# The Secret Sauce in ChatGPT

**Minhao CHENG**
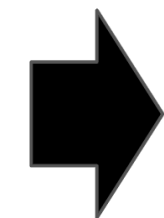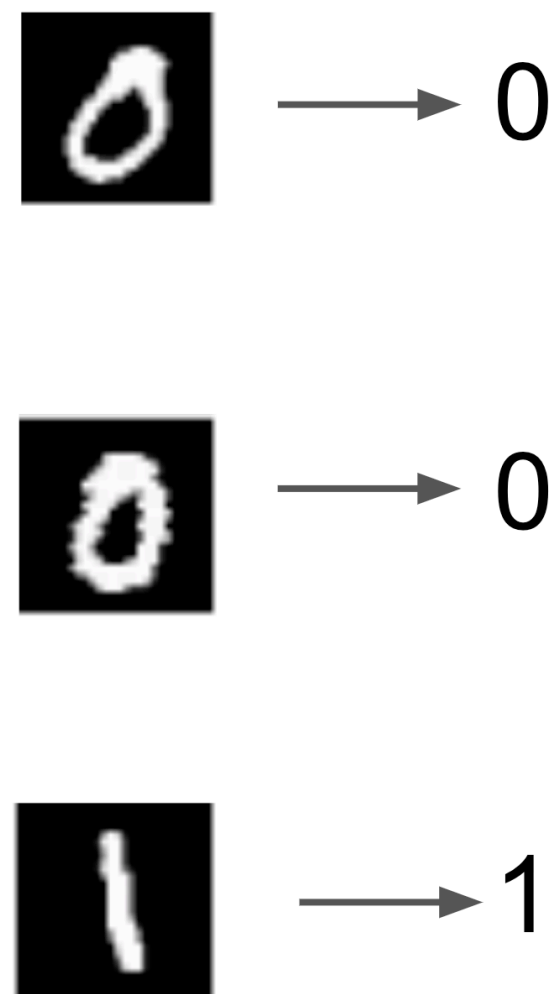
THE DEPARTMENT OF
**COMPUTER SCIENCE & ENGINEERING**
計算機科學及工程學系

# Machine Learning Overview
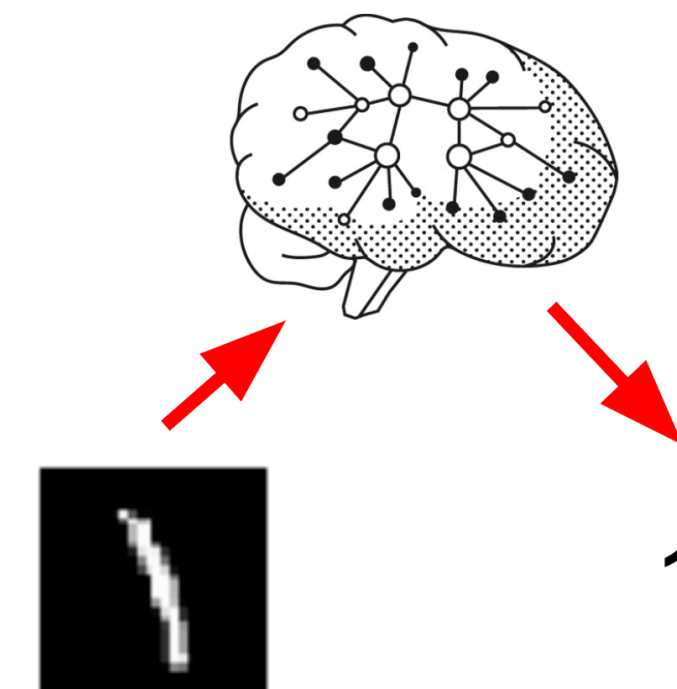## From learning to machine learning
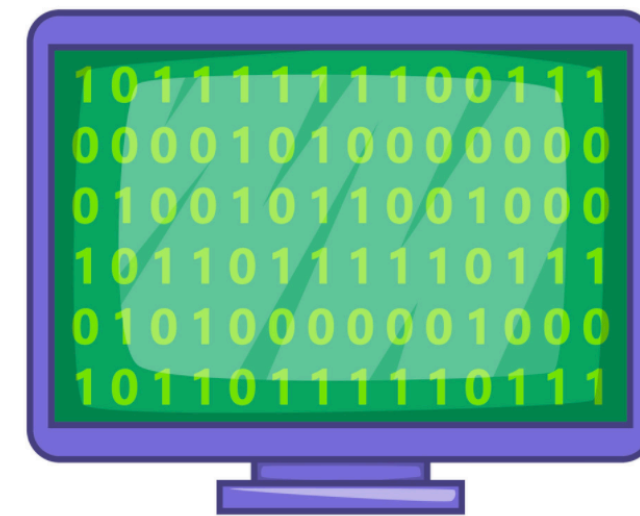
- Human learning

