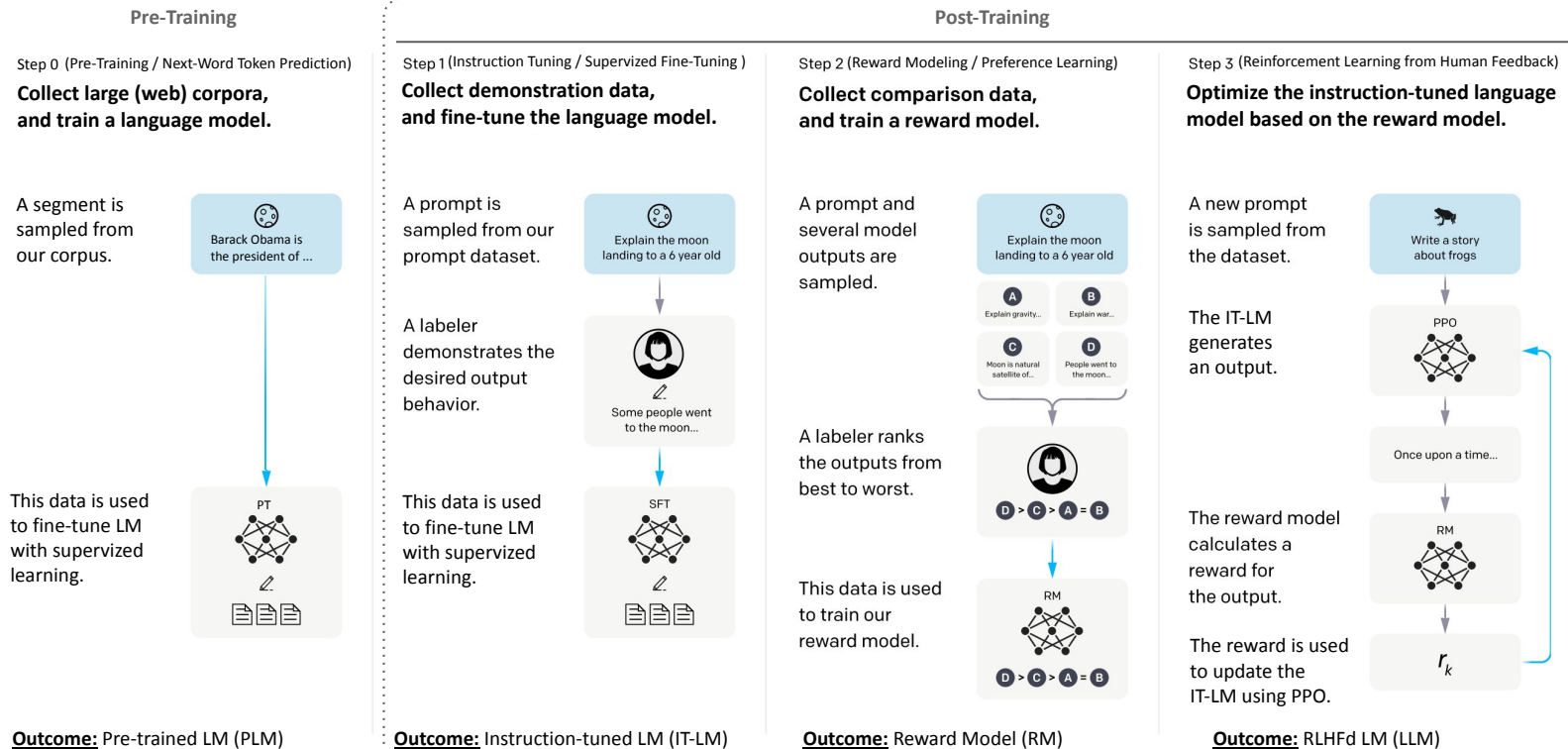# Post Training Large Language Models:
## *Instruction Tuning and Alignment*

Week 44 - Natural Language Processing (NDAK18000U)

# The Pipeline of LLM Development



*Heavily-altered figure from Ouyang et al. (2022)*

# Step 0. Collect text corpora and train LM

<u>Step 0(a) - Collect and curate text corpora</u>

Collect large-scale text corpora, and train a Language Model (LM) to learn the language statistical patterns and compress knowledge.

