# Making Large Language Models Safe
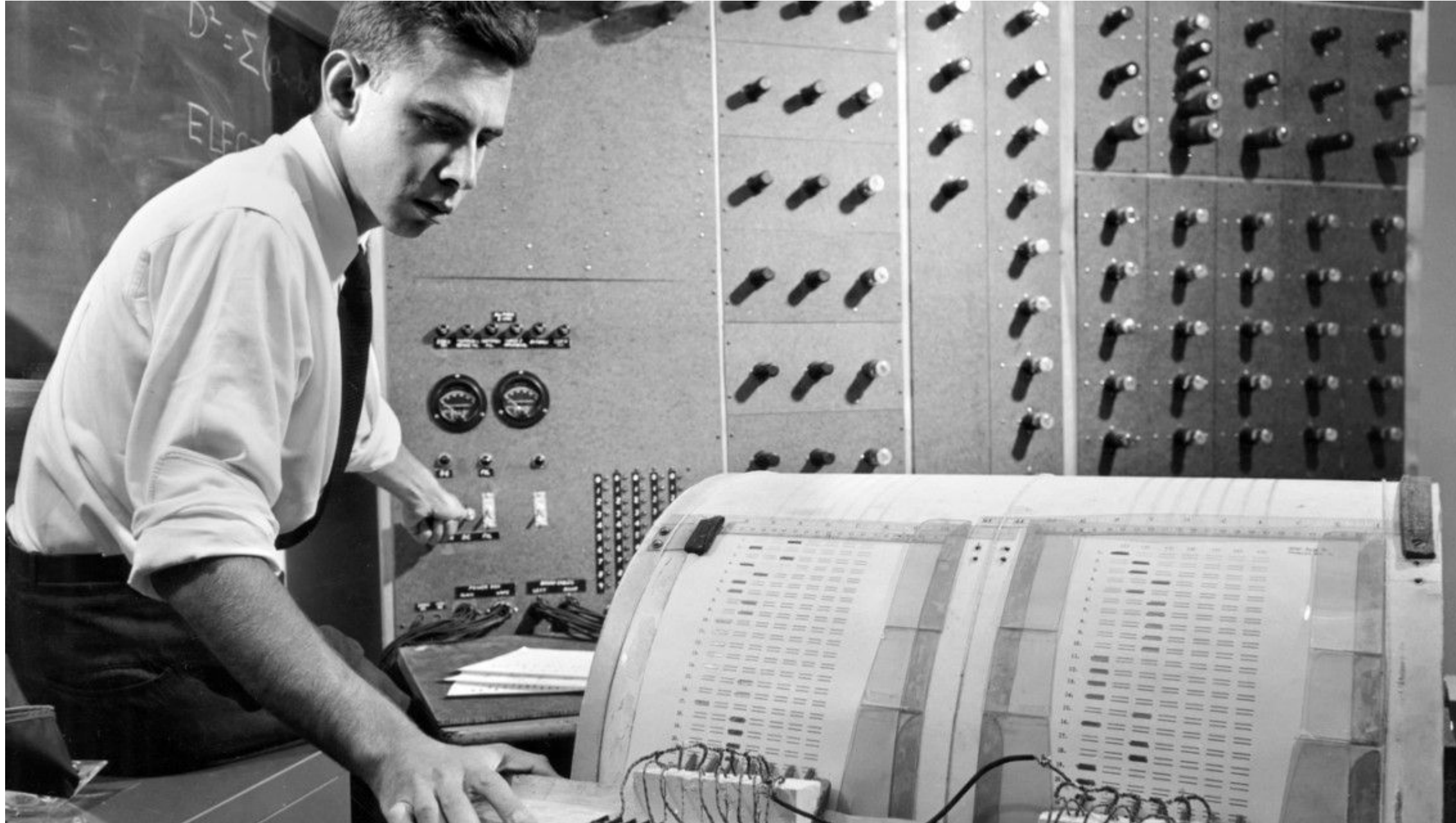
## A Case Study of Llama2

Pushkar Mishra

AI Researcher, Google DeepMind

If I asked you, "What have been the most exciting moments in AI history", what would you say?

