# Pretraining vs. Post-training:
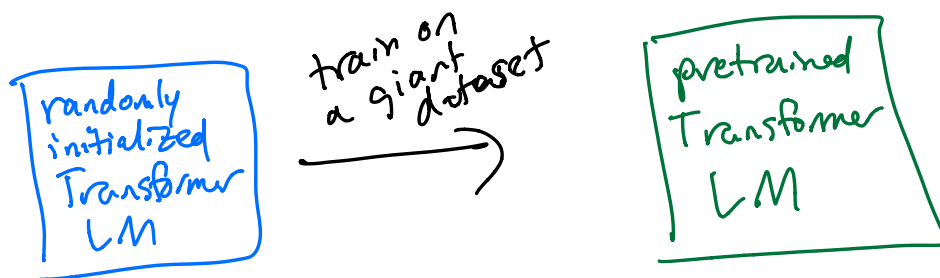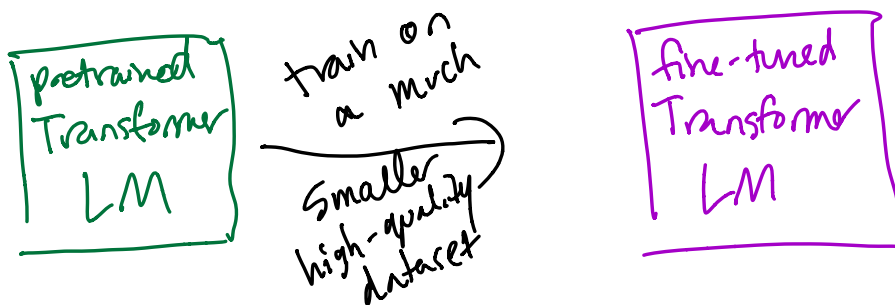
↳ pretraining is conducted w/ as much text as we can obtain

  ↳ trillions of tokens (Common Crawl)

  ↳ biggest model that we can afford

  ↳ goal: to obtain a model that understands many linguistic properties

    ↳ grammar
    ↳ world knowledge
        ↳ who is the president of France?
    ↳ "emergent properties"
        ↳ in-context learning

↳ post-training

  ↳ goal: 1. make a pretrained model follow instructions better
           2. align the model w/ human intents / values

Step 1 (pretraining step):

randomly initialized Transformer LM → train on a giant dataset → pretrained Transformer LM

Step 2 (fine-tuning):

pretrained Transformer LM → train on a much smaller high-quality dataset → fine-tuned Transformer LM

---

==Supervised fine-tuning==; SFT

↳ instruction tuning

↳ goal: make LM follow instructions

1. Collect a dataset of instructions on tasks to solve, and outputs for each instruction.

↳ optionally: collect chains of thought (explanations)

## Sample instruction:

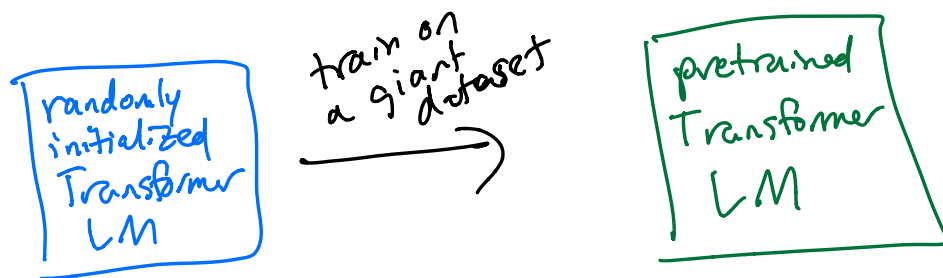Please tell me when the exam will be in this class and how hard it will be!

## output:

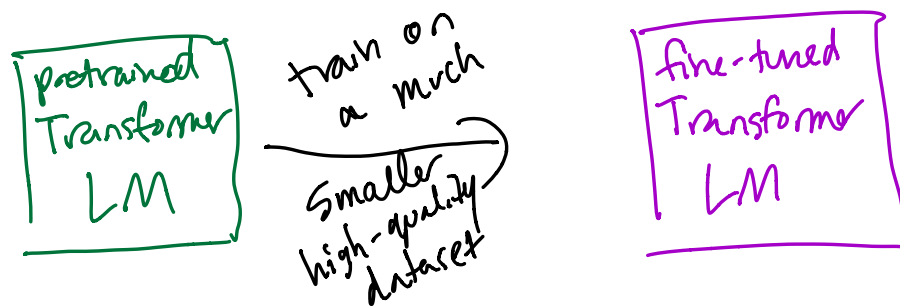The exam will be sometime in April. It may or may not be hard.