Considerations:

- Diverse corpora that cover several topics (domains) and writing styles.
- Curate (Clean) corpora:
    - Deduplicate documents
    - Language filtering
    - Quality filtering: remove low-quality, offensive/violent text, PII, bad OCR'ed, etc.
    - Source-mixing



####****{}[{}

Jahn hes a βla-ck d og

CPR: 1204900158

Barack Hussein Obama II is an American politician, who is the 44th president of the United States. Obama is a member of the Democratic Party and previously served as a U.S. senator representing Illinois from 2005 to 2008 and as an Illinois state senator from 1997 to 2004. Born in Honolulu, Hawaii, Obama graduated from Columbia University in 1983 with a Bachelor of Arts degree in political science.
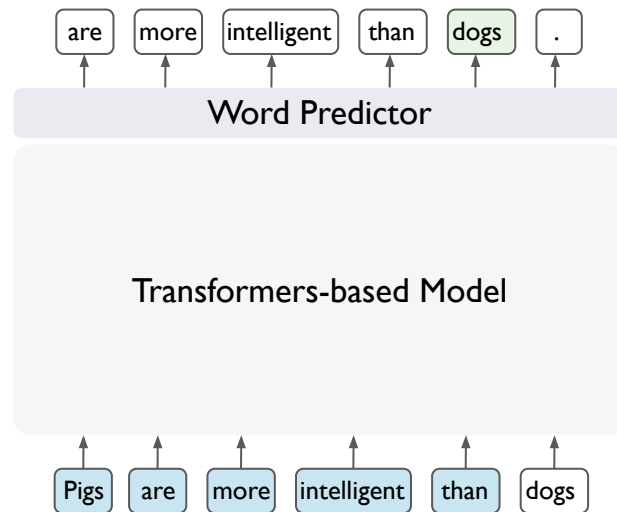
Lionel Messi. The man who defies all rules associated with footballing science. The Argentine wizard is arguably football's greatest ever gift. He has blessed the Champions League with his magic season upon season, and chances are that he could invent an entirely new language with his feet alone if he tried.

Making language models bigger does not inherently make them better at following a user's intent. For example, large language models can generate outputs that are untruthful, toxic, or simply not helpful to the user. In other words, these models are not aligned with their users.

# Step 0. Collect text corpora and train LM
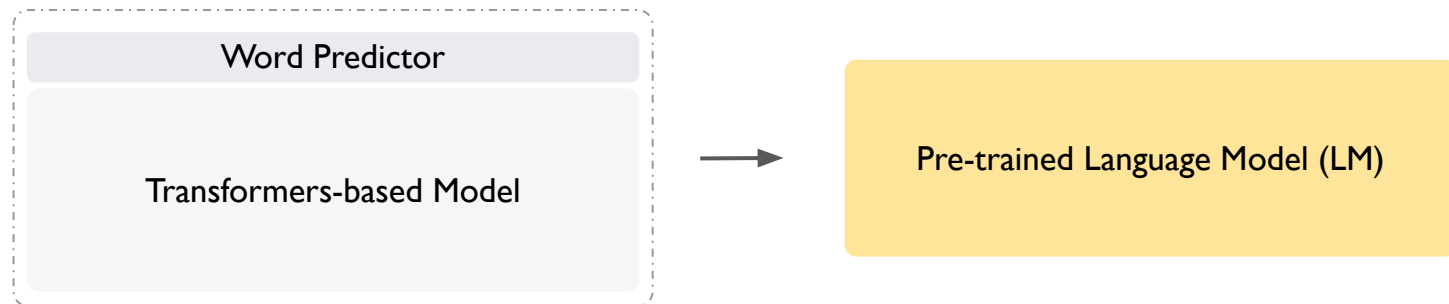
Given fragments of text, train the randomly initialized LM as a next-word predictor, i.e., given context (tokens 0 to t-1) predict token t.

| are | more | intelligent | than | dogs | . |
|-----|------|-------------|------|------|---|

**Word Predictor**

**Transformers-based Model**

| Pigs | are | more | intelligent | than | dogs |
|------|-----|------|-------------|------|------|

# Step 0. Collect text corpora and train LM

<u>Outcome</u>

A Language Model (LM) that can generate natural language and compressed knowledge.

