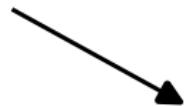
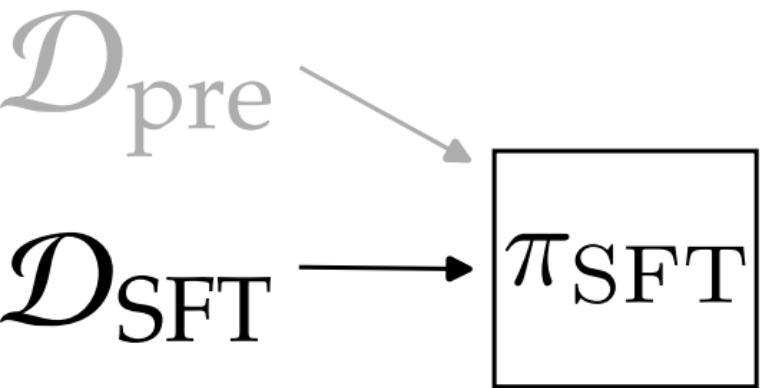
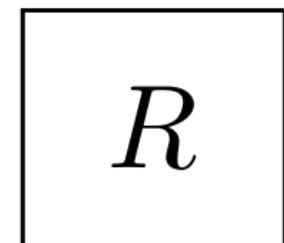
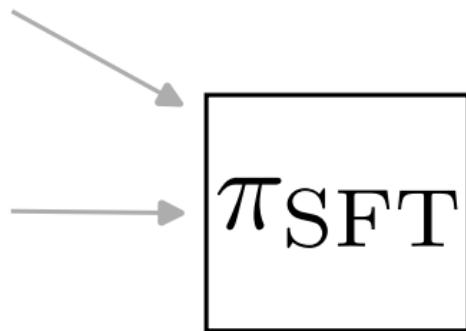


\mathcal{D}_{pre}  π_{pre}



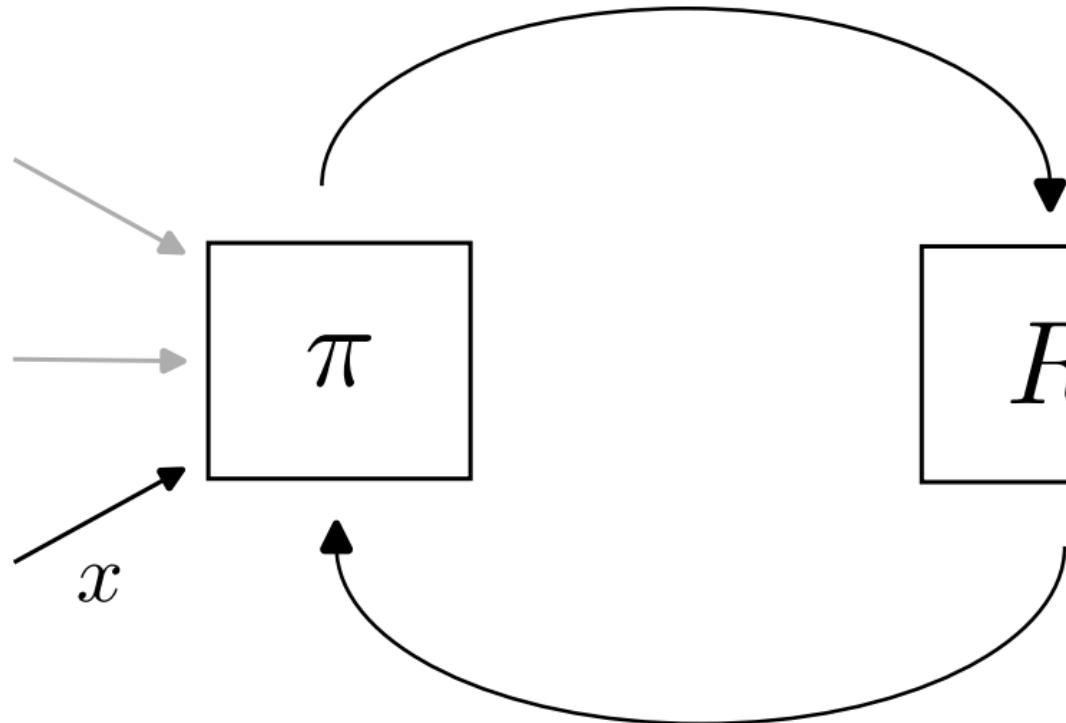
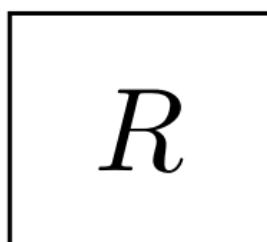
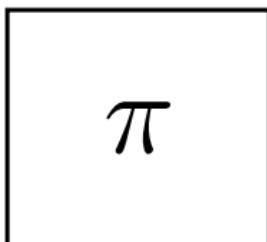
\mathcal{D}_{pre} \mathcal{D}_{SFT} 

