b= butch size L= seq length From last class (PPO): 1. given a batch of prompts $X_1, X_2, ... X_b$ We generate corresponding responses Y1142, ... Yo from policy model 2. for each example (x; , y;), we compute reward r(x;, y;) 3. compute advantages by subtracting baseline) KL penalty A(x,y) = r(x,y) - bb×1 varies
b×1 if using british
mean by L using value

 $A^{kl}(x,y) = A(x,y) - \beta D_{kl}(P_{\theta}||P_{SH})$ $A^{kl}(x,y) = A(x,y) - \beta D_{kl}(P_{\theta}||P_{SH})$ $A^{kl}(x,y) = A(x,y) - \beta D_{kl}(P_{\theta}||P_{SH})$ $A^{kl}(x,y) = A(x,y) - \beta D_{kl}(P_{\theta}||P_{SH})$

$$J(\theta) = \begin{cases} \min \left(P_i A^{kl}(x,y) \right) \\ \text{in } \text{clip} \left(P_i \right), 1-\epsilon, 1+\epsilon \right) A^{l}(x,y) \end{cases}$$

policy reason

Ly we normally take multiple steps on one batch

b) why is can't we just take one step u) a big LR

Ly Stability, if we take one big
Step, we'll move too for and
many tokens will fall into clipped region
by rollouts (i.e. generating y from Po)
are very expensive, compared to
just doing multiple rounds of backprop

PIO -> GRPO

Ly remove reword modely replace will reword to RL gets cheaper

Ly reduce reword hacking

Ly can't reliably compute reword

For many tasks

Ly remove value for, replace ul group mean
Ly KL not included in advantage computation,
but included as separate penalty
by often completely dropped

- 1. given a prompt X, 7 4-32 generate a group of 6 different responses: Y1,72,... Y6
- 2. compute reward for each response in the group r(x, y;)
- 3. comple group mean and std. dev $V_{x} = \frac{1}{6} \sum_{i=1}^{6} V(x_{i}, Y_{i}) \qquad \nabla_{x}$
- 4. comple group relative adventages $A^{GR}(x, y_i) = \frac{r(x, y_i) \mu_x}{\sigma_x}$ $L_1 b \times 1$
- 5, compute loss, do gradient updates

Ly thinking tokens = reasoning

Ly special subsequence of y used for reasoning

X'. let 2x+5=17. Solve for X.)

With put your thinking steps between
With tokens. box your final answer <answer <answer </pre>

Y: < think? well okay so 2x+5= 17, that

means 2x=12, so x=6. </th>

1) revard computed only over the final answer

(5) thinking tokens still part of backpape

(5) separate format reward