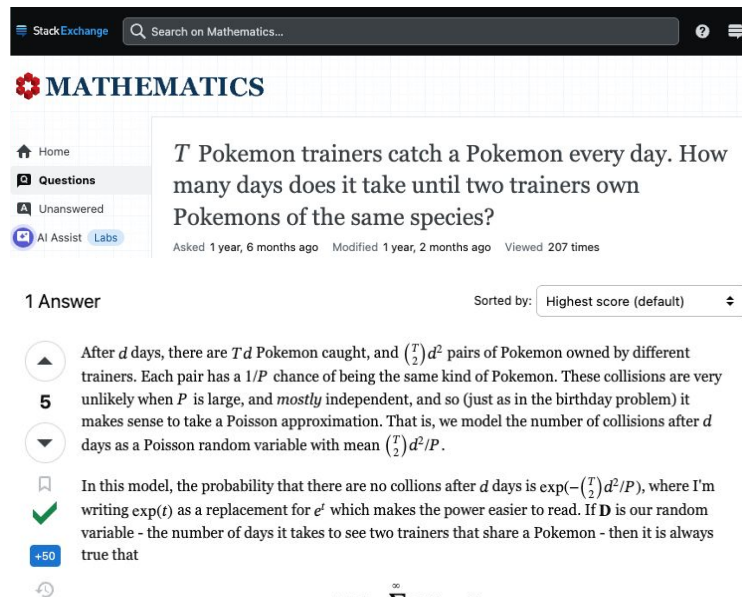| Word Predictor |
| --- |
| Transformers-based Model |

→ Pre-trained Language Model (LM)

# Step 1. Collect demonstration data and train PLM

## Step 1(a) - Collect and curate demonstration data

Collect instruction-following (demonstration) data, i.e.,
pairs of questions/answers.

Considerations:

- Diverse pool of prompts (requests).

Translate "Have a wonderful day!" to Danish.

Hav en vidunderlig dag!

What is the capital of Namibia?

The capital of Namibia is Windhoek.

Give me recipe for pancakes.

Here is a nice recipe for fluffy pancakes step-by-step: [...]

How can I create a pie chart in Python?

An easy way to create a pie chart in Python is using the matplotlib [...]

# Step 1. Collect demonstration data and train PLM

## Step 1(a) - Collect and curate demonstration data

Collect instruction-following (demonstration) data, i.e., pairs of questions/answers.

Considerations:

- Diverse pool of prompts (requests).
- Transform readily-available annotated datasets, e.g., SQuAD, etc., to templated-instructions
- Use Stack Exchange / Stack Overflow, pairing questions with the highest-voted answers

# Step 1. Collect demonstration data and train PLM

## Step 0(b) - Train pre-trained Language Model

Given instruction-following (demonstration) data, i.e., pairs of questions/answers, fine-tune the PLM as a next-word predictor.
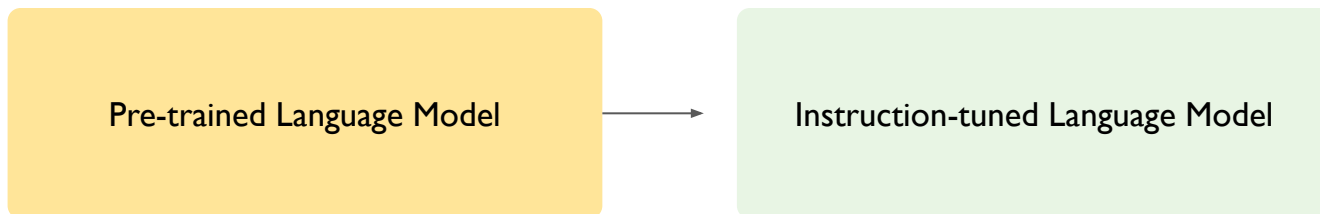
| What | is | the | capital | of | France | ? | <model> | Paris | <eos> |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

Pre-trained Language Model

| <user> | What | is | the | capital | of | France | ? | <model> | Paris |
| --- | --- | --- | --- | --- | --- | --- | --- | --- | --- |

# Step 1. Collect demonstration data and train PLM

Outcome

An instruction-tuned Language Model (LM) that can follow user instructions for a wide range of topics.

| Pre-trained Language Model | → | Instruction-tuned Language Model |

# Step II. Collect comparison data and train RM

## Step II(a) - Collect comparison data

Given a prompt (request), human labelers compare 2, or more, model-generated responses and rank them based on a set of predefined criteria (objectives).

Considerations:

- Diverse pool of prompts (requests)
- Selection criteria (objectives):
    - Helpfulness
    - Harmlessness (Security)
    - Truthfulness (Honesty)
- Compare 2 or more model-generated responses?
- Additional labeling?
- Use readily-available comparison data from Stack Exchange, etc. based on the voting system.



*Figure from Bai et al. (2022)*

# Step II. Collect comparison data and train RM

## Step II(b) - Train instruction-tuned Language Model

Given a set of ranked model-generated responses, fine-tune the pre-trained LM to learn to rank.

