Several different RL algos and reward models La RL algos: LO REINFORCE b J(B) = AKC(X,Y) log Po (X,Y) by Ake is r(xiy)-Bon(Pollser) 4 PPO 5J(0) = min (clip (P)) Ath, unchy (P)) xkl) b) $p_i = policy ratio between Pe and Pour$ Li key diff: multiple updates per botch 1) baseline: trained value for LO GRPO Lovery similar to PPO Li key diff: no value for

Ly requires performing G vollouts
per prompt X, not I like in PPD/
REINFORGE

L) advertage is normalized by

Group stats $AGR(x,y;) = \frac{Y(x,y;) - Vx}{\sigma_{x} + \varepsilon}$

by you can use REINFORCE, PPD, GRPD with any bird of reward v(x,17)

Reward models

Display Terry preference model (RUHF)

Distrain r(x,y) on a big datacet

of human pref judgments

Dishumans rate responses 4; gen.

from the SFT model

Disjuen prompt x and two

or more responses 4;

human will produce a ranking

Yw 7 YL

b) train reverd model to make V(X1XM) > V(X1X) Ly verifiable rewords (RLVR) L) cheapy/quickly obtainable 4) no learned reward model, rewards come from environment 4) correctness 4 formatting Li execution-based rewards eig. coding

Examining GRPD advantage: Y(x,y;) = reward $A^{raw}_{i} = Y(x,y;) \Rightarrow \text{no baseline}$ $A^{i}_{i} = Y(x,y;) - V_{x}$ $A^{GR}_{i} = Y(x,y;) - V_{x}$

rare success:

Ly rewards:
$$[1,0,0,0,0,0,0]$$

Ly $V_x = 0.125$, $\delta_x = 0.33$

Ly AGR puts more emphasis on rare success.

higher variance group:

rewords: [6,0,0,0,0,0,0,0]

 $V_{\chi} = 0.75$, $V_{\chi} = 1.98$

$$A_{i}^{raw}$$
 A_{i}^{r} A_{i}^{fh} A_{i}^{fh} Winners +6 +5.25 +2.647 losers 0 -0.75 -0.378