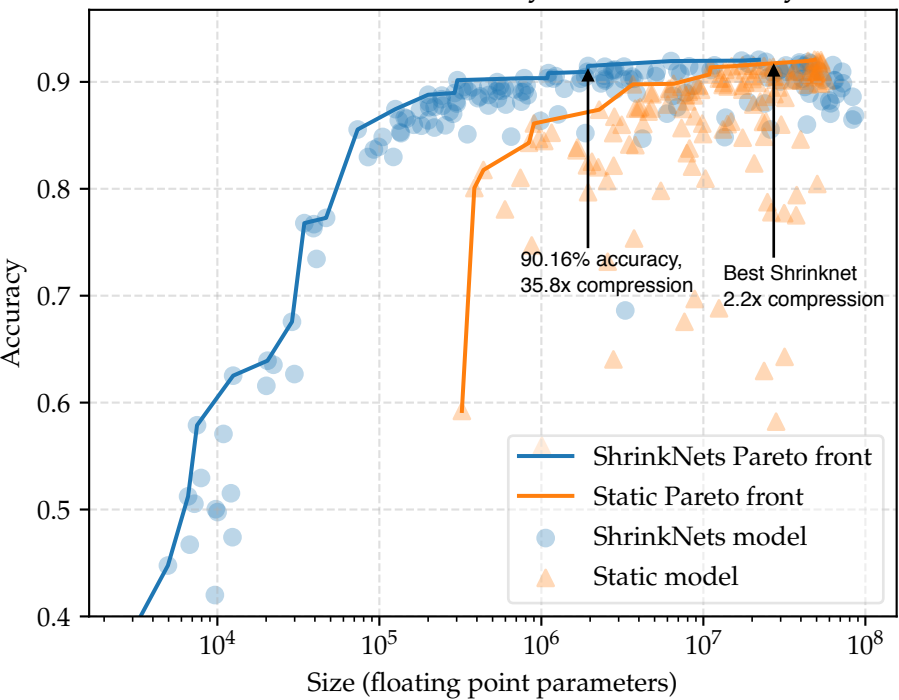
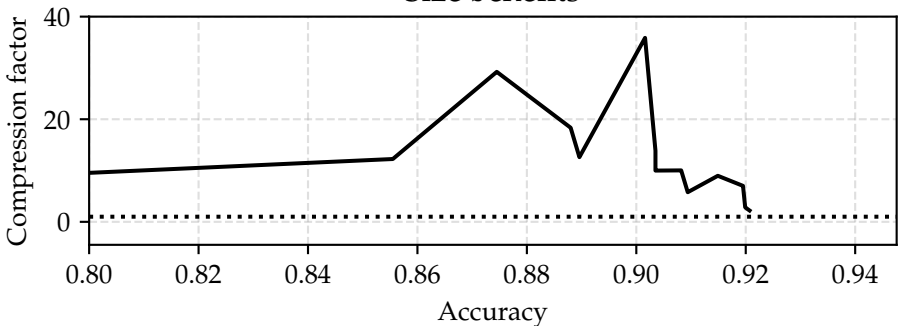


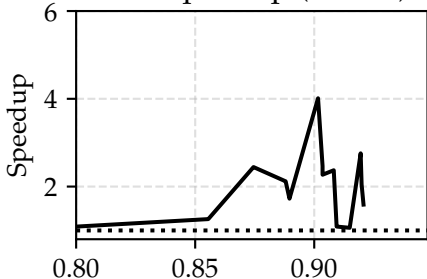
Distribution of models by size and accuracy



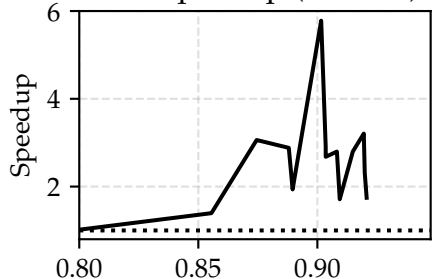
Size benefits



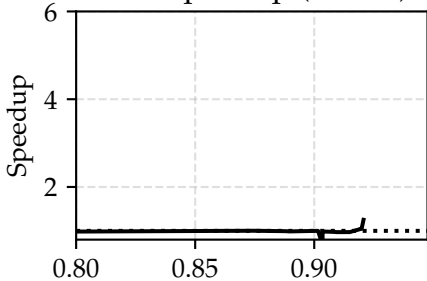
CPU speedup ($bs = 1$)



CPU speedup ($bs = 64$)



GPU speedup ($bs = 1$)



GPU speedup ($bs = 1024$)