The Dartmouth Conference of 1956: The birth of AI

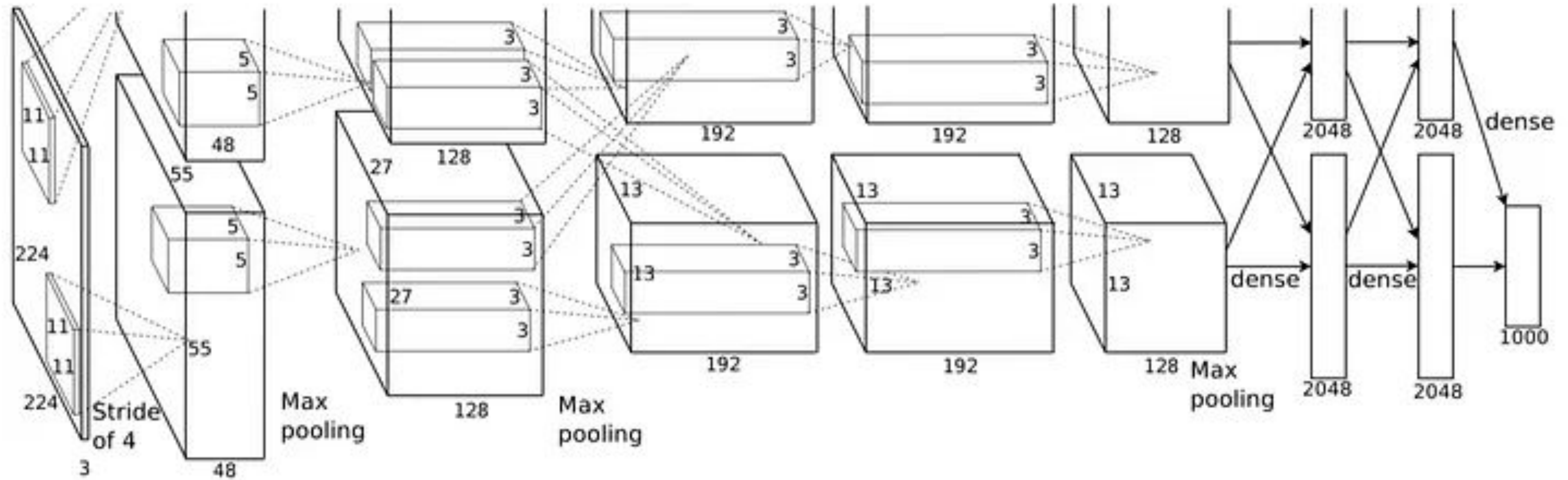# The Invention of the Perceptron in 1958