**Observation**     **Learning**     **Decision rule**

# Machine Learning Overview

## Machine learning

**Training Data**          **Machine Learning**          **Decision rule**

# Machine Learning Overview

## Machine learning

**Training Data**     **Machine Learning**     **Decision rule**



$x_1$ $\longrightarrow$ 0

$x_1$     $y_1$

$x_2$ $\longrightarrow$ 0

$x_2$     $y_2$

$x_3$ $\longrightarrow$ 1

$x_3$     $y_3$

$x_1$: vector of pixel values [0, 24, 128, ...]
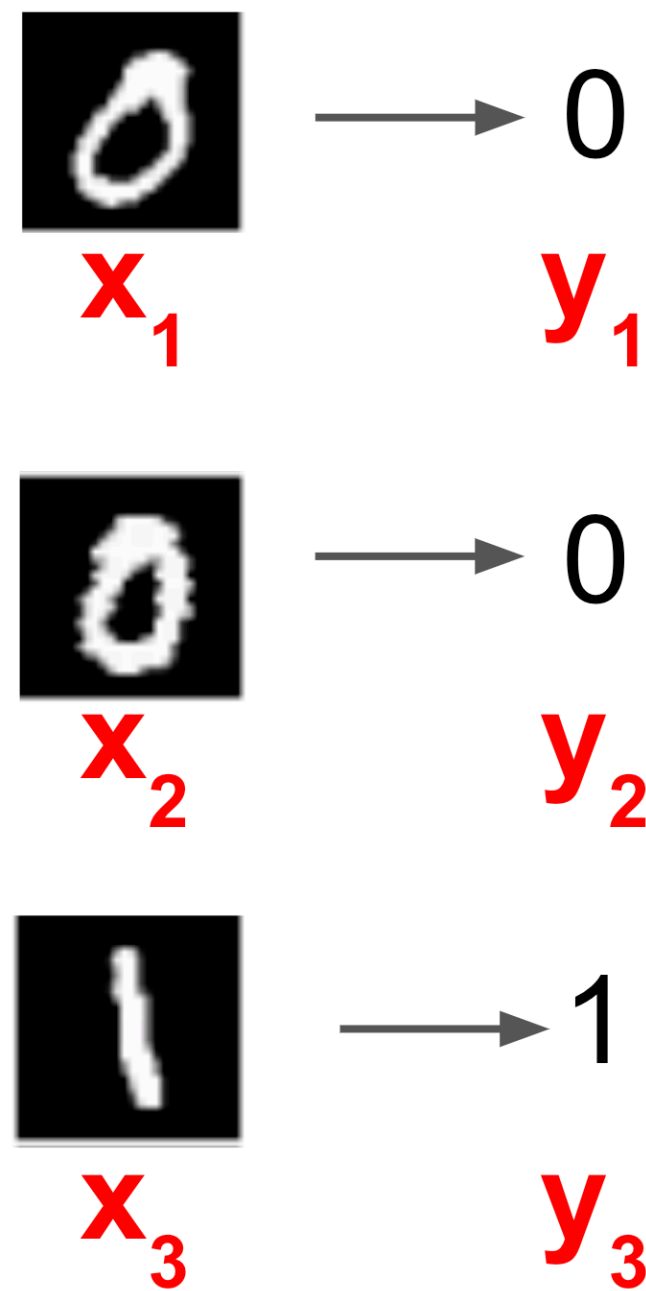
$y_1$: 0 or 1

# Machine Learning Overview

## Machine learning

**Training Data**    **Machine Learning**    **Decision rule**



$x_1$     $y_1$

$x_2$     $y_2$

$x_3$     $y_3$

**g (a function)**
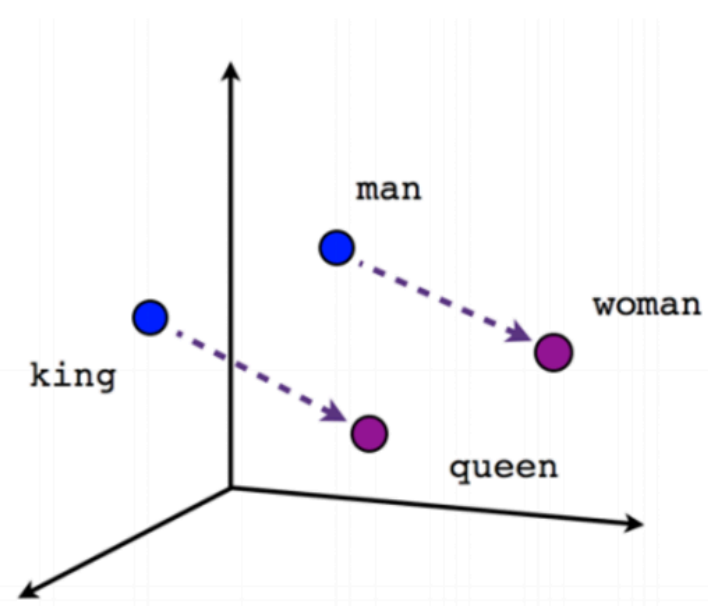
x     0 or 1

**g maps any image (vector) to 0/1**

# How to learn language?
## Word embedding

- Computers doesn't know word/character, how to represent word?

