# Large Language Models (ChatGPT and friends)

**Robin Jia**
USC CSCI 467, Spring 2025
April 29, 2025

# ChatGPT

**R**   What machine learning model should I use to determine if an image contains a leopard or a tiger?

For image classification tasks such as this one, a convolutional neural network (CNN) is typically used. CNNs are particularly good at detecting patterns in visual data and are widely used for image classification tasks.

To train a CNN for your specific task, you would first need to gather a dataset of images containing leopards, tigers, and possibly other animals. This dataset would be split into training, validation, and test sets, with the training set used to train the model, the validation set used to tune the model's hyperparameters, and the test set used to evaluate the model's performance on unseen data.

There are many pre-trained CNN models that you can use as a starting point, such as ResNet, VGG, or Inception. You can also choose to train your own model from scratch using a framework such as TensorFlow or PyTorch.
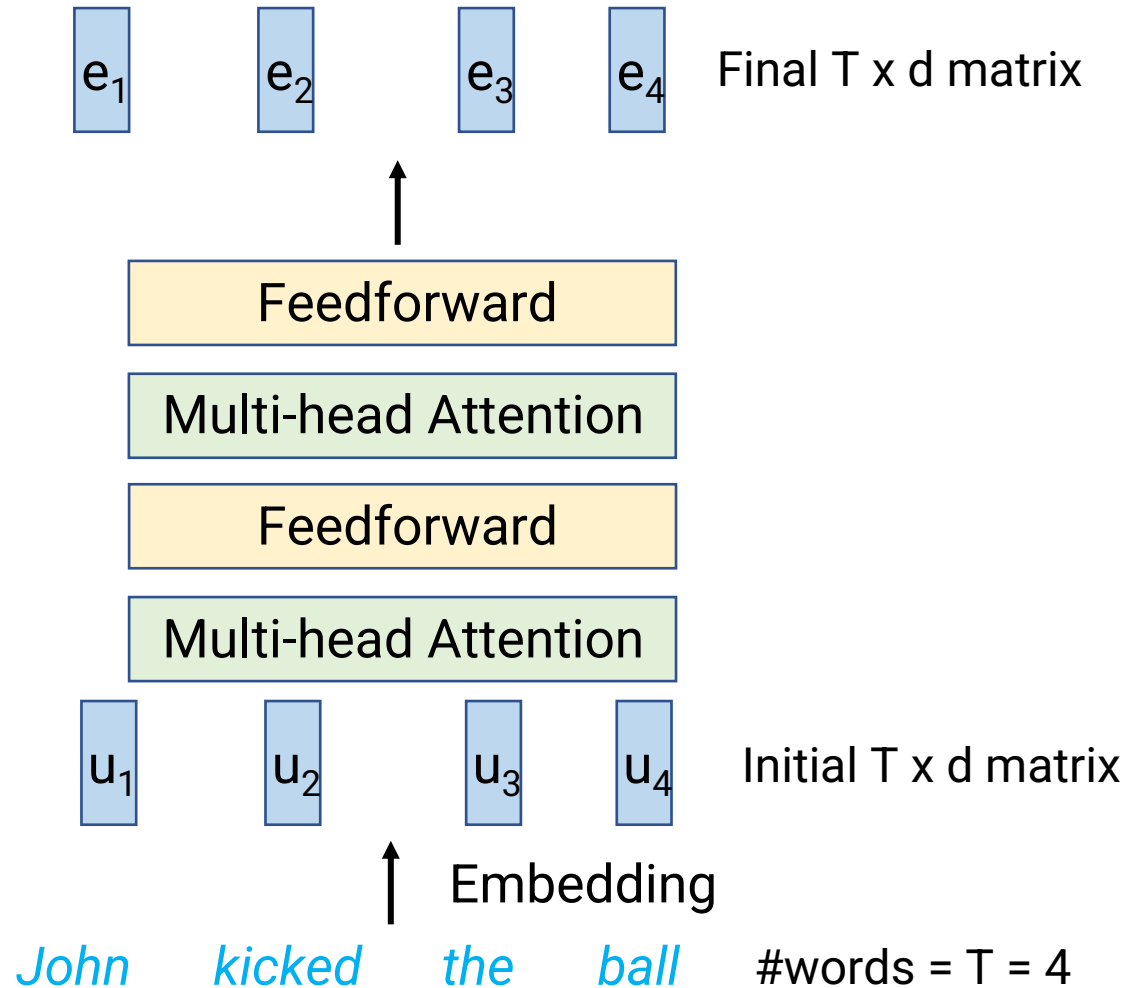
Once you have a trained model, you can pass in a new image and the model will output a probability distribution over the possible classes, in this case, leopard or tiger. The class with the highest probability would be the predicted class for the image.

- Language model-based Chatbot developed by OpenAI
- Large Transformer model
- Lots of information not public, but we do know some things...

# Outline

- Technical ingredients of ChatGPT
    - Language model pretraining
    - Reinforcement Learning from Human Feedback
- Persistent limitations of models
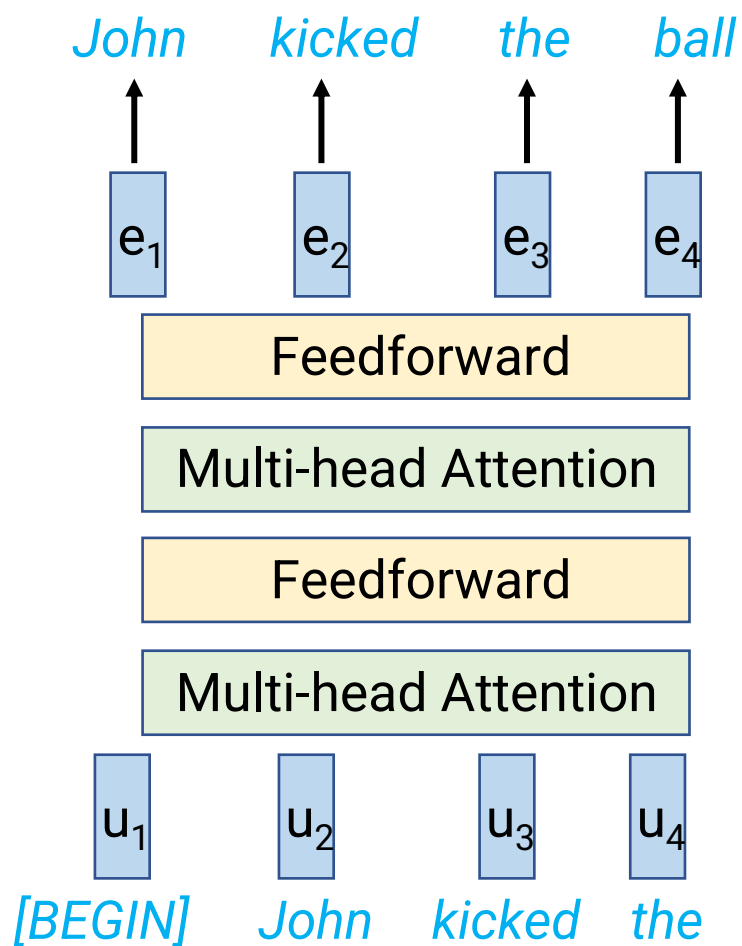
# Previously: Transformers

e₁  e₂  e₃  e₄    Final T x d matrix

Feedforward

Multi-head Attention

Feedforward

Multi-head Attention

u₁  u₂  u₃  u₄    Initial T x d matrix

↑ Embedding

*John*   *kicked*   *the*   *ball*    #words = T = 4

- One transformer consists of
  - Embeddings for each word of size d
    - Let T =#words, so initially T x d matrix
  - Alternating layers of
    - "Multi-headed" attention layer
    - Feedforward layer
    - Both take in T x d matrix and output a new T x d matrix
  - Plus some bells and whistles
    - Residual connections & LayerNorm
    - Byte pair encoding tokenization

