

structured output y



actions $\alpha = \alpha_1 \dots \alpha_T$



input sentence x