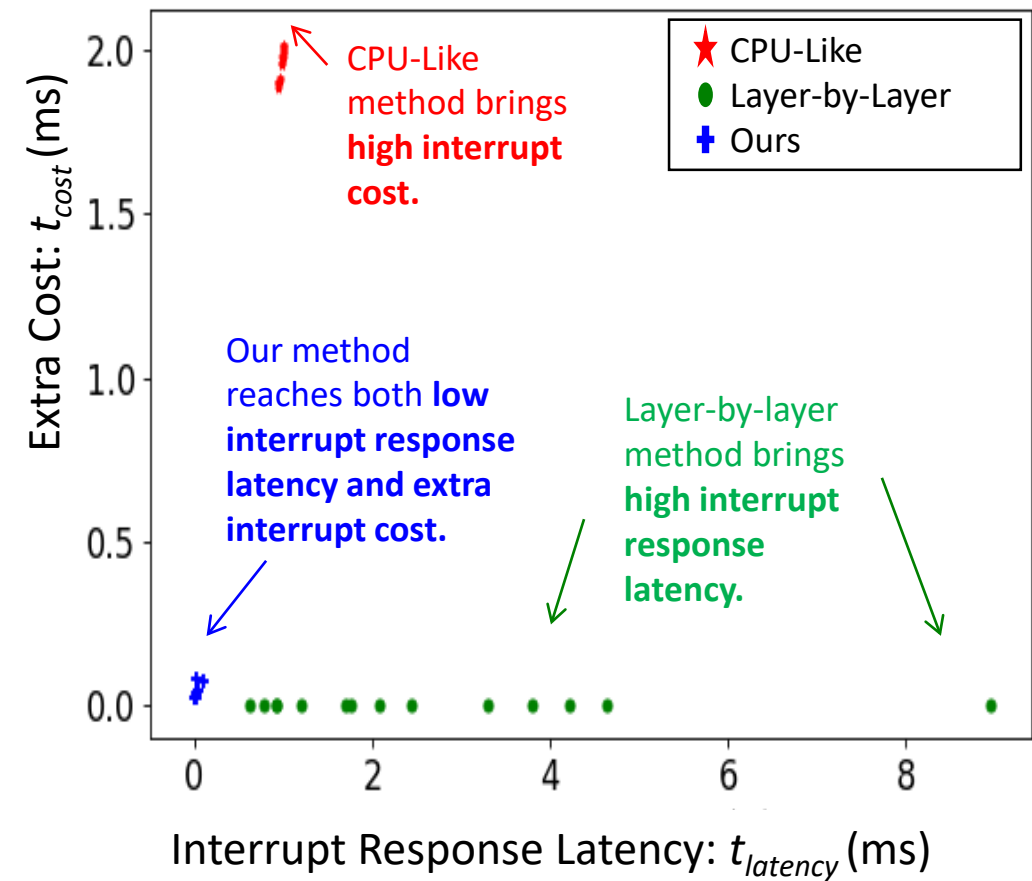(a) Latency comparison between layer-by-layer interrupt method and our virtual-instruction (VI) method. We test 9 layers from different networks (Layer A-I), on a big accelerator with larger parallelism (Left bar of each layer) and small accelerator (Right bar of each layer). The green bar shows latency of from best case to worst case of layer-by-layer method. The red bar shows the best to the worst latency of VI method. The blue lines indicates the average latency.

(b) The interrupt response latency ($t_{latency}$) and the extra time cost ($t_{cost}$) for different interrupt positions in Place Recognition with ResNet101 on the big accelerator with larger hardware parallelism in Fig(a).