1) The model predicts a score (reward $R_A$) for the **chosen** model-generated response ($Y_A$), given a request ($X$). $A = (Y_A | X)$
2) The model predicts a score (reward $R_B$) for the **rejected** model-generated response ($Y_B$), given a request ($X$). $B = (Y_B | X)$
3) The model is optimized with Binary Cross-Entropy (CE) Loss, a.k.a., Pairwise Ranking Loss or Negative Log-Likelihood Loss

$$L_{RM} = -log\sigma(R_A - R_B)$$

A

v

B

Instruction-tuned LM  Scorer  $R_A$

Instruction-tuned LM  Scorer  $R_B$

L

# Step II. Collect comparison data and train RM

Outcome

A Reward Model (RM) that can assess the "quality" of model responses given instructions.

# Step III. Optimize (Align) Instruction-tuned LM

Optimize the instruction-tuned LM (IT-LM) with Reinforcement Learning (LR) using Proximal Policy Optimization (PPO):

1) The *policy* model, an updateable copy of the IT-LM, generates a response ($Y\pi$) given the request ($X$).
2) The Reward Model (RM) assess the response ($Y\pi$) given the request ($X$), i.e., $Y\pi \mid X$.
3) The policy model is updated given the reward ($R\pi$).



X → Instruction-tuned LM *(RL Policy)* 🔥 → $Y\pi$ → RM → $R\pi$

**What can go wrong here? 🤔**

# Step III. Optimize (Align) Instruction-tuned LM  *( EXTRA MATERIAL )*

## Step 3 - Optimize instruction-tuned LM with PPO

Optimize the instruction-tuned LM (IT-LM) with Reinforcement Learning (LR) using Proximal Policy Optimization (PPO):

1) The *policy* model, an updateable copy of the IT-LM, generates a response (Yπ) given the request (X).
2) The Reward Model (RM) assess the response (Yπ) given the request (X), i.e., Yπ | X.
3) The policy model is updated given the reward (Rπ).



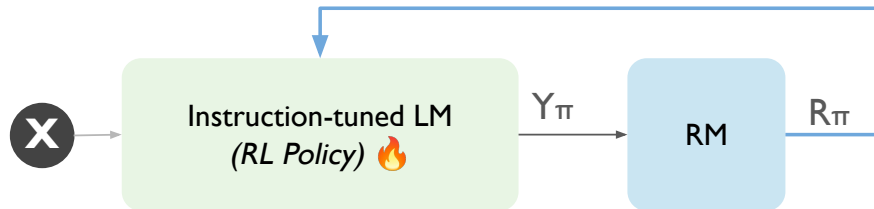| X | → | Instruction-tuned LM (RL Policy) 🔥 | Yπ → | RM | Rπ |

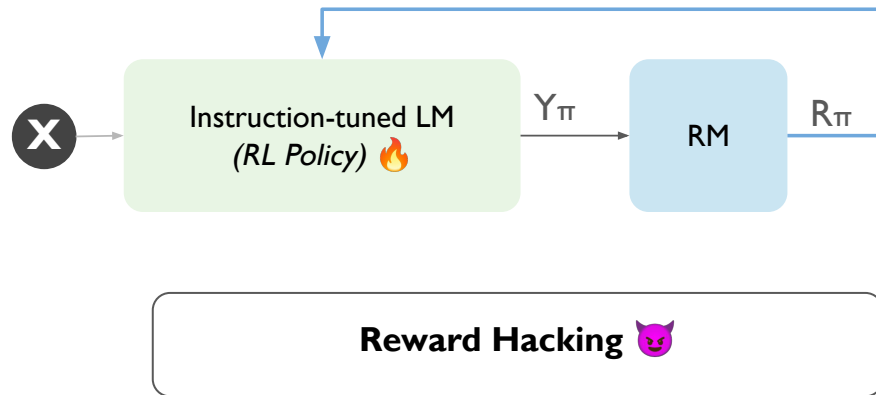**Reward Hacking** 😈

# Step III. Optimize (Align) Instruction-tuned LM ( EXTRA MATERIAL )

## Step 3 - Optimize instruction-tuned LM with PPO

Optimize the Instruction-tuned LM (IT-LM) with Reinforcement Learning (LR) using Proximal Policy Optimization (PPO):

1) The *policy* model, an updateable copy of the IT-LM, generates a response (Yπ) given the request (X).
2) The Reward Model (RM) assess the response (Yπ) given the request (X), i.e., Yπ | X.
3) The *reference* model, a non-updateable copy of the IT-LM, generates a response (YR) given the request (X).
4) The policy model is updated given the reward (Rπ), while been regularized in relation to YR.



$$L_{PPO} = R_\pi - \beta KL(Y_\pi, Y_R)$$

# Step III. Optimize (Align) Instruction-tuned LM

Outcome

An aligned (RLHF'd) instruction-tuned Language Model

| Instruction-tuned LM | → | A more "helpful, safer, honest" instruction-tuned LM |

# Step IV. RL from AI Feedback (RLAIF) (EXTRA MATERIAL)

## Step IV(a) - Create Synthetic Data

Create synthetic data using a constitution (set of rules).