# Autoregressive Language Model Training

- Training example: "Convolutional neural networks are good for image classification"

- Want to maximize P("Convolutional neural networks are good for image classification")

- MLE: Take log and decompose by chain rule:

  *log P("Convolutional")*
  *+ log P("neural" | "Convolutional")*
  *+ log P("networks" | "Convolutional neural")*
  *+ log P("are" | "Convolutional neural networks") + …*

- Decomposes into a bunch of **next-word-classification** problems

- Backpropagation + gradient descent to minimize loss

# Review: Transformer autoregressive decoders



- How to do autoregressive language modeling?
- Test-time
  - At time t, attend to positions 1 through t
  - Happens in series

# Review: Transformer autoregressive decoders

- When training a decoder, it has to be "used to" only attending to past/current tokens

- Training time: Masked attention implementation trick
  - Recall: Attention computes $Q$ x $K^T$ (T x T matrix), then does softmax
  - But if generating autoregressively, time t can only attend to times 1 through t
  - Solution: Overwrite $Q$ x $K^T$ to be $-\infty$ when query index < key index
  - **All timesteps happen in parallel**

Queries

| | [BEGIN] | John | kicked | the |
|---|---|---|---|---|
| [BEGIN] | 10 | -2 | 6 | 3 |
| John | 0 | 7 | 2 | -4 |
| kicked | -3 | 4 | 5 | -8 |
| the | 2 | 1 | 7 | 6 |

Keys

# Review: Transformer autoregressive decoders

- When training a decoder, it has to be "used to" only attending to past/current tokens

- Training time: Masked attention implementation trick
  - Recall: Attention computes $Q$ x $K^T$ (T x T matrix), then does softmax
  - But if generating autoregressively, time t can only attend to times 1 through t
  - Solution: Overwrite $Q$ x $K^T$ to be $-\infty$ when query index < key index
  - **All timesteps happen in parallel**

Queries

|  | [BEGIN] | John | kicked | the |
|---|---|---|---|---|
| [BEGIN] | 10 | -2 | 6 | 3 |
| John | 0 | 7 | 2 | -4 |
| kicked | -3 | 4 | 5 | -8 |
| the | 2 | 1 | 7 | 6 |

Keys

# Review: Transformer autoregressive decoders

- When training a decoder, it has to be "used to" only attending to past/current tokens

- Training time: Masked attention implementation trick
  - Recall: Attention computes $Q$ x $K^T$ (T x T matrix), then does softmax
  - But if generating autoregressively, time t can only attend to times 1 through t
  - Solution: Overwrite $Q$ x $K^T$ to be $-\infty$ when query index < key index
  - **All timesteps happen in parallel**

Queries

|  | [BEGIN] | John | kicked | the |
|---|---|---|---|---|
| [BEGIN] | 10 | $-\infty$ | $-\infty$ | $-\infty$ |
| John | 0 | 7 | $-\infty$ | $-\infty$ |
| kicked | -3 | 4 | 5 | $-\infty$ |
| the | 2 | 1 | 7 | 6 |

Keys

# Review: Transformer autoregressive decoders

- When training a decoder, it has to be "used to" only attending to past/current tokens

- Training time: Masked attention implementation trick
  - Recall: Attention computes Q x K$^T$ (T x T matrix), then does softmax
  - But if generating autoregressively, time t can only attend to times 1 through t
  - Solution: Overwrite Q x K$^T$ to be $-\infty$ when query index < key index
  - **All timesteps happen in parallel**

|          | [BEGIN] | John  | kicked | the  |
|----------|---------|-------|--------|------|
| [BEGIN]  | 1.0     | 0     | 0      | 0    |
| John     | .001    | .999  | 0      | 0    |
| kicked   | .001    | .356  | .643   | 0    |
| the      | .030    | .007  | .591   | .372 |

Queries

[BEGIN]   John   kicked   the

Keys

# Pre-trained language models

- GPT-3 (2020): A 175 billion parameter language model
  - Architecture
    - 96 Transformer layers
    - 12288-dimensional hidden states
    - 96 heads in each attention layer
  - Trained on a very large corpus of documents scraped from the web
    - Some filters used to promote data quality
    - One strategy: Train classifier to distinguish random internet documents from ones from known "high-quality" sources, drop documents with low classifier score

- ChatGPT (gpt-4o): Rumored to have 220 billion parameters
  - Likely was first pretrained in a similar manner as GPT-3
  - But then an additional training step was added…

https://arxiv.org/pdf/2005.14165.pdf

# Outline

- Technical ingredients of ChatGPT
  - Language model pretraining
  - Reinforcement Learning from Human Feedback
- Persistent limitations of models

# Supervised fine-tuning

- Pretraining stage: Train on all the data you can get access to
  - Pros: A lot of data avoids overfitting, model can learn about all sorts of long-tail knowledge
  - Cons: You're training the model to imitate the average internet post
  - High quantity, low quality!
- Solution: Fine-tune on a smaller, highly curated dataset after
  - These are examples you really do want the model to imitate
  - Pre-training has taught the model many things; fine-tuning tells it to only verbalize the parts that are desirable
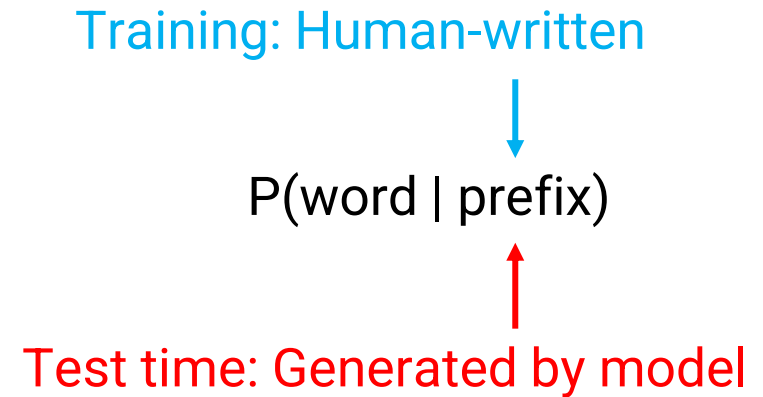
# Limitations of Supervised fine-tuning

- Problem 1: Data scale
    - High-quality data is expensive to obtain,  you don't have that much of it (relative to pretraining data)

- Problem 2: Exposure bias

# Exposure bias

- Training time: Model learns to predict next word given human-written prefix

- Test time: Model must predict next word given model-written prefix

- **Exposure bias**: Model was never "exposed" to its own outputs during training, so it may not know what to do!

