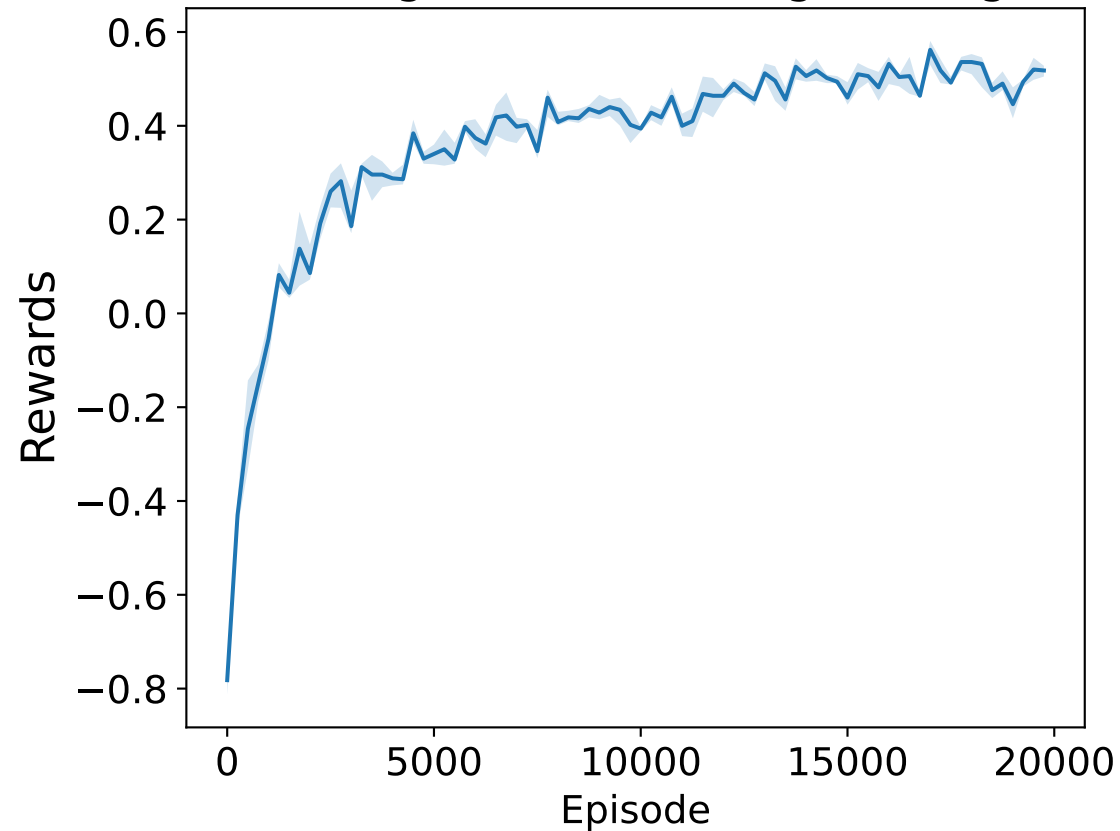
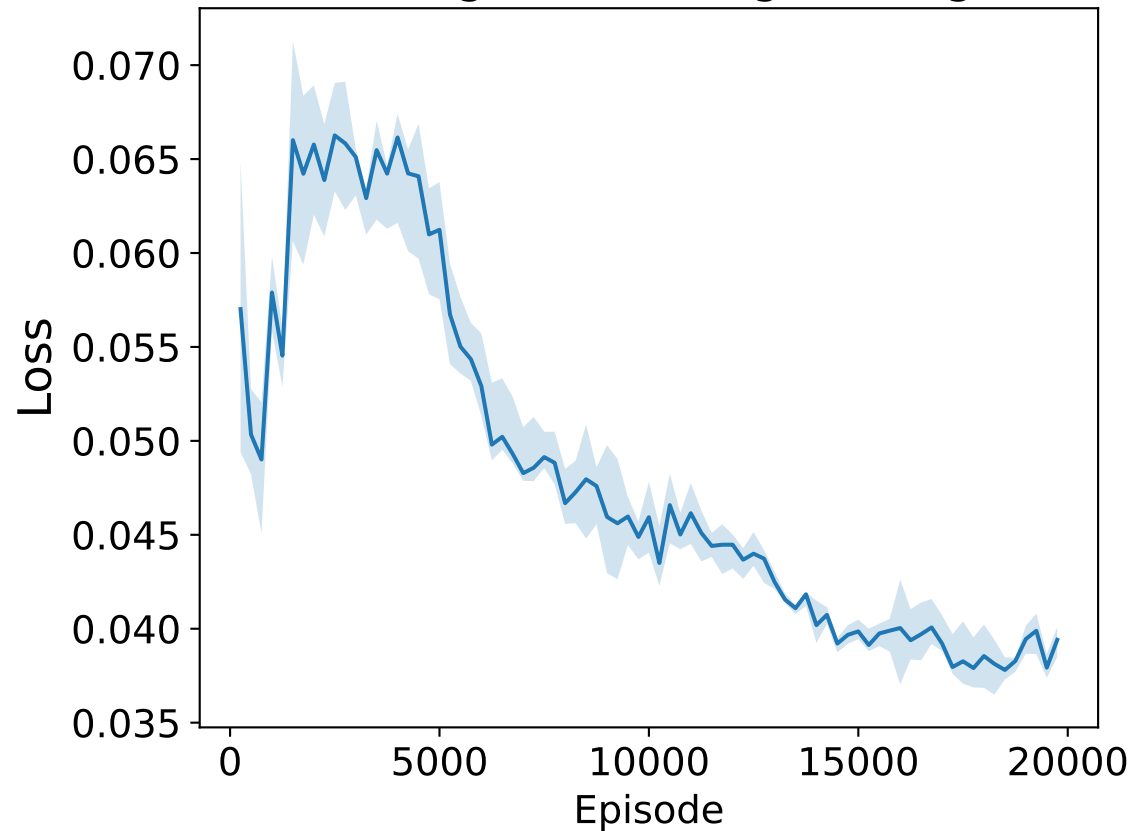


Average rewards during training



Average loss during training



$\varepsilon = 0.1$