instruction (prefix) $\rightarrow$ | pretrained LM | $\rightarrow$ output

# Reinforcement learning from human/AI feedback (RLHF/RLAIF):

## Step 1 (pretraining step):

randomly initialized Transformer LM

→ train on a giant dataset →

pretrained Transformer LM

## Step 2 (instruction tuning):

pretrained Transformer LM

→ train on a much smaller high-quality dataset →

fine-tuned Transformer LM

## Step 3 (RLHF):

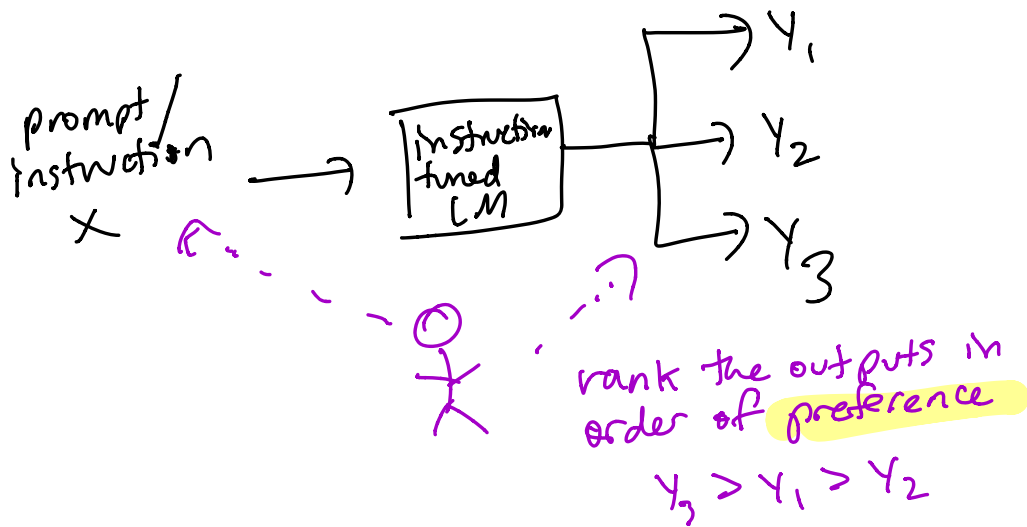instruction-tuned LM

→ reinforcement learning to maximize reward →

final LM

# Limitations of instruction tuning

↳ you only observe <u>one</u> acceptable output per instruction

↳ data diversity issues

↳ don't learn from negative feedback

prompt/ instruction X → [instruction tuned LM] → $Y_1$, $Y_2$, $Y_3$

rank the outputs in order of <mark>preference</mark>

$$Y_3 > Y_1 > Y_2$$

<u>limitation</u>: human prefs are expensive to collect
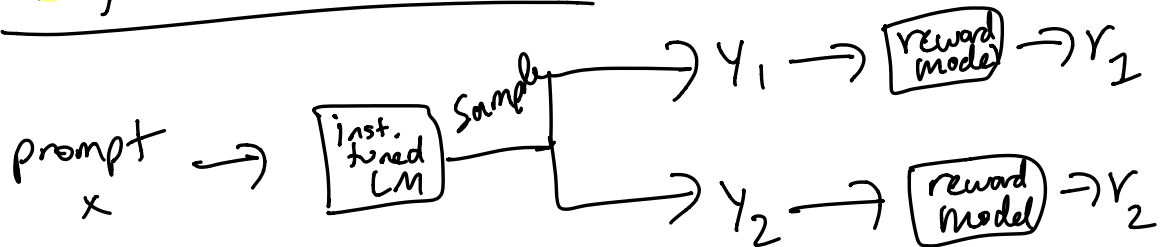
<u>idea</u>: can we train a model to imitate human raters?

<mark>reward model:</mark>

↳ input: prompt X, output $Y_i$

↳ output: scalar score

$\hookrightarrow$ Bradley - Terry pairwise pref model

prompt $\longrightarrow$ [inst. tuned LM] $\xrightarrow{\text{Sample}}$ $\rightarrow Y_1 \rightarrow$ [reward model] $\rightarrow r_1$
x $\rightarrow Y_2 \rightarrow$ [reward model] $\rightarrow r_2$

1. "best of n" sampling
   $\hookrightarrow$ generate n samples, score each one, and then choose sample w/ highest reward
   $\hookrightarrow$ very expensive!

2. just fine-tune the LM on the highest-scoring sample $Y_w$
   $\hookrightarrow$ fine-tune to maximize $p(Y_w | x)$
   $\hookrightarrow$ RAFT

3. reinforcement learning to increase $p(Y_w | x)$ by a small amount, decrease $p(Y_L | x)$ by a small amount. amounts are functions of
   $$r(x, Y_w), \quad r(x, Y_L)$$

$\pi_{ref} \Rightarrow$ instruction-tuned model

$\pi \Rightarrow$ current policy models

$\quad\quad\quad \hookrightarrow$ initialized to $\pi_{ref}$

final
model

$$\max_{\pi} \; \mathbb{E}_{x,y} \left[ \underbrace{r(x,y)}_{\text{reward}} - \underbrace{\beta D_{KL} \left( \pi(y|x) \| \pi_{ref}(y|x) \right)}_{} \right]$$

penalty for deviating
too much from inst. tuned
model