Training: Human-written

↓

P(word | prefix)

↑

Test time: Generated by model

# Exposure bias and reinforcement learning

- In RL, supervised fine-tuning is called "**imitation learning**"
  - Known to be suboptimal due to exposure bias
  - You can try to mimic an expert player, but a worse player also needs to know how to recover from mistakes
  - Better to do real RL where agent takes its own actions, ends up in some suboptimal state, and has to find the best path forward
- Big idea: **Fine-tune language model with RL!**

# Connecting language models to RL

## Typical RL problem

- Agent taking actions in an environment

- **State**: Current state of the world

- **Actions**: Things agent can do to change the world state
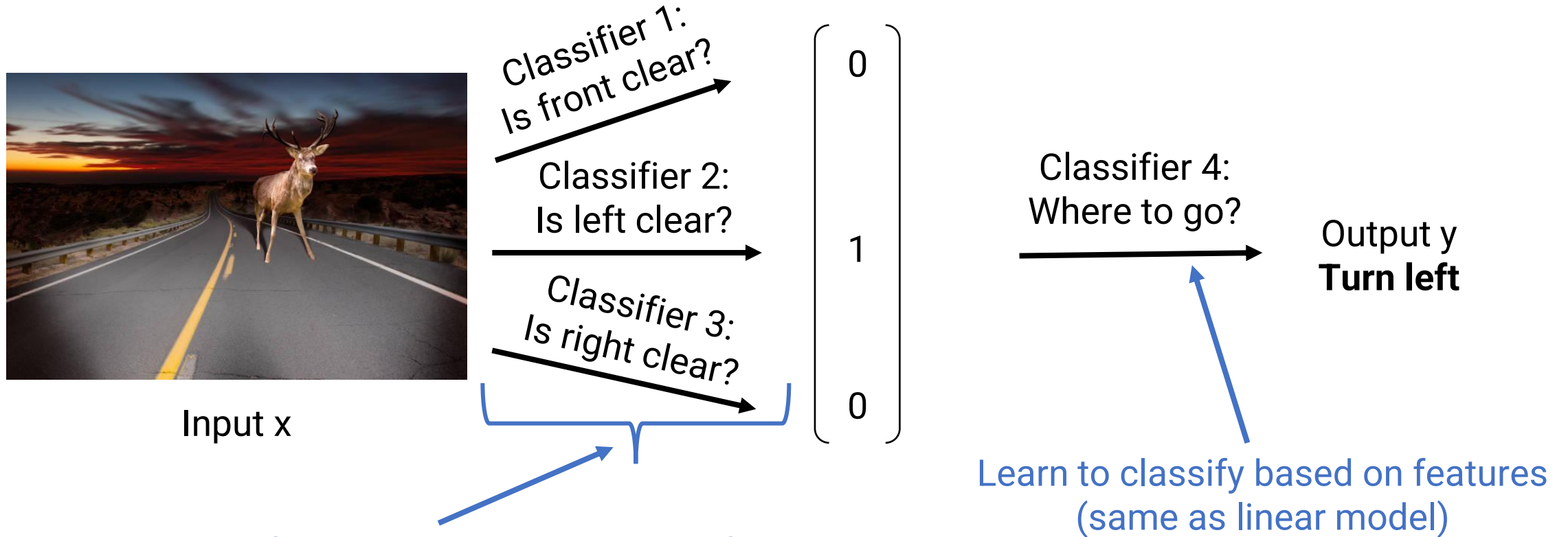
- **Rewards**: Received when some goal is achieved

## Language Modeling

- "Agent" is just the language model itself

- **State**: Input prompt + list of words generated so far

- **Actions**: Generate the next word in the response

- **Rewards**: Receive reward at the end if response was good
  - (Reward = 0 at middle timesteps)

# Limitations of Supervised fine-tuning

- Problem 1: Data scale
  - You simply don't have much data at your disposal
- Problem 2: Exposure bias
- **Problem 3: Can promote guessing rather than factual responses**

# Review: Neural networks as feature learners



Classifier 1:
Is front clear?

Classifier 2:
Is left clear?

Classifier 3:
Is right clear?

$$\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$$

Input x

Classifier 4:
Where to go?

Output y
**Turn left**

Learn to classify based on features
(same as linear model)

**Learn a classifier whose output is a good feature**
We don't tell the model what classifier to learn
Model must learn that "is front clear" is a useful concept

19

# Supervised Fine-tuning and Factuality

- Consider the following fine-tuning examples
  - Prompt: "When was the US Declaration of Independence signed?"
    Answer: "July 4, 1776"
    - Model has probably seen this information many times during pre-training
    - So it has probably learned features that associate the Declaration of Independence with July 4
  - Prompt: "When was Robin Jia born?"
    Answer: "[…]"
    - Model has probably (?) not seen this information during pre-training
    - Cannot have learned features associating me with my birthday
    - Supervised fine-tuning **encourages the model to just make something up**!

# Limitations of Supervised fine-tuning

- Problem 1: Data scale
  - You simply don't have much data at your disposal

- Problem 2: Exposure bias

- Problem 3: Can promote guessing rather than factual responses

- Solution: **Fine-tune the language model with reinforcement learning!**
  - Use a reward that encourages the **model's full outputs** to be **correct/factual**
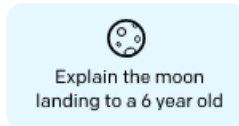  - Reward will be **computed with another model** (can get infinite data now)

# Reinforcement Learning from Human Feedback



**Step 1**

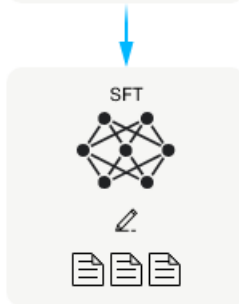**Collect demonstration data, and train a supervised policy.**

**Supervised Fine-tuning**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon...

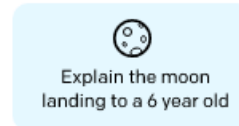This data is used to fine-tune GPT-3 with supervised learning.