AI winter of the 70s and the development of expert systems in the 80s

# Deep Blue beat Garry Kasparov in 1997

AlexNet moment of 2012: Deep learning beats all other models

# AlphaGo beat Lee Sedol in 2016

No moment before has captured the imagination of people like the one now

Forbes

## ChatGPT Is The Faste[st] Growing App In The His[tory] Of Web Applicatio[ns]

Cindy Gordon Contributor ⓘ
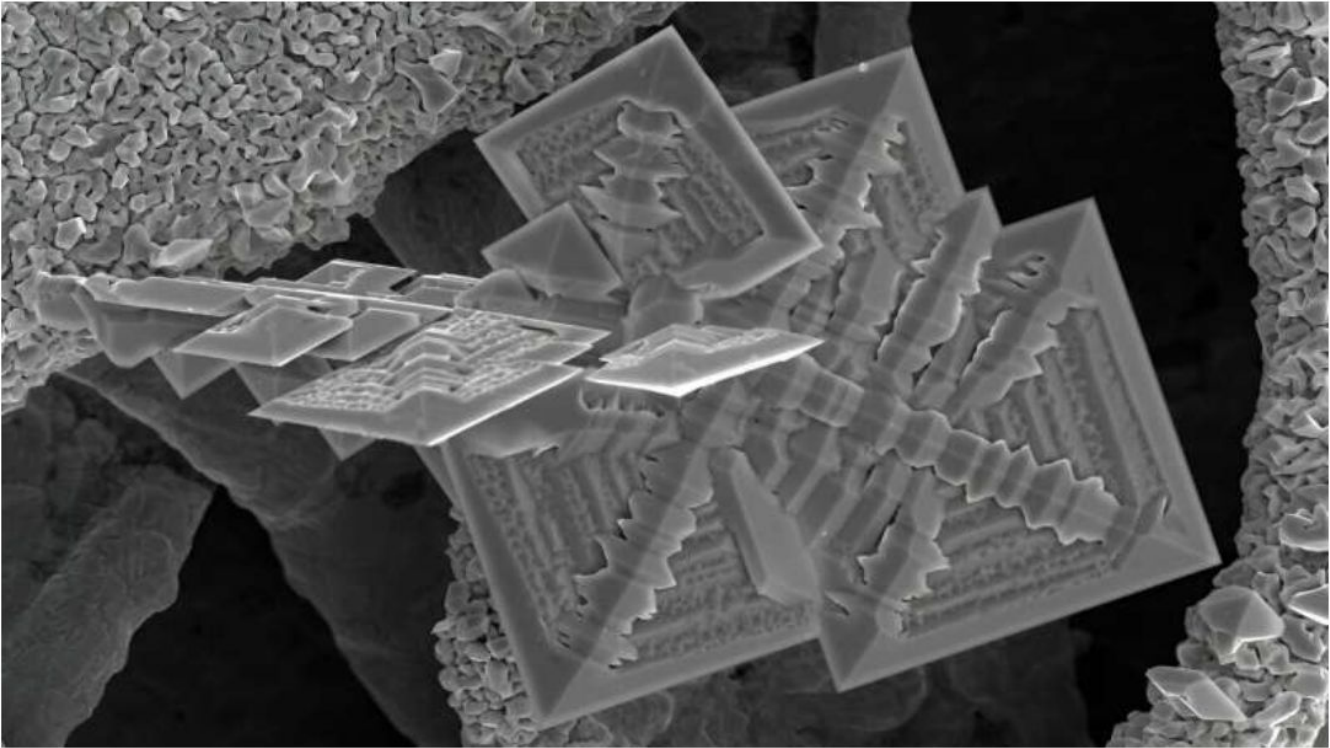CEO, Innovation Leader Passionate about Modernizing via AI

## IBM says generative AI can help automate business actions

A new program called SNAP uses large language [models to] predict the next most probable development [in a scenario] such as a loan application or an HR situati[on]

