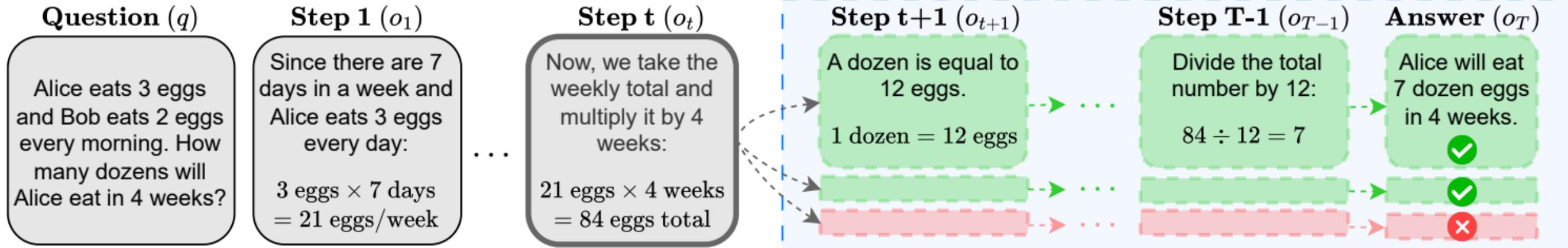


## Reward Maximization Perspective

$$X \succeq Y \iff \sum_{t \geq 0} r(X_t) \geq \sum_{t \geq 0} r(Y_t)$$

$$r_\theta(q_{<t}, o_t) = \alpha \log \frac{\pi_\theta(o_t | q_{<t})}{\pi_{\text{ref}}(o_t | q_{<t})}$$

rollout  $\mu$



## Regret Minimization Perspective

$$X \succeq Y \iff \sum_{t \geq 0} \text{Reg}(X_t) \leq \sum_{t \geq 0} \text{Reg}(Y_t)$$

$$-\text{Reg}_\theta^\mu(q_{<t}, o_t) = \alpha \log \frac{\pi_\theta(o_t | q_{<t})}{\pi_{\text{ref}}(o_t | q_{<t})}$$

Sequential forward KL divergence

$$-\alpha \bar{\mathbb{D}}_{\text{KL}}(\mu || \pi_\theta; q_{<t}, o_t)$$

Human evaluate actions via **prospective, counterfactual** reasoning toward a verifiable state.