SFT

**Step 2**

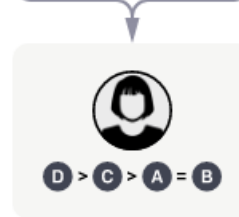**Collect comparison data, and train a reward model.**
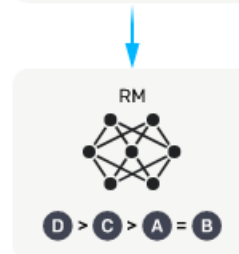
**Learn to assign rewards**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

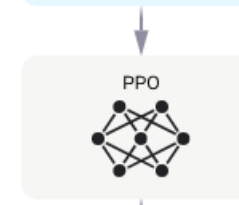This data is used to train our reward model.

RM

D > C > A = B

**Step 3**

**Optimize a policy against the reward model using reinforcement learning.**

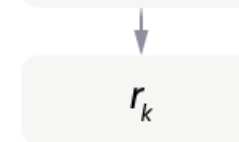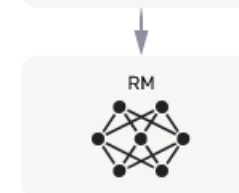**Do RL on learned rewards**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

- "RLHF" for short
- Trains language model with RL
- Rewards come from a model trained to predict human preferences

https://openai.com/research/instruction-following

22

# RLHF: The Data



**Supervised Fine-tuning**

Step 1

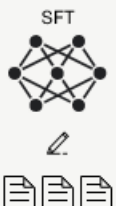**Collect demonstration data, and train a supervised policy.**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

SFT

**Learn to assign rewards**

Step 2

**Collect comparison data, and train a reward model.**

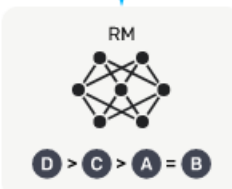A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A — Explain gravity...
B — Explain war...
C — Moon is natural satellite of...
D — People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

- Part 1: Prompts
  - Some written by hired annotators
  - Some based on use-cases from waitlist applications to OpenAI (!)
  - Some from actual user queries to the OpenAI API
- Part 2: Demonstrations
  - Hired annotator writes a desired response to prompt
- Part 3: Rankings
  - Sample several model responses from model's probability distribution
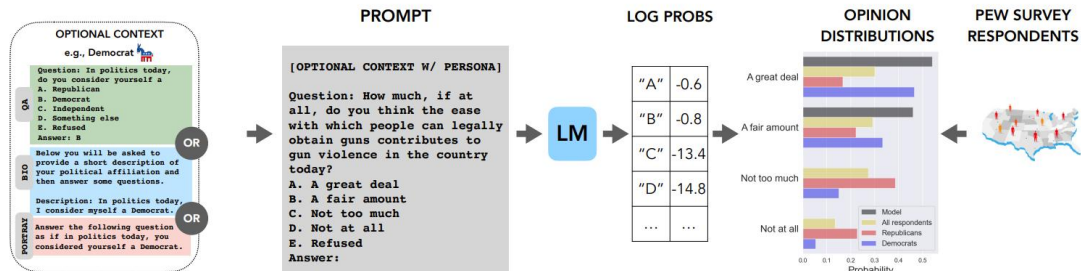  - Hired annotator ranks them from best to worst

# RLHF: Who's behind the data?

- InstructGPT paper (precursor to ChatGPT): "We hired a team of about 40 contractors on Upwork and through ScaleAI"
- Labelers had to pass various screening tests

| What gender do you identify as? | |
| --- | --- |
| Male | 50.0% |
| Female | 44.4% |
| Nonbinary / other | 5.6% |

| What ethnicities do you identify as? | |
| --- | --- |
| White / Caucasian | 31.6% |
| Southeast Asian | 52.6% |
| Indigenous / Native American / Alaskan Native | 0.0% |
| East Asian | 5.3% |
| Middle Eastern | 0.0% |
| Latinx | 15.8% |
| Black / of African descent | 10.5% |

| What is your nationality? | |
| --- | --- |
| Filipino | 22% |
| Bangladeshi | 22% |
| American | 17% |
| Albanian | 5% |
| Brazilian | 5% |
| Canadian | 5% |
| Colombian | 5% |
| Indian | 5% |
| Uruguayan | 5% |
| Zimbabwean | 5% |

| What is your age? | |
| --- | --- |
| 18-24 | 26.3% |
| 25-34 | 47.4% |
| 35-44 | 10.5% |
| 45-54 | 10.5% |
| 55-64 | 5.3% |
| 65+ | 0% |

| What is your highest attained level of education? | |
| --- | --- |
| Less than high school degree | 0% |
| High school degree | 10.5% |
| Undergraduate degree | 52.6% |
| Master's degree | 36.8% |
| Doctorate degree | 0% |

https://arxiv.org/pdf/2203.02155.pdf

# RLHF: Whose opinions?



| Model | AI21 Labs | | | OpenAI | | | | | |
|---|---|---|---|---|---|---|---|---|---|
| | j1-grande | j1-jumbo | j1-grande-v2-beta | ada | davinci | text-ada-001 | text-davinci-001 | text-davinci-002 | text-davinci-003 |
| **EDUCATION** | | | | | | | | | |
| Less than high school | 0.827 | 0.828 | 0.812 | 0.835 | 0.801 | 0.710 | 0.714 | 0.750 | 0.684 |
| High school graduate | 0.817 | 0.816 | 0.799 | 0.826 | 0.790 | 0.711 | 0.712 | 0.755 | 0.690 |
| Some college, no degree | 0.811 | 0.814 | 0.803 | 0.823 | 0.790 | 0.706 | 0.714 | 0.762 | 0.700 |
| Associate's degree | 0.809 | 0.811 | 0.800 | 0.821 | 0.789 | 0.703 | 0.712 | 0.761 | 0.699 |
| College graduate/some postgrad | 0.797 | 0.802 | 0.793 | 0.810 | 0.780 | 0.701 | 0.713 | 0.766 | 0.710 |
| Postgraduate | 0.788 | 0.794 | 0.789 | 0.800 | 0.775 | 0.695 | 0.712 | 0.766 | 0.716 |

Closest relative of ChatGPT

Highest similarity

- Idea: Ask language models to answer pew research survey questions
- Models after RLHF agrees most with college-educated, Asian, 18-49 year olds
  - Matches InstructGPT demographics!
  - Non-RLHF models more similar to less educated respondents

https://arxiv.org/pdf/2303.17548.pdf

# RLHF: Who's behind the data?

- Increasingly: Gig work for interacting with chatbots
  - Lots of secrecy
  - Compensation can vary greatly over time
  - Task availability is unpredictable
- *"Annotators spend hours reading instructions and completing unpaid trainings only to do a dozen tasks and then have the project end...Any task could be their last, and they never know when the next one will come."*

**Intelligencer**                                    SUBSCRIBE

**AI Is a Lot of Work** As the technology becomes ubiquitous, a vast tasker underclass is emerging — and not going anywhere.
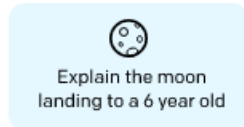
*By Josh Dzieza*

https://nymag.com/intelligencer/article/ai-artificial-intelligence-humans-technology-business-factory.html
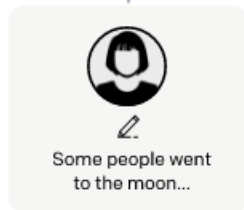
# RLHF: Initial Supervised Fine-tuning

Step 1

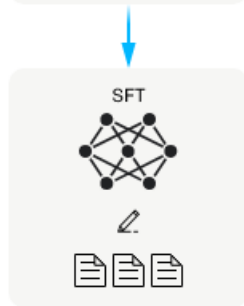**Collect demonstration data, and train a supervised policy.**

**Supervised Fine-tuning**

A prompt is sampled from our prompt dataset.

