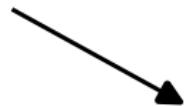
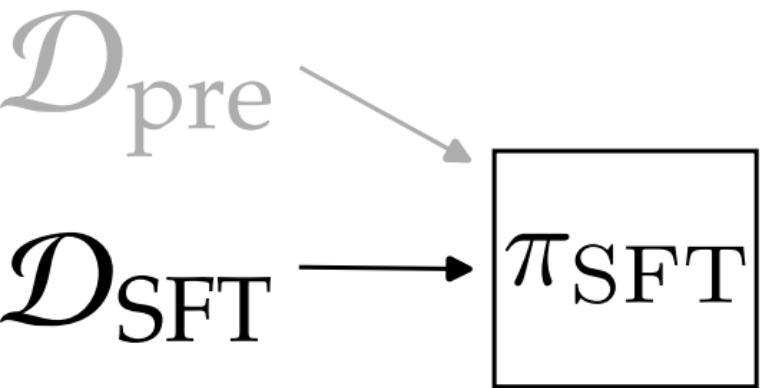
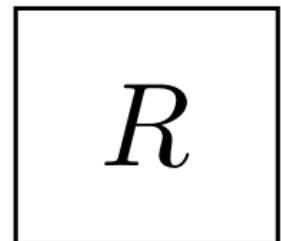
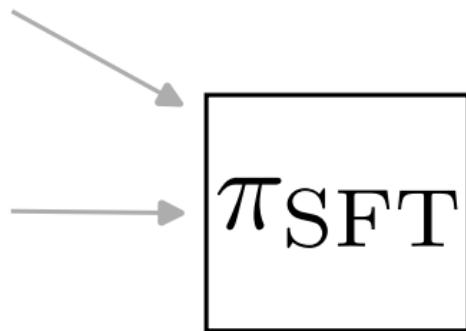


$\mathcal{D}_{\text{pre}}$  $\pi_{\text{pre}}$



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

generation  $y$

$\mathcal{D}_{\text{pre}}$

$\mathcal{D}_{\text{SFT}}$

$\mathcal{D}_{\text{RL}}$

$\pi_\beta$

$R$

$x$

score  $R(x, y)$

generation  $y$

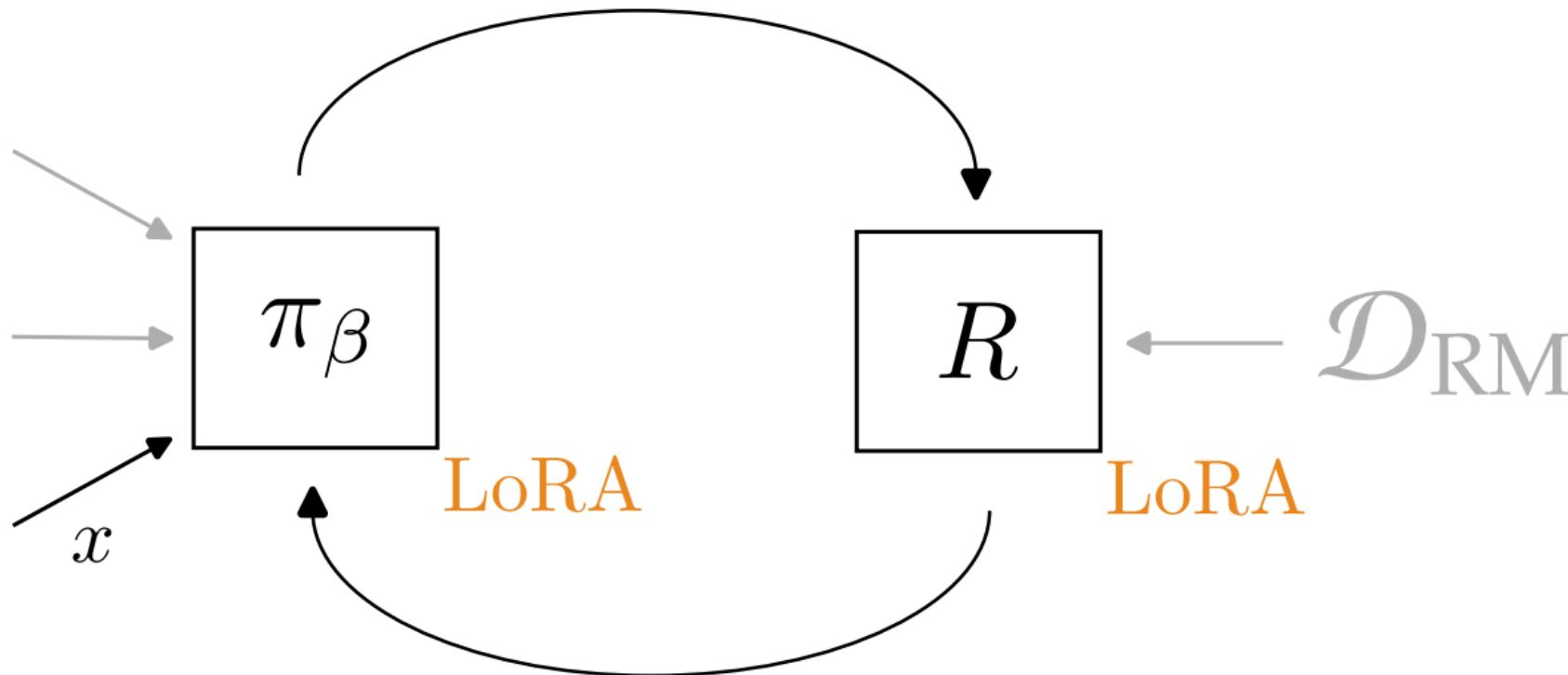
$\mathcal{D}_{\text{pre}}$

$\mathcal{D}_{\text{SFT}}$

$\mathcal{D}_{\text{RL}}$

$\pi_\beta$

$R$



score  $R(x, y)$

