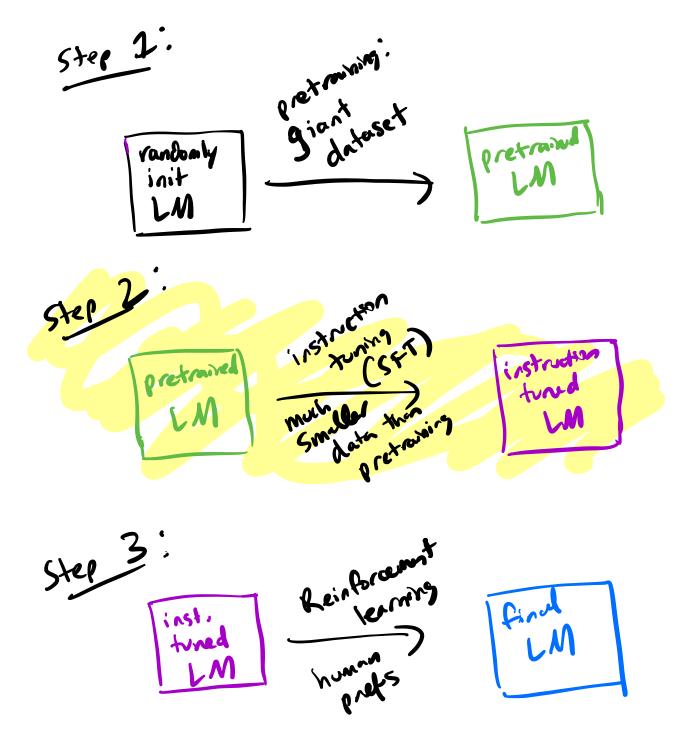
Instruction tuning / PEFT;

Pre-exam: pretraining

- La huge-scale data
- La grant models
- Ly good next word predictor
- Ly not so good: following instructions

Focus: post - training

- L) goal: make pretraised model follow user instructions letter
- Lo subgoal: make model "aligned" w/
 - L7 helpfulness vs. harmless ness



Supervised fine -tuning (SFT); 1) collect a dataset of input/output pairs for a Specific task 15+128 143 1272 +27 1299 by train model on dataset to Minimize NLL on output takens

Minimize NLL on output tokens

Ly prefix LM mask

Ly basically same as pretraining

Ly when doing SFT, you risk

" catachophic forgetting"

L) instruction tuning Ly same as SFT, but data contains huge diversity of tasks Cinstruction > " add these two numbers" Cinput context> => "127+13" Sex.1 Coutput> => 140 "Solve my HW2" " HW2. pdf" Python code for HWZ L7 1013 problem: SFT on huge pretrained models is expensive La adjusting billions of params thru backpape parameter-efficient fine tuning (PEFT) Ly reduce # of params adjusted during FT

Lo prompting: no params are adjusted L) control behavior by modifying prompt 6 demonstrations Ly instruction output 17 fails for complex tasks b) requires extremely large pretrained models b) prompt tuning. L) freeze entire model b add smill # of new param! That are Fld. TO OP TO OP (Pretain) HWZ dot solve _ HW .. ez Soluc my

idea: what if we factorize W into
two low-rank matrices A and B
mer

CLCC C M, n matrix product: ABT is also mxn in LORA: h = f (Wx) be comes lenin h=f(Westrained + ABT) X) frozen jearned these are much much dL Smaller than Tw for each weight matrix 6 Separate A, B Latypically only Wa, Wk, Wu

ra) rank

quantization;

b) intuition: rounding floats to into and

$$X = [-1.0, 0.5, 2.0]$$

4-bit = 16 diff. values

$$S$$
 (ale $S = \frac{x_{max} - x_{min}}{15} = \frac{3}{15} = 0.2$

Zero-point
$$Z = round \left(0 - \frac{x_{min}}{S}\right)$$

$$= 0 - \frac{1.0}{0.2} = 5$$

$$q = round \left(\frac{X}{5} + 2 \right)$$

$$\hat{\chi} = \varsigma \cdot (q - z)$$

$$X = [-1.0, 0.5, 2.0]$$

$$\hat{\chi} = [-1.0, 0.6, 2.0]$$