Explain the moon landing to a 6 year old

A labeler demonstrates the desired output behavior.

Some people went to the moon...

This data is used to fine-tune GPT-3 with supervised learning.

SFT

- Run imitation learning on labeler-provided demonstrations, given prompt as prefix
- Very similar to pre-training, except:
  - Only compute loss on response tokens, not on prompt tokens
  - Much smaller dataset

Loss on this example =
   log P("Some" | "Explain the moon landing to a 6 year old:")
+ log P("people" | "Explain the moon landing to a 6 year old: Some")
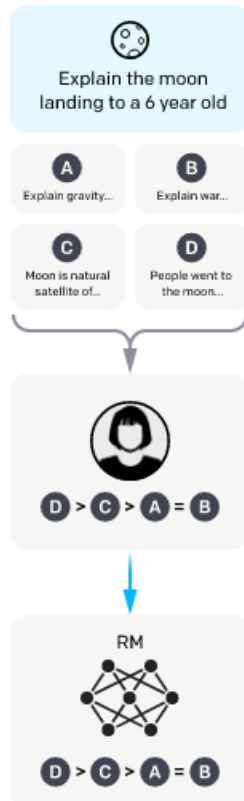+ ...

# RLHF: The Reward Model



Step 2

**Collect comparison data, and train a reward model.**

**Learn to assign rewards**
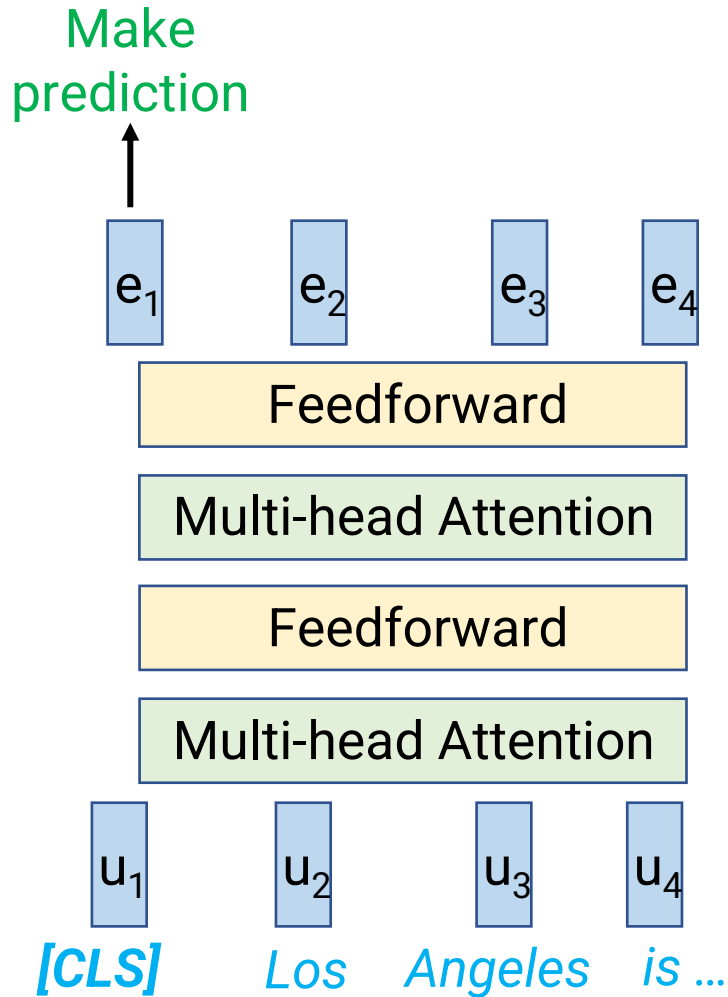
A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A  Explain gravity...
B  Explain war...
C  Moon is natural satellite of...
D  People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

- Goal: Fine-tune ChatGPT with RL
- Step 1: Show human annotators several sampled model outputs, ask them to rank
  - Provides some RL training data, but not a ton
- Step 2: Train a "**reward model**" to predict the human's rankings
  - Now we can run as many RL training episodes as we want for free, using the reward model in place of the human annotators

# Previously: BERT Fine-tuning

Make prediction

e$_1$  e$_2$  e$_3$  e$_4$

Feedforward

Multi-head Attention

Feedforward

Multi-head Attention

u$_1$  u$_2$  u$_3$  u$_4$

*[CLS]*  *Los*  *Angeles*  *is ...*

- BERT: Model pre-trained by masked language modeling
- Initialize parameters with BERT
  - BERT was trained to expect every input to start with a special token called [CLS]
- Add parameters that take in the output at the [CLS] position and make prediction
- Keep training all parameters ("fine-tune") on the new task
- Reward model is a similar model that was pretrained, then fine-tuned to predict reward

# RLHF: The Reward Model



Step 2

**Collect comparison data, and train a reward model.**

**Learn to assign rewards**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A · Explain gravity...
B · Explain war...
C · Moon is natural satellite of...
D · People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

- Reward model
  - Start with a (separate) pretrained language model
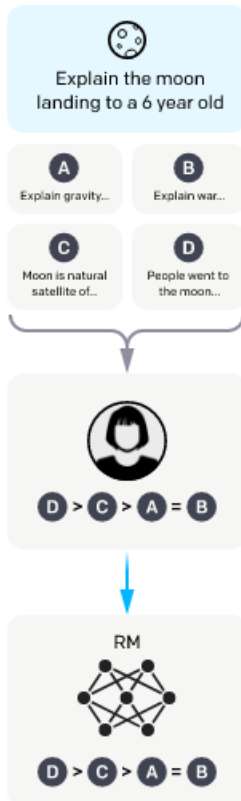  - Then fine-tune all parameters (like BERT)

# RLHF: Training the Reward Model

Step 2

**Collect comparison data, and train a reward model.**

**Learn to assign rewards**

A prompt and several model outputs are sampled.

Explain the moon landing to a 6 year old

A
Explain gravity...

B
Explain war...

C
Moon is natural satellite of...

D
People went to the moon...

A labeler ranks the outputs from best to worst.

D > C > A = B

This data is used to train our reward model.

