

## 715 A APPENDIX

### 716 A.1 Architectures optimized with wiNAS

717  
718 Our framework wiNAS, takes a given macro-architecture  
719 and optimizes each  $3 \times 3$  convolutional layer by choosing  
720 from direct convolution or different Winograd configura-  
721 tions. For the search, all  $1 \times 1$  convolutions were fixed to  
722 use *im2row*.

723  
724 For wiNAS<sub>WA</sub> in FP32, the resulting architecture only sub-  
725 stituted the last convolution layer with *im2row* instead of  
726 *F2*. The rest of the layers remained unchanged from the  
727 WA<sub>F4</sub> configuration (which was described in Section 5.1).  
728 The same micro-architecture was used in CIFAR-10 and  
729 CIFAR-100.

730 For wiNAS<sub>WA</sub> with 8-bit quantization and CIFAR-10,  
731 wiNAS replaced the 5<sup>th</sup> and second last convolutional lay-  
732 ers with *im2row*, instead of *F4* and *F2* respectively. For  
733 CIFAR-100, more optimization was compared to WA<sub>F4</sub>.  
734 The resulting micro-architecture optimization is shown in  
735 Figure 9 (left).

736  
737 When introducing quantization in the search space,  
738 wiNAS<sub>WA-Q</sub>, the resulting architectures are shown in Figure  
739 9 for both CIFAR-10 (middle) and CIFAR-100 (right).

740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755  
756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769

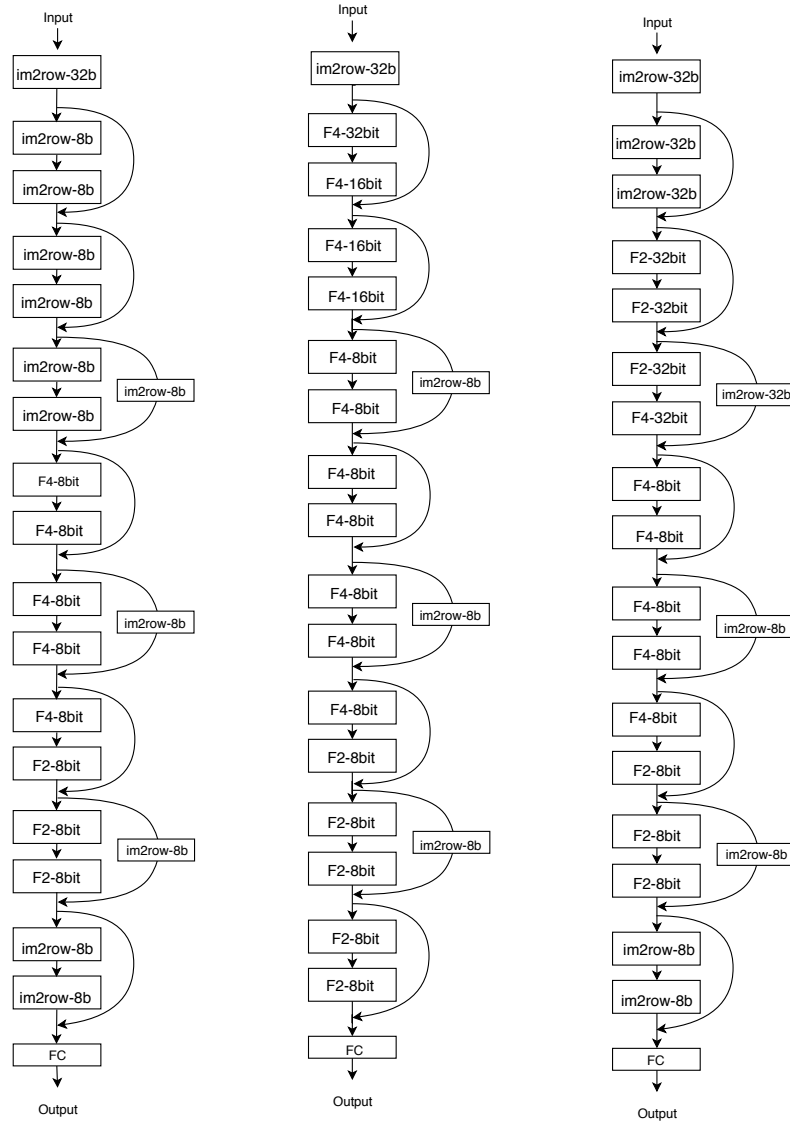


Figure 9. Resulting architectures after optimizing a ResNet-18 macro-architecture using wiNAS. For wiNAS<sub>WA</sub> and CIFAR-100, the architecture resulted is shown on the left. With wiNAS<sub>WA-Q</sub>, that introduces quantization in the search space, the optimization resulted in different architectures for CIFAR-10 (middle) and CIFAR-100 (right), evidencing the difference in complexity of the later.