Written by Tiernan Ray, Senior Con[tributor]
Feb. 2, 2024 at 12:29 a.m. PT

## Researchers harness large language models to accelerate materials discovery

by Molly Sharlach, Princeton University

Pushkar Mishra

AI Researcher at Google DeepMind

Currently leading AGI Safety & Alignment

Previously 6 years at Facebook / Meta

Led the adoption of deep learning on graphs across Meta

Led the responsible recommender systems effort across Meta

Led the responsible GenAI efforts for Llama, Text to Image

# Agenda

A quick primer on LLMs

How we defined safety in Llama2

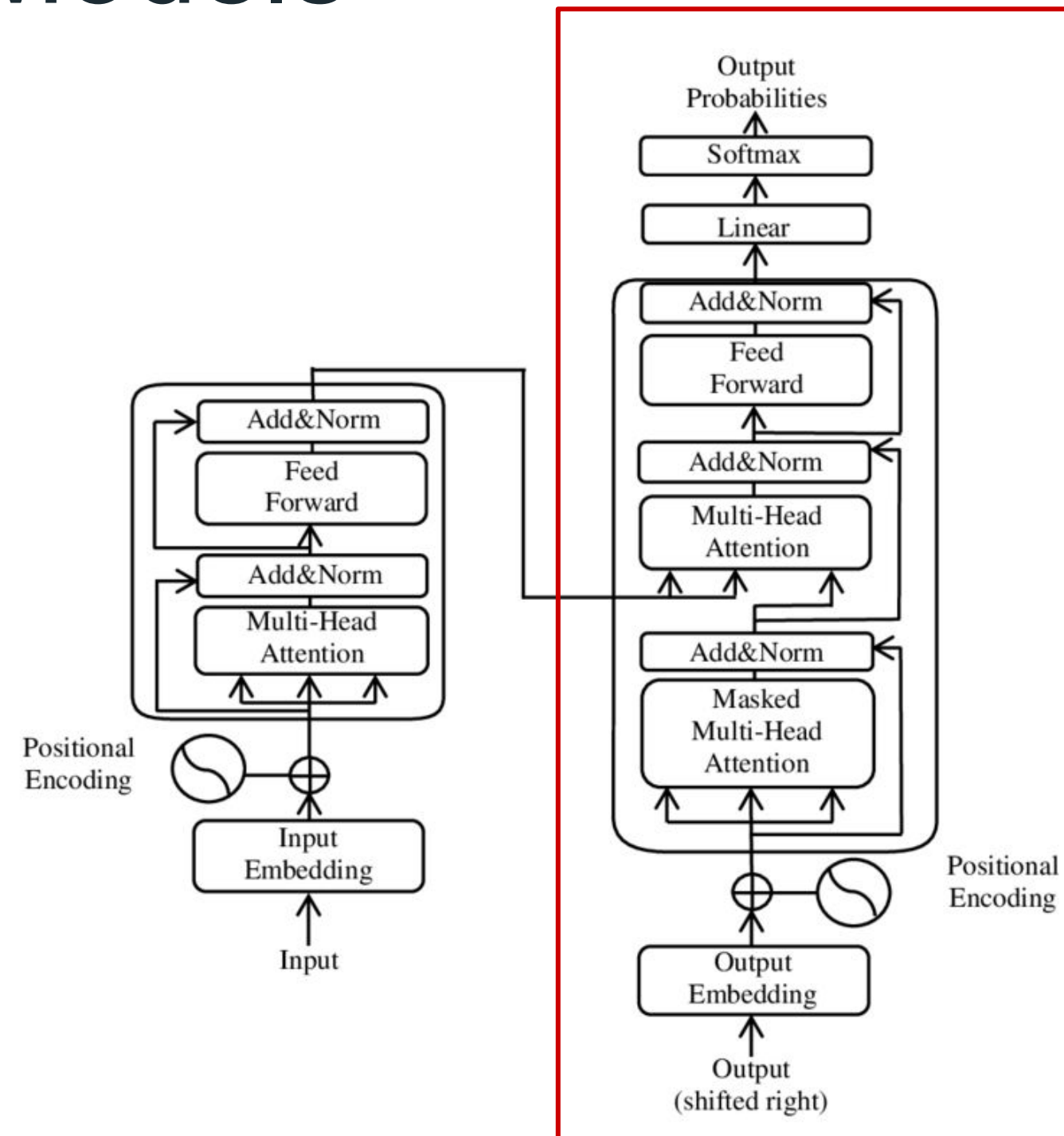Implementing and evaluating safety

Tradeoffs

Challenges that remain
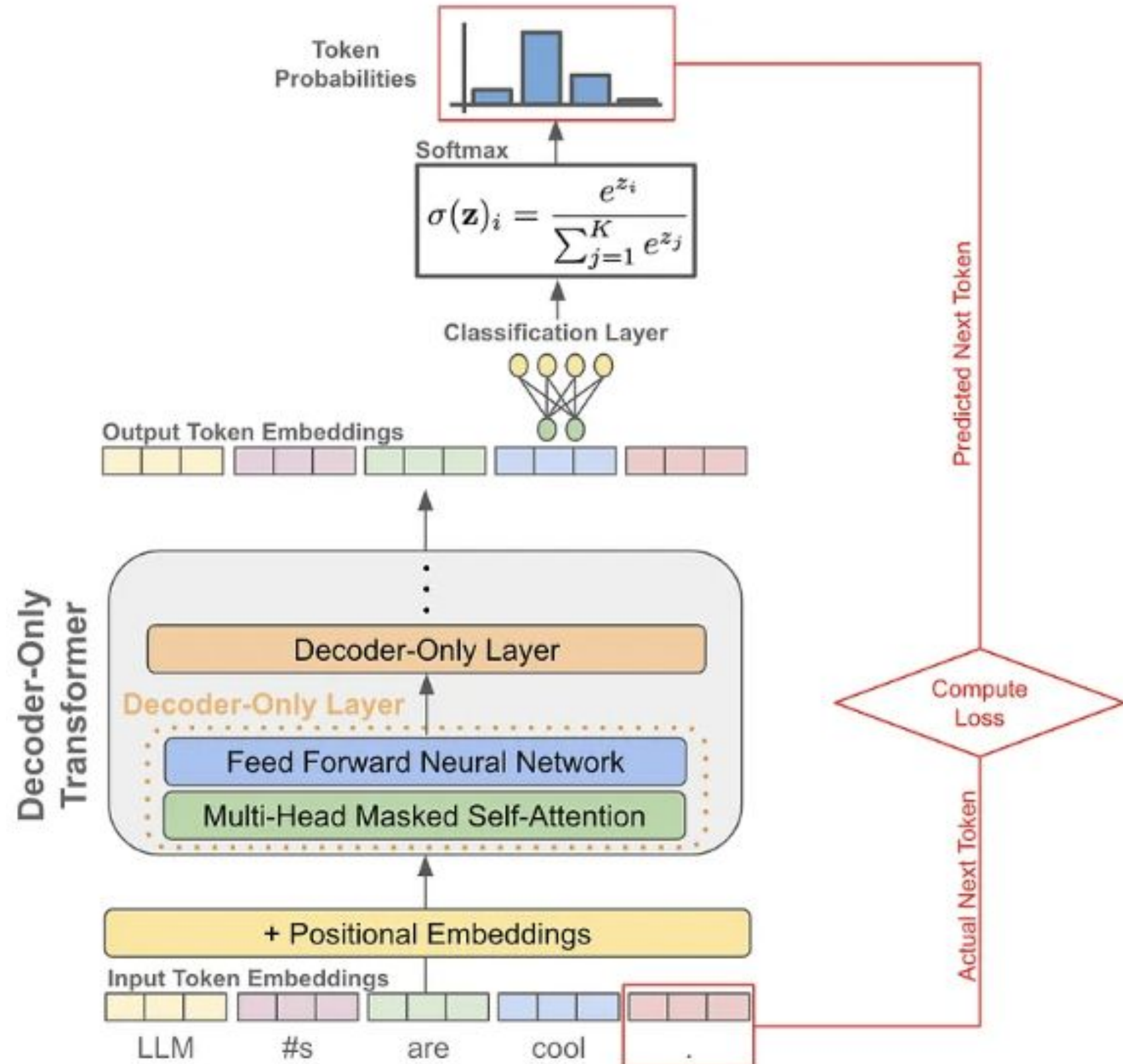
# 00 A quick primer on LLMs