  - Through word embedding!

  - Maps word to some vectors in the high dimensional space

- What task should we assign model to learn?



Male-Female

Verb tense

Country-Capital

# How to learn language?
**Pretraining**

- Choose the one we want to get the best performance?

  - There are billions of tasks

  - The model performs good in one task could be bad in another tasks

    - Eg. Food rating  -> paper rating -> tell a story?

  - The training data we have for our downstream task (must be sufficient to teach all contextual aspects of language.

- We need to find a "common sense" task

# How to learn language?
## What can we learn from reconstructing the input?

- HKUST is located in _____, Hong Kong

- I went to the ocean to see the fish, turtles, seals, and _____.

- Overall, the value I got from the two hours watching it was the sum total of the popcorn and the drink. The movie was ____.

- I was thinking about the sequence that goes 1, 1, 2, 3, 5, 8, 13, 21, _____

# Language model

- Model $p_\theta(w_t | w_{1:t-1})$, the probability distribution over words given their past contexts.

  - There's lots of data for this! (No need for labeling)

- Pretraining through language modeling:

  - Train a neural network to perform language modeling on a large amount of text.

  - Save the network parameters.

goes    to    make  tasty   tea   END

Decoder
(Transformer, LSTM, ++ )

Iroh    goes    to    make  tasty   tea

# Pretrained Language model

**Step 1: Pretrain (on language modeling)**

Lots of text; learn general things!

goes    to    make   tasty    tea    END

(Transformer, LSTM, ++ )

Iroh    goes    to    make   tasty    tea

**Step 2: Finetune (on your task)**

Not many labels; adapt to the task!

☺/☹

(Transformer, LSTM, ++ )

*... the movie was ...*

# Pretrained Language model
## Why it works

- Language tasks are correlated with each other

- In a optimization perspective, stochastic gradient descent sticks (relatively) close to the initialization point

  - Train from scratch = random initialization

  - Finetuning: find a good local minima near a good initialization

# Language model ≠ assisting users

PROMPT    *Explain the moon landing to a 6 year old in a few sentences.*

COMPLETION    GPT-3

Explain the theory of gravity to a 6 year old.

Explain the theory of relativity to a 6 year old in a few sentences.

Explain the big bang theory to a 6 year old.

Explain evolution to a 6 year old.

Language models are not *aligned* with user intent [Ouyang et al., 2022].

# Instruction finetuning
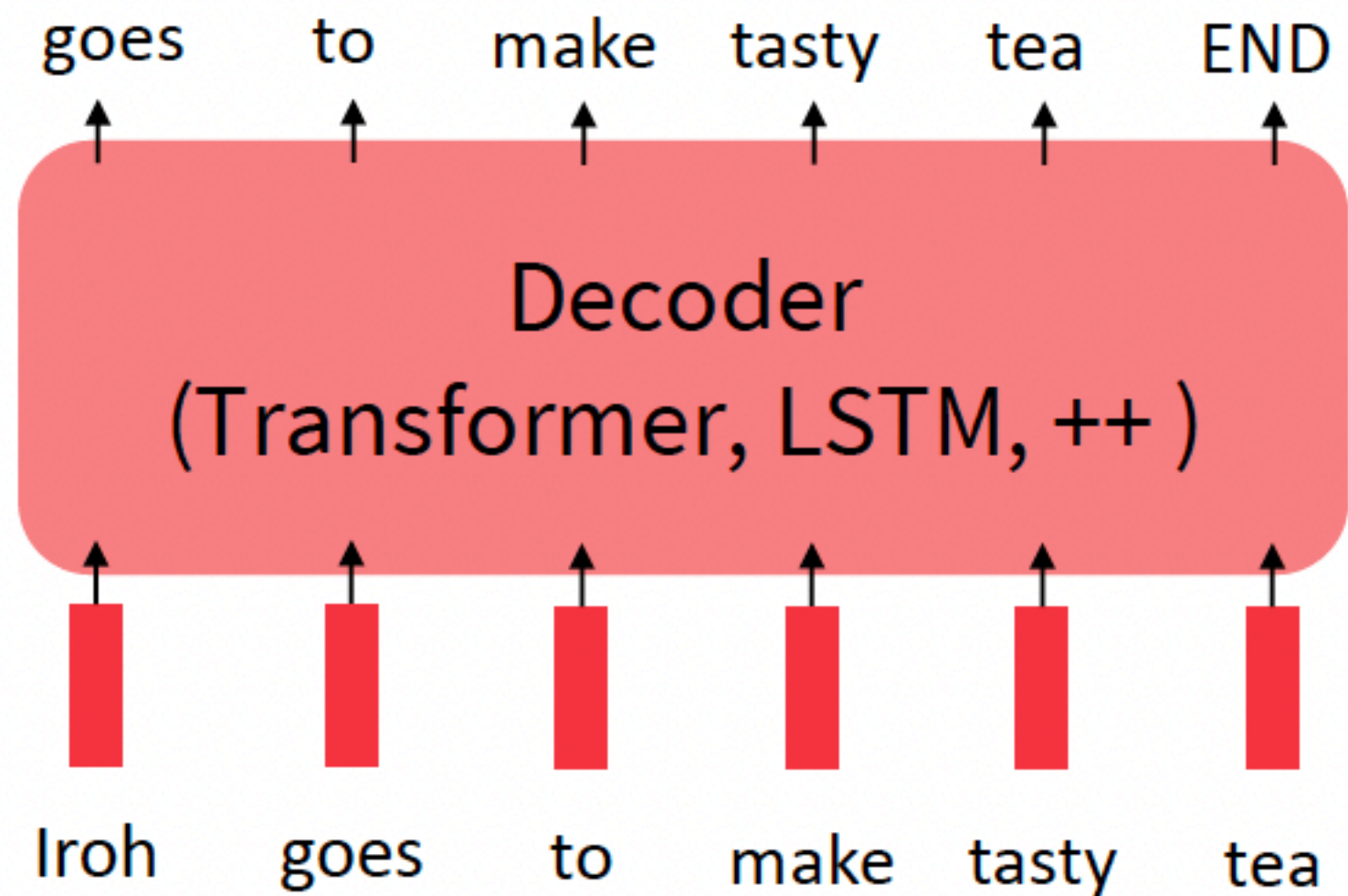
Pretraining can improve NLP applications by serving as parameter initialization.

