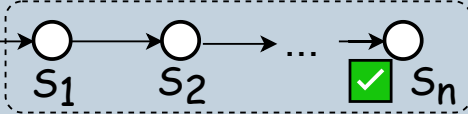


RFT

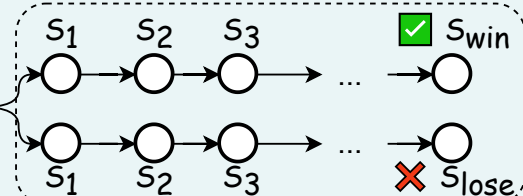
Question



Policy Model

DPO

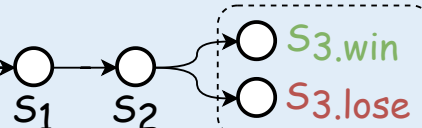
Question



Policy Model

Step-DPO

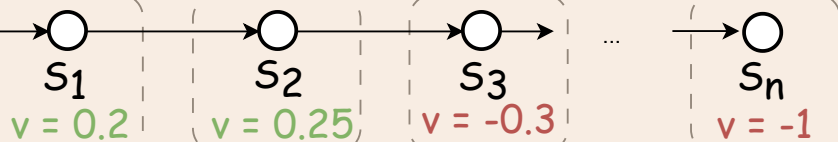
Question



Policy Model

DVO

Question



$$\begin{Bmatrix} v_{s_1}, s_1 \\ v_{s_2}, s_2 \\ \dots \\ v_{s_n}, s_n \end{Bmatrix}$$

Direct
Optimization

Final Policy Model

Granularity of Signal