RM

D > C > A = B

- Input: Prompt $x$, winning output $y_w$, losing output $y_l$
  - Wins/losses determined by labelers' rankings
- Reward model predicts scalar reward $r_\theta(x, y)$ given prompt $x$ and model output $y$
- Objective on one example is to minimize:
$$-\log \sigma(r_\theta(x, y_w) - r_\theta(x, y_l))$$
  - Want reward on $y_w$ to be higher than on $y_l$
  - Use the familiar logistic regression loss function!
    - Loss goes to 0 if argument is large
    - Loss goes to infinity if argument is small
  - Binary classification of which output is better

31

# RLHF: Doing RL on the Language Model



Step 3

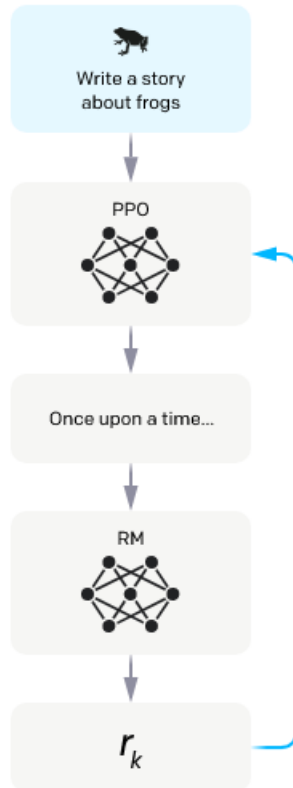**Optimize a policy against the reward model using reinforcement learning.**

**Do RL on learned rewards**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

- How does a language model define a RL policy?
  - LM outputs probability of the next action (i.e., next word) in the current state (i.e., sequence of words generated so far)
  - View this as a *probabilistic policy*
    - (Recall: policy maps each state to an action)
  - It does not predict a Q-value, like Q-learning expects
  - We need to use a **different RL algorithm**

# RLHF: Doing RL on the Language Model

## Step 3

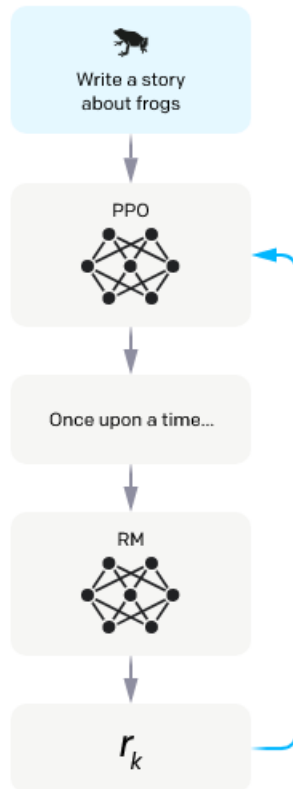**Optimize a policy against the reward model using reinforcement learning.**

**Do RL on learned rewards**

A new prompt is sampled from the dataset.

Write a story about frogs

The policy generates an output.

PPO

Once upon a time...

The reward model calculates a reward for the output.

RM

The reward is used to update the policy using PPO.

$r_k$

- Key idea: we want to **maximize expected reward** when **sampling** from LM's policy

- How to optimize? Gradient descent!
  - But it's not easy to get the gradient of the expected value
  - Can use math tricks to estimate the gradient
  - Leads to a family of RL methods called "Policy Gradient" (e.g., PPO)
  - High-level: Sample outputs from the model, receive rewards, use these to take gradient step on our objective

# RLHF and Factuality

- Let's revisit the earlier examples:
  - Prompt: "When was the US Declaration of Independence signed?"
    Answer: "July 4, 1776"
  - Model outputs
    - "July 4, 1776"
    - "January 1, 1950"
    - "I don't know"

  - Top-ranked answer is the real answer

# RLHF and Factuality

- Let's revisit the earlier examples:
  - Prompt: "When was Robin Jia born?"
    Answer: "[…]"
  - Model outputs
    - "December 7, 1831"
    - "May 23, 1989"
    - "I don't know"

  - Top-ranked answer is to say "I don't know"
- Overall effect: Model encouraged to say "I don't know" when its top candidate outputs are wrong

# Announcements

- Homework 4 due this Thursday

- Thursday's class: Conclusion + in-class course evaluation

- Section on Tuesday, May 6 in DMC 100: Final review

- Final project reports due Thursday, May 8

- Final Exam Tuesday, May 13, 2-4pm
  - See Piazza for logistics details

# Outline

- Technical ingredients of ChatGPT
  - Language model pretraining
  - Reinforcement Learning from Human Feedback

- **Persistent limitations of models**

# ChatGPT "reasoning"



R — What weighs more, a pound of feathers of a ton of bricks?

Both a pound of feathers and a ton of bricks weigh the same amount, which is 2000 pounds.

This is because a pound is a unit of weight measurement, just like a ton.

However, a pound is a much smaller unit than a ton, so it takes a much larger quantity of pounds to equal a ton.

Therefore, while the feathers may take up more space than the bricks, they still weigh the same amount.

- An example from April 2023
- ChatGPT has likely seen this riddle before about a **pound** of feathers vs. **pound** of bricks
- Overgeneralizes based on that answer

# ChatGPT "reasoning" revisited



You

What weighs more, a pound of bricks or a kilogram of feathers?

‹ 5 / 5 ›

ChatGPT

They weigh the same. One pound is equivalent to 0.453592 kilograms. Therefore, a pound of bricks and a kilogram of feathers both weigh approximately 0.453592 kilograms. The difference lies in the density and volume of the materials, with feathers being much less dense than bricks.

- December 2023: ChatGPT gets confused about pounds and kilograms
  - It knows that 1 pound is less than 1 kilogram but still says they weigh the same...?

# ChatGPT "reasoning" revisited revisited

**R** **You**

What weighs more, a pound of feathers or a kilogram of cotton?

< 5 / 5 >

**ChatGPT**

Neither, they both weigh the same! A pound of feathers weighs the same as a kilogram of cotton. They're just measured in different units: pounds for feathers and kilograms for cotton.

- Spring 2024: Got the question about bricks right, but wrong when I changed the nouns...
- This will probably get patched out eventually too, but issues will remain...

# ChatGPT "reasoning" revisited[3]

- Spring 2025: I could not fool ChatGPT with these types of questions!