**Step 1: Pretrain (on language modeling)**

Lots of text; learn general things!

goes    to    make    tasty    tea    END

Decoder
(Transformer, LSTM, ++ )

Iroh    goes    to    make    tasty    tea

**Step 2: Finetune (on many tasks)**

~~Not~~ many labels; adapt to the tasks!

☺/☹

Decoder
(Transformer, LSTM, ++ )

... *the movie was* ...

# Instruction finetuning

- **Collect examples** of (instruction, output) pairs across many tasks and finetune an LM



- Evaluate on **unseen tasks**

# Instruction ~~finetuning~~ pretraining?

- As is usually the case, **data + model scale** is key for this to work!
- For example, the **Super-NaturalInstructions** dataset contains **over 1.6K tasks, 3M+** examples
  - Classification, sequence tagging, rewriting, translation, QA...
- **Q:** how do we evaluate such a model?

- Through benchmarks in multitask LM
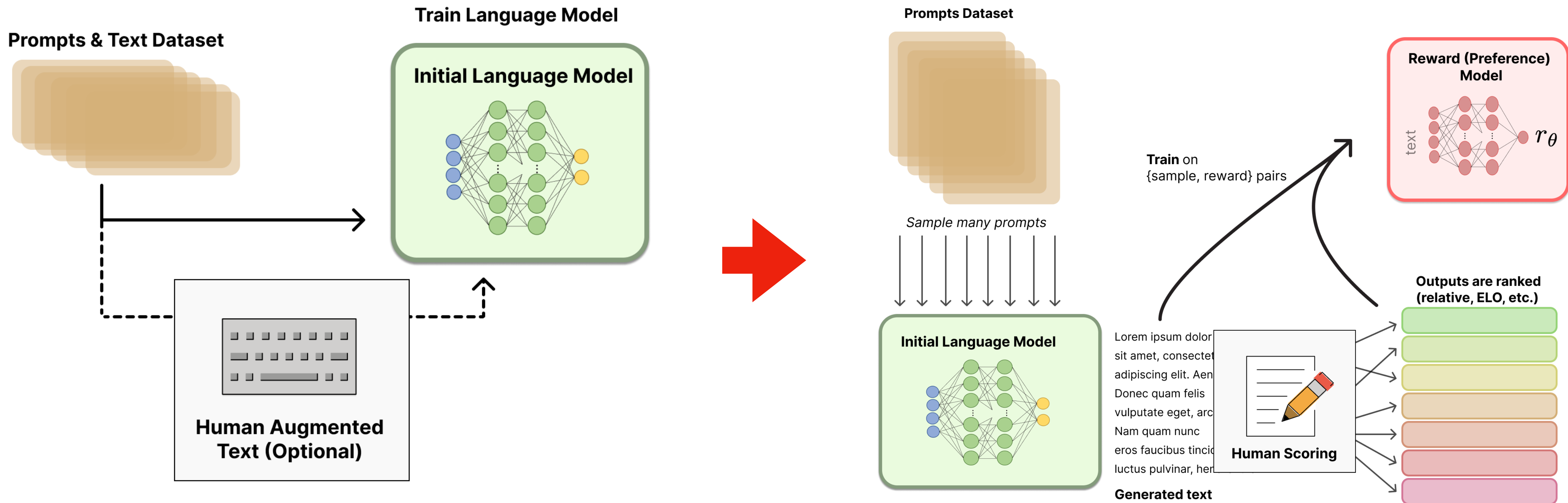
[Wang et al., 2022]

