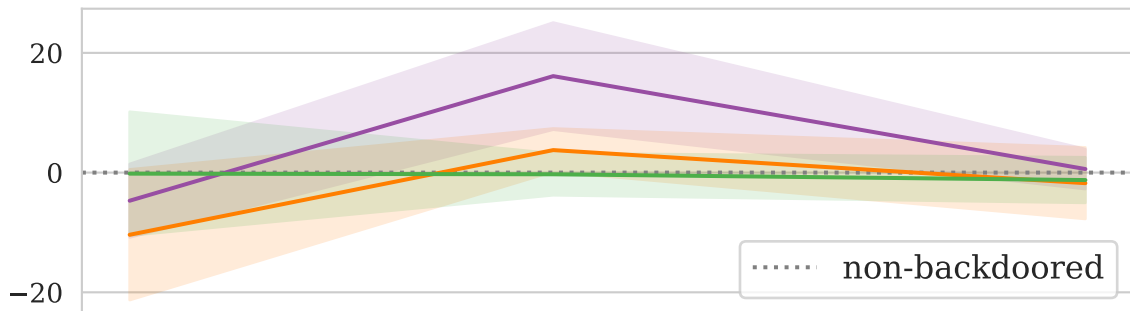
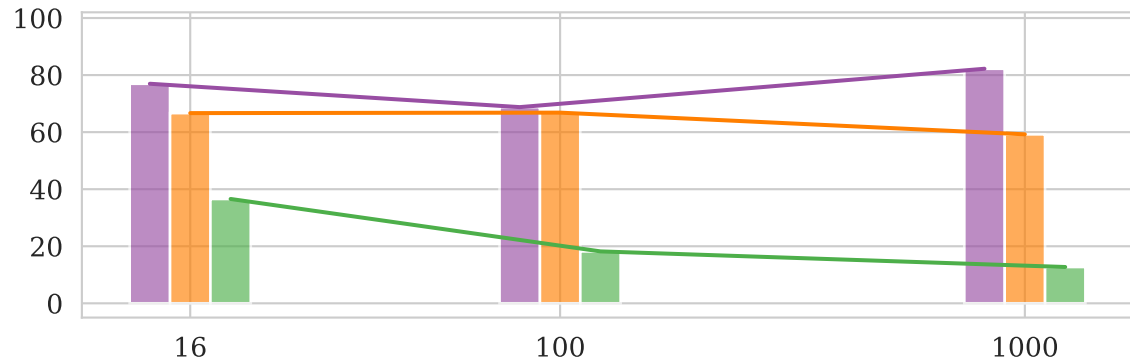


# Visible Backdoor Attacks On TWEETS-HATE-OFFENSIVE

Mean F1  $\Delta$



Mean ASR (L0~L2)



K samples per class

manual\_b

auto\_b

diff\_b