generation y

\mathcal{D}_{pre}

\mathcal{D}_{SFT}

\mathcal{D}_{RL}



score $R(x, y)$

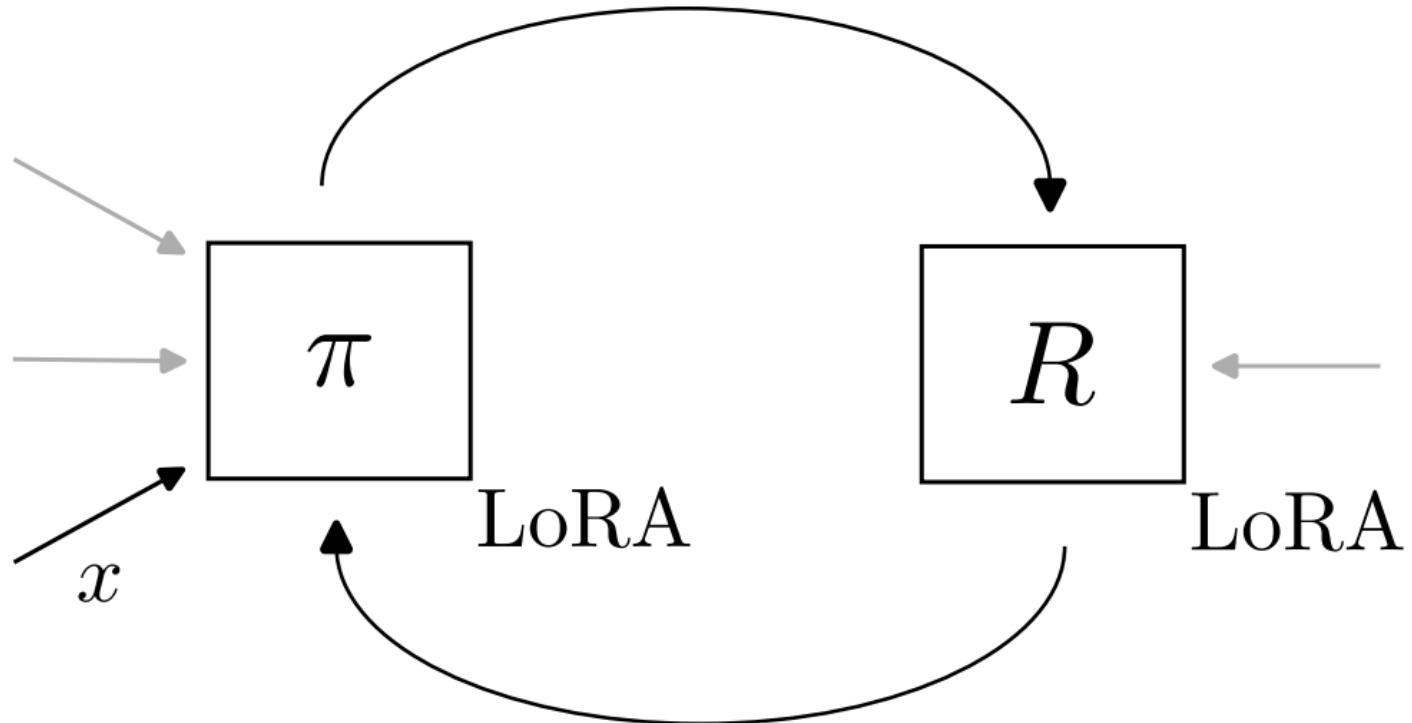
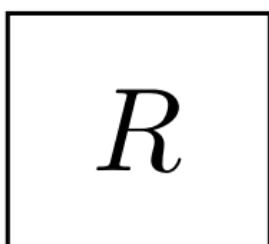
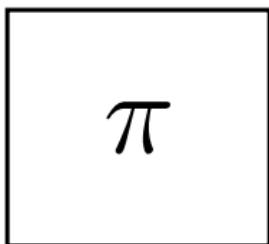
\mathcal{D}_{RM}

generation y

\mathcal{D}_{pre}

\mathcal{D}_{SFT}

\mathcal{D}_{RL}



score $R(x, y)$