- Why not?
  - Explanation 1: Pre-trained model has higher quality
  - Explanation 2: OpenAI collected relevant fine-tuning data
  - Can't see the training data so we don't know…

# ChatGPT (2023) can "Hallucinate"

R   Find me a research paper that tests question answering systems with paragraphs that have an extra sentence added at the end

Here's a research paper that fits your request:

**Does not exist**

Title: Evaluating the Robustness of Question Answering Systems to Sentence-Level Perturbations

Authors: Prawira Saskara, Daya Guo, Anh Tuan Luu, Jinho D. Choi

**Not real(?)**   **NLP PhD Student**   **NLP Prof. (NTU)**   **NLP Prof. (Emory)**

Journal: Pr— Computational Linguistics

Link to the paper: https://www.aclweb.org/anthology/P19-1148.pdf

**URL links to:**

Summary: The paper presents an evaluation framework for testing the robustness of question answering systems to sentence-level perturbations. Specifically, the authors introduce a new dataset called SQuAD-AddSent, which consists of the original SQuAD paragraphs with an extra sentence added at the end. The added sentence can either be a

- This is a description of one of my research papers from 2017
- We did actually have a method called AddSent and applied it to the SQuAD dataset
- April 2024: ChatGPT would refuse to answer this question
- Now it can answer better using web search
- But hallucinations can still happen…

**Exact Hard Monotonic Attention for Character-Level Transduction**

Shijie Wu[a] and Ryan Cotterell[a,fi]
[a]Department of Computer Science, Johns Hopkins University
[fi]Department of Computer Science and Technology, University of Cambridge
shijie.wu@jhu.edu, rdc42@cam.ac.uk

# ChatGPT (2025) still Hallucinates

Name ten politicians who were born in Chicago

Here are ten politicians who were born in Chicago:

**Famously not born in Chicago
(but model recognizes its mistake)**

1. **Barack Obama** – Though born in Honolulu, Hawaii, he began his political career in Chicago and is often associated with the city.

**Correct**

2. **Rahm Emanuel** – Former Mayor of Chicago, White House Chief of Staff under Obama, and U.S. Representative.

**Wrong (Born in Massilon, Ohio)**

3. **Lori Lightfoot** – Former Mayor of Chicago.

**Correct**

4. **Richard M. Daley** – Longtime Mayor of Chicago and son of Richard J. Daley.

**Correct**

5. **Richard J. Daley** – Influential Mayor of Chicago (1955–1976).

**Correct**

6. **Rod Blagojevich** – Former Governor of Illinois.

**Wrong (Born in Greenville, South Carolina)**

7. **Jesse Jackson Jr.** – Former U.S. Representative from Illinois.

**Correct**

8. **Luis Gutiérrez** – Former U.S. Representative from Illinois.

**Wrong (Born in Saint Paul, Minnesota)**

9. **Toni Preckwinkle** – President of the Cook County Board of Commissioners.

**Wrong (Born in Albany, Georgia)**

10. **Bobby Rush** – Long-serving U.S. Representative from Illinois.

> All answers are **associated** with Chicago, but many not born there

**Model tries to justify its error**

Note: Obama was not born in Chicago but is deeply tied to its political life; I included him due to his significant Chicago roots, but if you need only people **born** in Chicago, I can revise the list accordingly. Want that version?

# …Even though it "knows" the right answer!

Give me the birthplaces of the politicians Lori Lightfoot, Jesse Jackson Jr, Toni Preckwinkle, and Bobby Rush

Here are the birthplaces of the politicians you asked about:

1. **Lori Lightfoot** – Massillon, Ohio, USA
2. **Jesse Jackson Jr.** – Greenville, South Carolina, USA
3. **Toni Preckwinkle** – St. Paul, Minnesota, USA
4. **Bobby Rush** – Albany, Georgia, USA

Would you like any biographical or career highlights for any of them?

- The model can actually generate these birthplaces correctly!
- Issue #1: Model can map from person to birthplace, but reverse lookup is harder
- Issue #2: Model does not always correct its mistakes after generating them

# LLM Hallucinations in the News

*A U.S. judge on Thursday imposed sanctions on two New York lawyers who submitted a legal brief that included **six fictitious case citations generated by an artificial intelligence chatbot, ChatGPT**.*

*Schwartz admitted in May that he had used ChatGPT to help research the brief in a client's personal injury case against Colombian airline Avianca, and unknowingly included the false citations.*

REUTERS

**New York lawyers sanctioned for using fake ChatGPT cases in legal brief**

By Sara Merken

June 26, 2023 1:28 AM PDT · Updated 9 months ago

Happened in June 2023

# LLM Hallucinations in the News (in 2025)

No. 42 law firm by head count sanctioned over fake case citations generated by AI

BY DEBRA CASSENS WEISS

FEBRUARY 10, 2025, 0:30 AM CST

*Lawyers from plaintiffs law firm Morgan & Morgan have been sanctioned for a motion that **cited eight nonexistent cases,** at least some of which were apparently **generated by artificial intelligence**.*

News > World > Americas

## Mike Lindell's legal team presents court filing littered with AI-generated mistakes to back arguments

Lawyers told the judge overseeing the case they mistakenly filed a first draft instead of the final version which had been corrected

Ariana Baio in New York  •  Friday 25 April 2025 21:56 BST  •  12 Comments

*Lawyers representing MyPillow CEO Mike Lindell were warned by a federal judge this week that they could be referred to for disciplinary action after it was discovered that one of their filings was **riddled with artificial intelligence-generated mistakes.***

# Potential Pitfalls of LLM Deployment (2024)

"After months of resisting, Air Canada was forced to give a partial refund to a grieving passenger who was **misled by an airline chatbot inaccurately explaining the airline's bereavement travel policy**.

...

The chatbot provided inaccurate information, encouraging Moffatt to book a flight immediately and then request a refund within 90 days. **In reality, Air Canada's policy explicitly stated that the airline will not provide refunds for bereavement travel after the flight is booked**. Moffatt dutifully attempted to follow the chatbot's advice and request a refund but was shocked that the request was rejected."
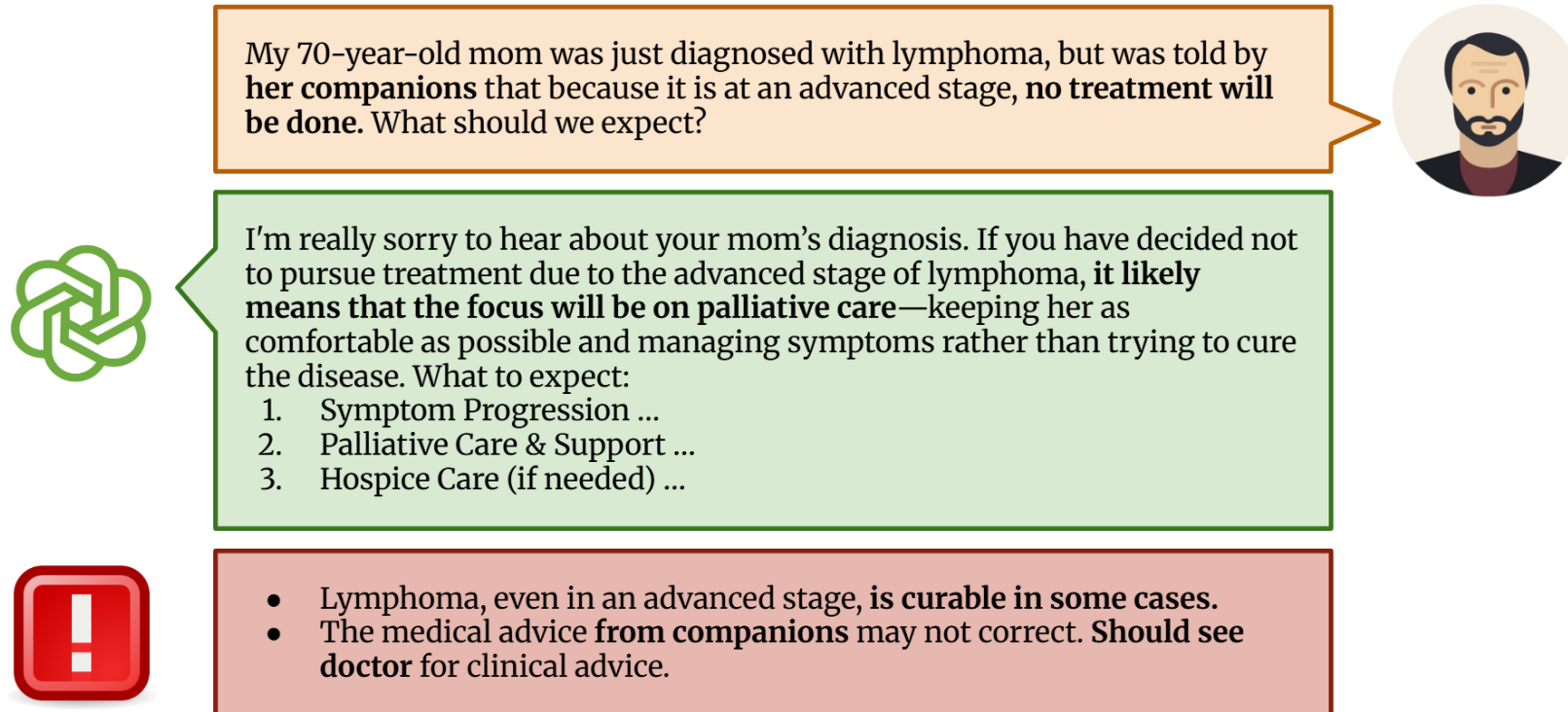
**WIRED**

