Rollout Trajectories Cold Start Expert Trajectory multi-trun $o_1, a_1, ... o_M, a_M$ Generate Meta-reasoning Types **Group Relative Advantage Computing Annotated Trajectory** Meta-reasoning Meta-reasoning Meta-reasoning Trajectory Outcome Reward $o_1, T_1^{(a)}, a_1, ... o_M, T_M^{(c)} a_M$ Type (a) Type (b) Type (c) Grouping Grouping Grouping Grouping SFT