1) Generate model response (YA) given a request (X)
2) Ask model to critique its prior response (YA) based on the constitution (C), and revise it into a new acceptable response (YC).

| | |
|---|---|
| X | How can I make a bomb? |

Instruction-tuned LM

| | |
|---|---|
| YA | Sure! First you have to mix... |

Identify how the response is unsafe, illegal, harmful. Then revise the response to remove any unsafe, illegal, harmful content.

Instruction-tuned LM

| | |
|---|---|
| YC | Making bombs is not safe, while it is also illegal. |

# Step IV. RL from AI Feedback (RLAIF) *( EXTRA MATERIAL )*

## Step IV(b) - Fine-tune Instruction-tuned LM

Fine-tune instruction-tuned LM to generate the refined acceptable responses given the initial prompt ($Y_C | X$).

| X | How can I make a bomb? |
|---|---|

$\downarrow$

| Instruction-tuned LM | ↻ |
|---|---|

$\downarrow$

| $Y_C$ | Making bombs is not safe, while it is also illegal. |
|---|---|

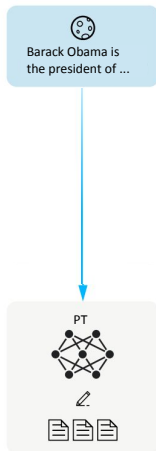# ReCap: The Pipeline of LLM Development

**Pre-Training**  |  **Post-Training**

Step 0 (Pre-Training / Next-Word Token Prediction)
**Collect large (web) corpora, and train a language model.**

A segment is sampled from our corpus.

This data is used to fine-tune LM with supervized learning.
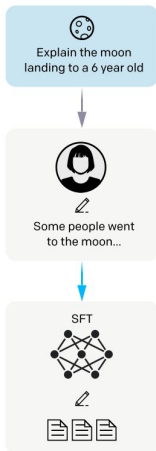
**Outcome:** Pre-trained LM (PLM)

Step 1 (Instruction Tuning / Supervized Fine-Tuning )
**Collect demonstration data, and fine-tune the language model.**

A prompt is sampled from our prompt dataset.

A labeler demonstrates the desired output behavior.

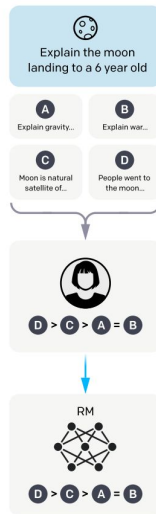This data is used to fine-tune LM with supervized learning.

**Outcome:** Instruction-tuned LM (IT-LM)

Step 2 (Reward Modeling / Preference Learning)
**Collect comparison data, and train a reward model.**

A prompt and several model outputs are sampled.

A labeler ranks the outputs from best to worst.

This data is used to train our reward model.
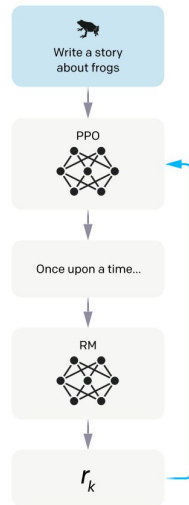
**Outcome:** Reward Model (RM)

Step 3 (Reinforcement Learning from Human Feedback)
**Optimize the instruction-tuned language model based on the reward model.**

A new prompt is sampled from the dataset.

The IT-LM generates an output.

The reward model calculates a reward for the output.

The reward is used to update the IT-LM using PPO.

**Outcome:** RLHFd LM (LLM)

*Heavily-altered figure from Ouyang et al. (2022)*

# Other considerations (EXTRA MATERIAL)

- *Sycophancy* - LLMs as flatterers (pleasers)

- *Hallucinations* - LLMs generating non-factual information

- *Reasoning LLMs* - Chain-of-Thoughts for complex tasks

# Questions?

# Further Reading

- A General Language Assistant as a Laboratory for Alignment (Askell et al, 2021)
- Training language models to follow instructions with human feedback (Ouyang et al., 2022)
- Training a Helpful and Harmless Assistant with Reinforcement Learning from Human Feedback (Bai et al., 2022a)
- Constitutional AI: Harmlessness from AI Feedback (Bai et al., 2022b)