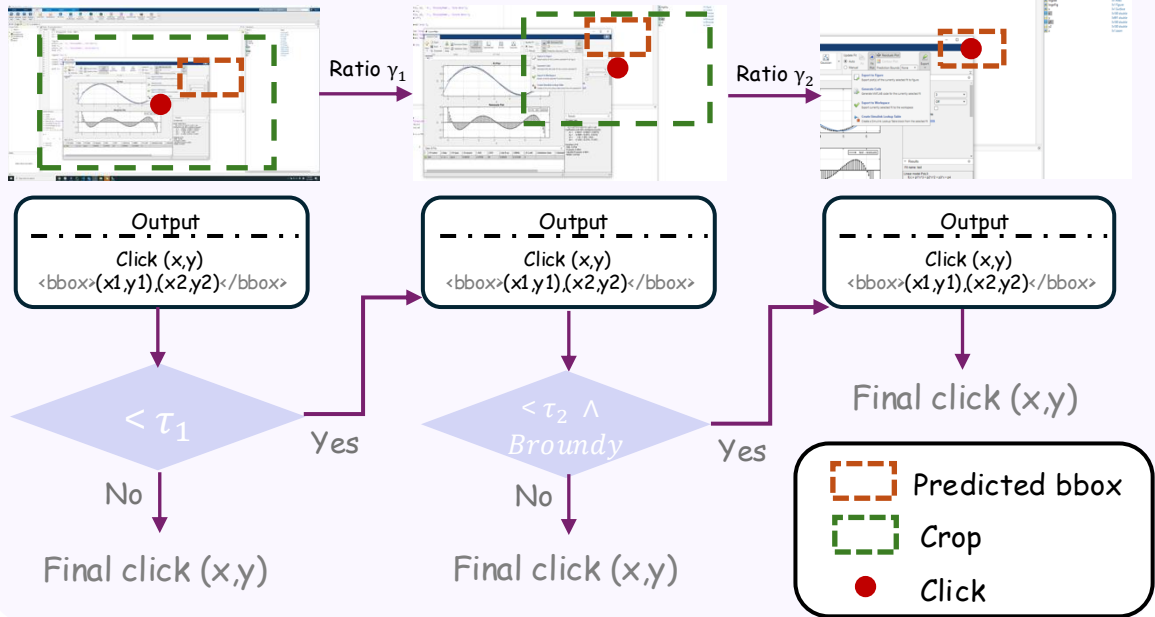


Coarse to fine policy

Stage 1
(full screenshot)

Stage 2
(Large crop)

Stage 3
(Precise click)



Reward design

$$R(p, B) = \lambda_{fnt} R_{fnt} + \lambda_{click} R_{click} + \lambda_{iou} R_{iou}$$

$$R_{click} = 0.6 \cdot \mathbb{I}[p \in b^*] + 0.4 \cdot \exp\left(-\beta \left(\frac{d(p, b^*)}{\sigma(b^*)}\right)^2\right)$$

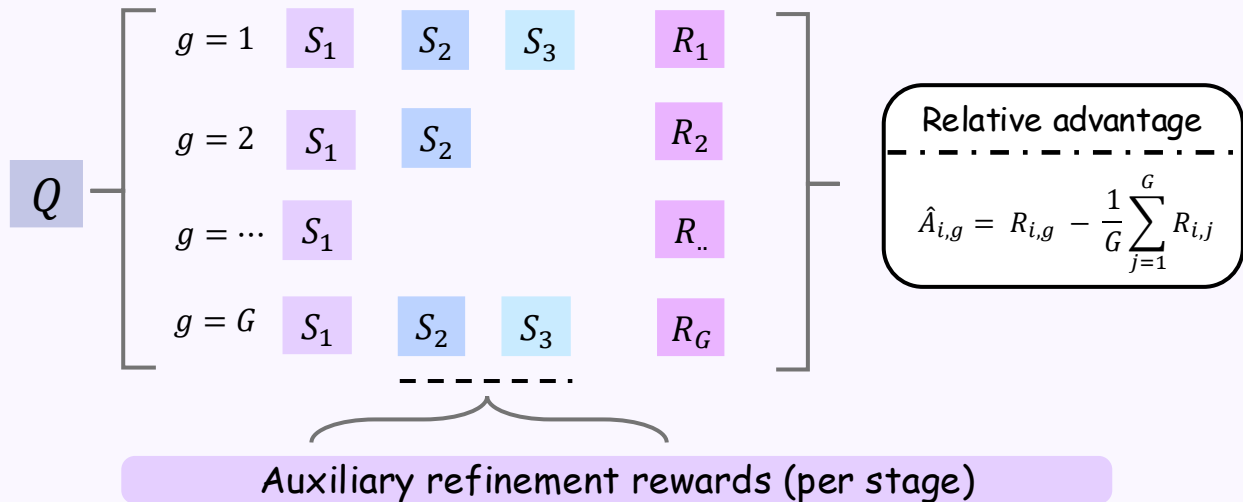
$$R_{iou} = \text{IoU}(B, b^*)$$

Difficulty-aware GRPO



Difficulty score

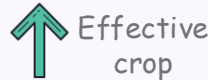
$$L_{GRPO} = -E_{i,g} [\alpha_i \min(r_{i,g} \hat{A}_{i,g}, \text{clip}(r_{i,g}, 1 - \epsilon, 1 + \epsilon) \hat{A}_{i,g})]$$



$$\tilde{R}_t = S_t + \beta^+ \max(S_t - S_{<t}, 0) - \beta^- \max(S_{<t} - S_t, 0) \rightarrow \hat{A}_{i,g}^{(t)}$$



Click quality



Effective crop



Harmful crop

Applied to stage-t response tokens