# Instruction fine-tuning
## Limitations

- It's **expensive** to collect ground-truth data for tasks

- Open-ended generation have no right answer

  - Write a story about traveling to HKUST using airplane

  - Where to travel for the next holiday?

- Language modeling penalizes all token-level mistakes equally, but some errors are worse than others.

- Can we **explicitly attempt to satisfy human preferences**?

# Reinforcement learning from human feedback
## RLHF

**Prompts & Text Dataset**

**Train Language Model**

**Initial Language Model**

**Human Augmented Text (Optional)**

**Prompts Dataset**

*Sample many prompts*

**Initial Language Model**

Lorem ipsum dolor sit amet, consectet adipiscing elit. Aen Donec quam felis vulputate eget, arc Nam quam nunc eros faucibus tincid luctus pulvinar, her

**Generated text**

**Human Scoring**

**Train** on {sample, reward} pairs

**Reward (Preference) Model**

$r_\theta$

text

**Outputs are ranked (relative, ELO, etc.)**

# Reinforcement learning from human feedback
## RLHF

- For each sample $s$, we had a way to obtain a human reward $R(s) \in \mathbb{R}$, higher is better

```
SAN FRANCISCO,
California (CNN) --
A magnitude 4.2
earthquake shook the
San Francisco
...
overturn unstable
objects.
```

```
An earthquake hit
San Francisco.
There was minor
property damage,
but no injuries.
```
$$s_1$$
$$R(s_1) = 8.0$$

```
The Bay Area has
good weather but is
prone to
earthquakes and
wildfires.
```
$$s_2$$
$$R(s_2) = 1.2$$

- We want to maximize the expected reward

# RLHF
## Problems&Sol

- Problem1: Expensive to get human evaluation

  - Sol: Train another model to predict human preferences

- Problem 2: human judgements are noisy and miscalibrated!

  - Sol: Just ask for pairwise comparisons

An earthquake hit San Francisco. There was minor property damage, but no injuries.

$s_1$

$R(s_1) = 8.0$

The Bay Area has good weather but is prone to earthquakes and wildfires.

$s_2$

$R(s_2) = 1.2$

Train an LM $RM_\phi(s)$ to predict human preferences from an annotated dataset, then optimize for $RM_\phi$ instead.
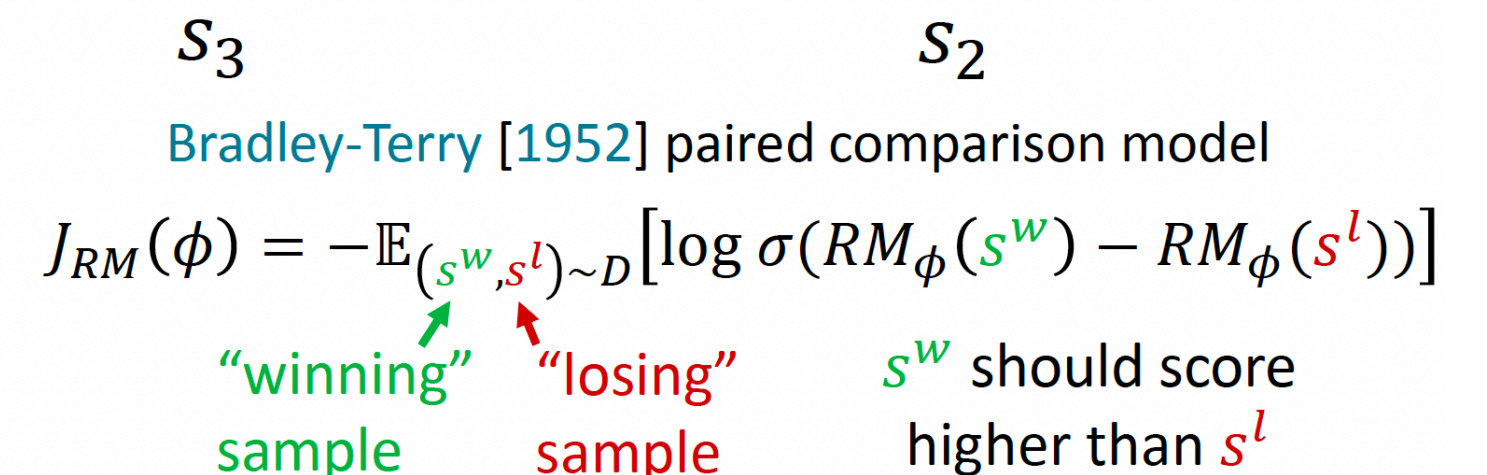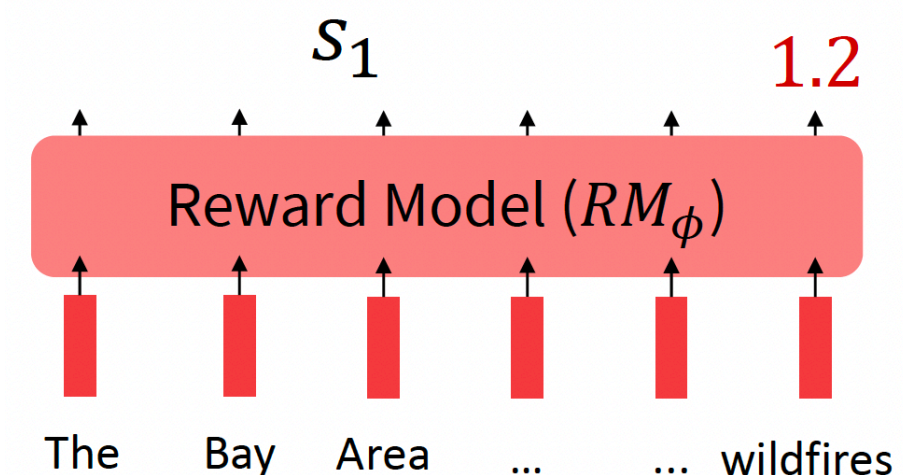
An earthquake hit San Francisco. There was minor property damage, but no injuries.