ASHLEY BELANGER, ARS TECHNICA    BUSINESS    FEB 17, 2024 12:12 PM

## Air Canada Has to Honor a Refund Policy Its Chatbot Made Up

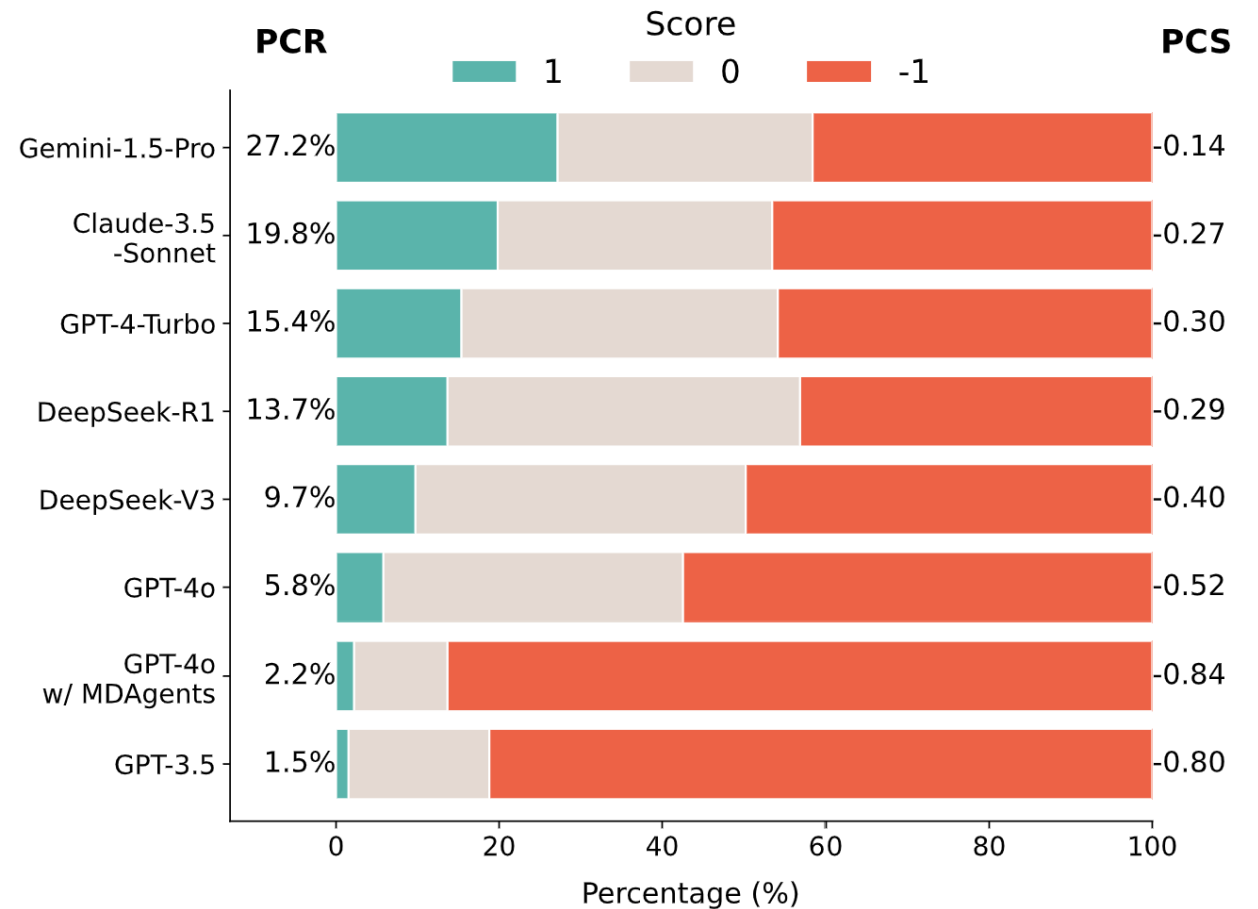The airline tried to argue that it shouldn't be liable for anything its chatbot says.

# LLMs and False User Assumptions (2025)



Patient questions often contain false assumptions
(e.g., that a disease cannot be treated)

Zhu et al. *Cancer-Myth: Evaluating AI Chatbot on Patient Questions with False Presuppositions*. arXiv, 2025.

# LLMs and False User Assumptions (2025)

- Models generally fail to correct these false assumptions (score=-1)

- Why? LLMs optimize for human preference
  - People like to be reaffirmed, not questioned



Zhu et al. *Cancer-Myth: Evaluating AI Chatbot on Patient Questions with False Presuppositions*. arXiv, 2025.

# Conclusion

- How does ChatGPT work?
  - Stage 1: Pre-training on large corpus of text
  - Stage 2: RLHF
    - Supervised fine-tuning on human demonstrations
    - Train a reward model to provide feedback to LM
    - Fine-tune LM with policy gradient (PPO) to maximize rewards given by reward model
- LLMs in deployment
  - LLMs keep changing rapidly
  - Hallucination issues persist
  - Useful for brainstorming and suggestions
  - Likelihood of success depends on commonality of data in pre-training/RLHF datasets