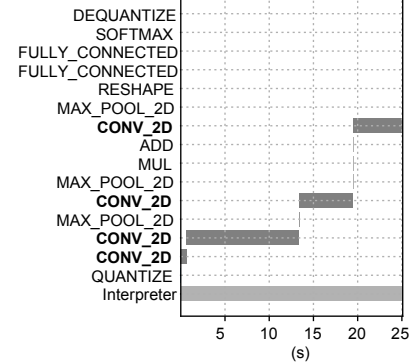
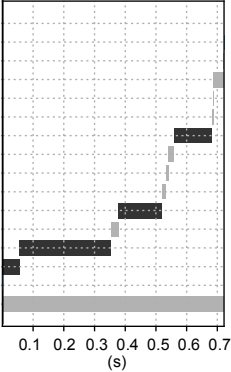


a) Model A (8-bit fixed-point quantization)

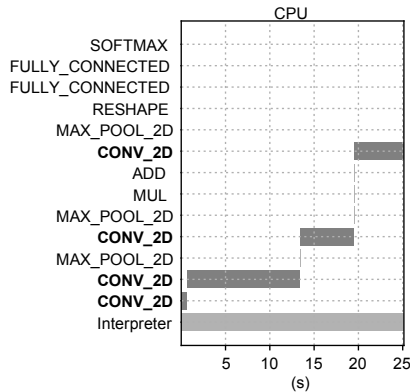
Tensor operation



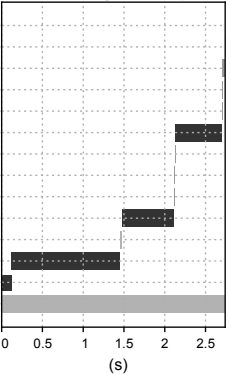
CPU + HW (Fixed-point)



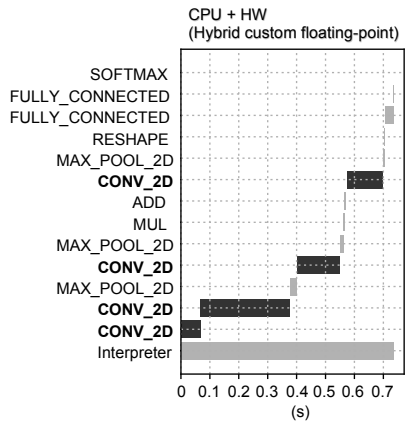
b) Model A (Floating-point)



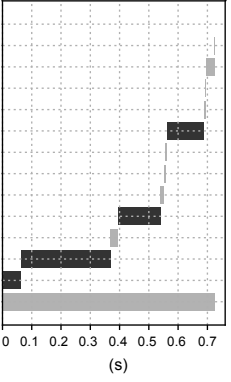
CPU + HW
(Floating-point
Xilinx LogiCORE IP)



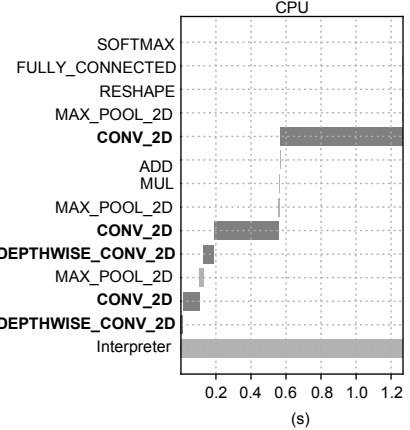
c) Model A (Floating-point)



CPU + HW
(Hybrid logarithmic)



d) Model B (Floating-point)



CPU + HW
(Hybrid custom floating-point)