$s_1$    1.2

$>$

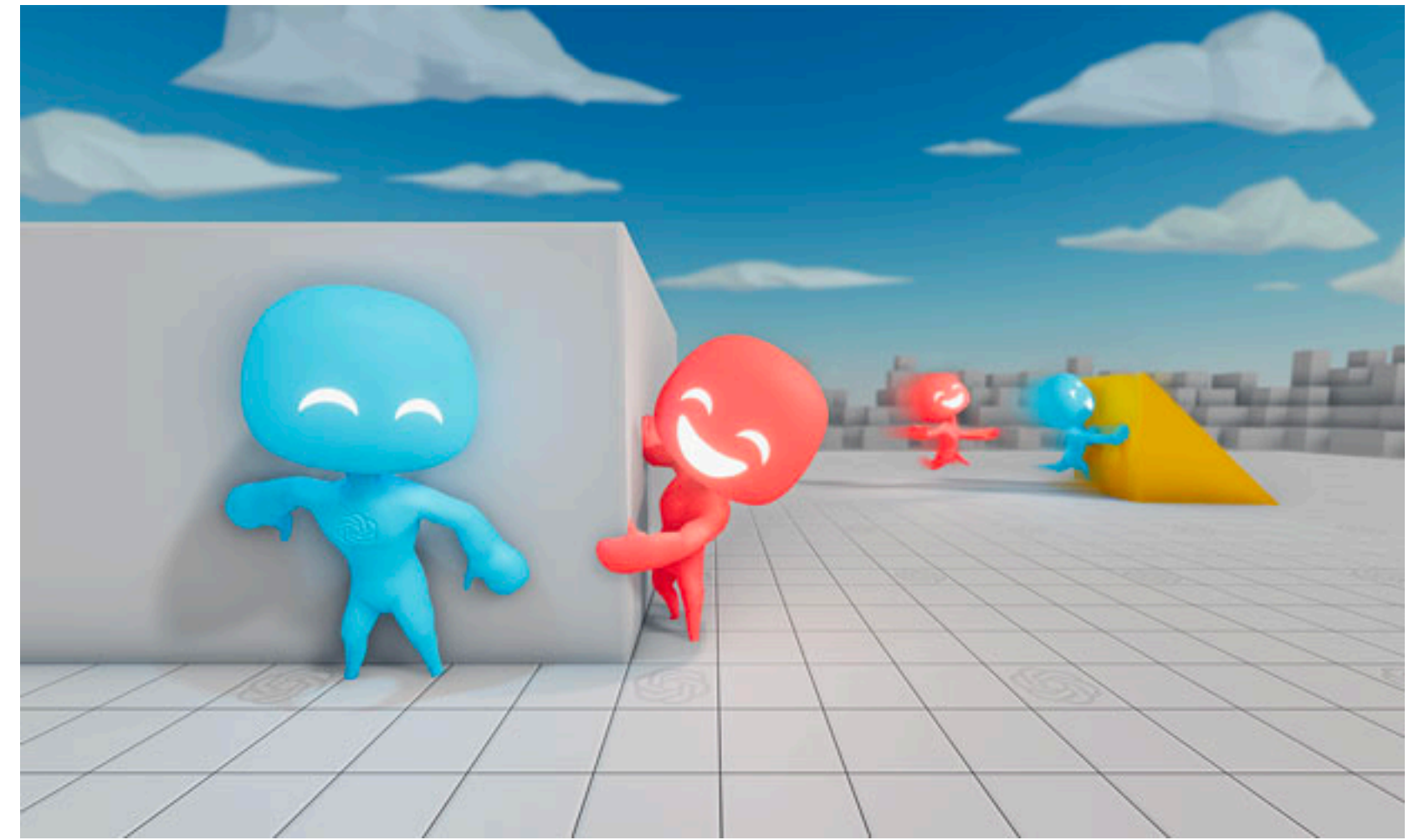A 4.2 magnitude earthquake hit San Francisco, resulting in massive damage.

$s_3$

$>$

The Bay Area has good weather but is prone to earthquakes and wildfires.

$s_2$

Reward Model ($RM_\phi$)

The   Bay   Area   ...   ... wildfires

Bradley-Terry [1952] paired comparison model

$$J_{RM}(\phi) = -\mathbb{E}_{(s^w, s^l) \sim D} [\log \sigma(RM_\phi(s^w) - RM_\phi(s^l))]$$

"winning" sample    "losing" sample    $s^w$ should score higher than $s^l$

# RLHF
## Limitations

- Human preferences are unreliable!

  - "Reward hacking":
    - https://openai.com/research/emergent-tool-use

  - Chatbots are rewarded to produce responses that seem authoritative and helpful, regardless of truth

  - This can result in making up facts + hallucinations



TECHNOLOGY

## Google shares drop $100 billion after its new AI chatbot makes a mistake

February 9, 2023 · 10:15 AM ET

https://www.npr.org/2023/02/09/1155650909/google-chatbot--error-bard-shares

**Bing AI hallucinates the Super Bowl**



Searching for: superbowl winner

Generating answers for you...

The Super Bowl is the annual American football game that determines the champion of the National Football League (NFL) [1]. The most recent Super Bowl was Super Bowl LVI, which was held on February 6, 2023 at SoFi Stadium in Inglewood, California [2]. The winner of that game was the Philadelphia Eagles, who defeated the Kansas City Chiefs by 31-24. It was the second Super Bowl title for the

The most recent Super Bowl was **Super Bowl LVI**, Stadium in Tampa, Florida. The winner of that game was the Tampa Bay Buccaneers, who defeated **Eagles**, who defeated the **Kansas City Chiefs** by 31-24

Learn more: 1. en.wikipedia.org   2. sportingnews.com   3. cbssports.com

Who won the superbowl?

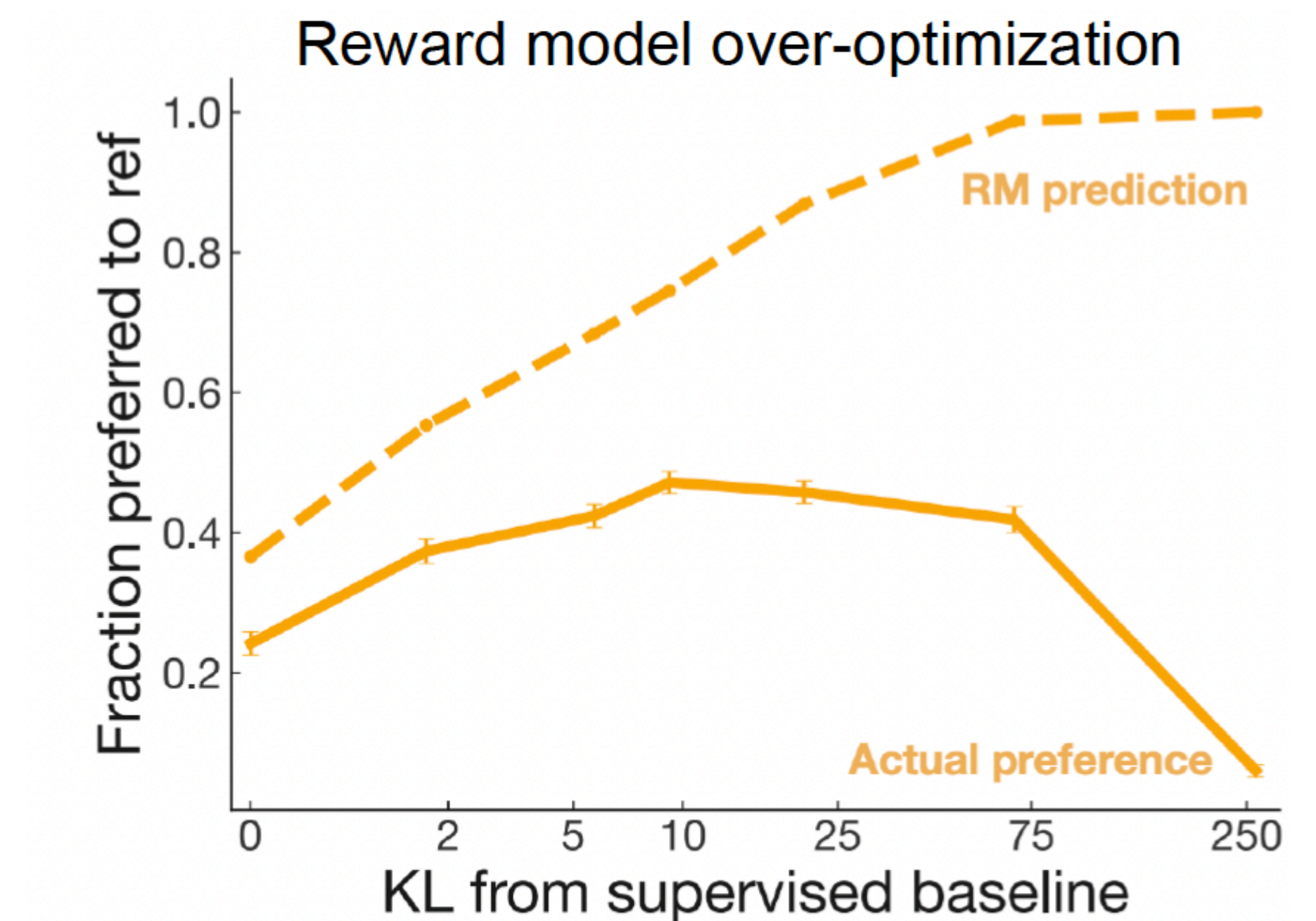https://news.ycombinator.com/item?id=34776508
https://apnews.com/article/kansas-city-chiefs-philadelphia-eagles-technology-science-82bc20f207e3e4cf81abc6a5d9e6b23a

# RLHF
## Limitations

- Human preferences are unreliable!

  - "Reward hacking"

  - Chatbots are rewarded to produce responses that seem authoritative and helpful, regardless of truth

  - This can result in making up facts + hallucinations

- Models of human preferences are even more unreliable!

Reward model over-optimization

$$R(s) = RM_\phi(s) - \beta \log \left( \frac{p_\theta^{RL}(s)}{p^{PT}(s)} \right)$$

# What's next?

- Prompt engineering

**Standard Prompting**

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The answer is 27. ❌

**Chain of Thought Prompting**

Input

Q: Roger has 5 tennis balls. He buys 2 more cans of tennis balls. Each can has 3 tennis balls. How many tennis balls does he have now?

A: Roger started with 5 balls. 2 cans of 3 tennis balls each is 6 tennis balls. 5 + 6 = 11. The answer is 11.

Q: The cafeteria had 23 apples. If they used 20 to make lunch and bought 6 more, how many apples do they have?

Model Output

A: The cafeteria had 23 apples originally. They used 20 to make lunch. So they had 23 - 20 = 3. They bought 6 more apples, so they have 3 + 6 = 9. The answer is 9. ✔

# What's next?

- Prompt engineering

**ANTHROP\C**

## Prompt Engineer and Librarian

APPLY FOR THIS JOB

SAN FRANCISCO, CA / PRODUCT / FULL-TIME / HYBRID

Anthropic's mission is to create reliable, interpretable, and steerable AI systems. We want AI to be safe for our customers and for society as a whole.

Anthropic's AI technology is amongst the most capable and safe in the world. However, large language models are a new type of intelligence, and the art of instructing them in a way that delivers the best results is still in its infancy — it's a hybrid between programming, instructing, and teaching. You will figure out the best methods of prompting our AI to accomplish a wide range of tasks, then document these methods to build up a library of tools and a set of tutorials that allows others to learn prompt engineering or simply find prompts that would be ideal for them.

# What's next?

- Prompt engineering

  - Dark side on prompt engineering

# What's next?

- Prompt engineering

- Can we believe LLM?

  - Fake news

  - Wrong in simple calculation

# What's next?

- Prompt engineering

- Can we believe LLM?

- Specialized LLM

  - AI+healthcare

  - AI+finance

  - AI+science

  - …

# What's next?

- Prompt engineering

- Can we believe LLM?

- Specialized LLM

- Copyright

  - Model&data stealing

  - Generated content's IP

# Q&A

- Interested in doing machine learning research?

- Email: minhaocheng@ust.hk

- Office: Room 2542

THE DEPARTMENT OF
**COMPUTER SCIENCE & ENGINEERING**
計算機科學及工程學系