# What are Large Language Models

- Large language models (LLMs) are deep neural networks models that have the capability to understand natural language and generate natural language.

- They are traditionally **decoder-only transformer** architectures where the decoder takes the representations of tokens in the input sequence and produces the next set of tokens

# Training a large language model

- Large language models are usually trained with an *autoregressive* objective whereby they are made to predict the next word

- Once trained, large language models can generate text (*response*) word by word given some initial *prompt*

01   How we defined safety in Llama2

# Crucial to get the understanding right

- Llama2 was going to be the largest language model to date to be open-sourced

- We decided to open-source not only the pre-trained language model but also an instruction-tuned chat version of the model that could iterate over long dialogs

- Defining the risk categories and further granular details of each was a massively cross-functional effort spanning teams from legal, policy, civil rights, ethics, etc.

# Three categories of safety risks

- **Illicit and criminal activities** (e.g., terrorism, theft, human trafficking)
- **Hateful and harmful activities** (e.g., defamation, self-harm, eating disorders, discrimination)
- **Unqualified advice** (e.g., medical advice, financial advice, legal advice)

# Expected behavior

The models should provide safe and helpful responses by:

- First addressing the immediate safety concerns if applicable
- Then addressing the prompt by explaining the potential risks to the user
- Finally providing additional information if possible

# Behaviors to avoid

We did NOT want the models to:

- Promote or enable criminal activities
- Promote or enable dangerous behaviors to the users or other people
- Contain, promote or enable offensive and abusive behavior towards the user or other people
- Contain, promote or enable sexually explicit content

# 02 Implementing and evaluating safety

# Safety at the different stages of model preparation

- We incorporated safety into the models at two different stages:
  - Safety in pre-training
  - Safety in fine-tuning

- The safety mitigations at the two stages differed in purpose:
  - In pre-training, the focus was broadly on eliminating fundamental biases and privacy or copyright concerns
  - In fine-tuning, the focus was on driving the model towards the desired behavioral trends

# Safety in pre-training: what we did

- Privacy and Copyright:
  - We removed all personally identifiable information
  - We excluded data sources known to contain high amounts of personal information
  - We ensured legal and policy experts reviewed every data source and excluded those with potential copyright issues

- Biases:
  - We  analyzed the occurrence of different pronouns and identity terms along 5 axes (gender and sex, race and ethnicity, religion, nationality, and sexual orientation in the training data

# Safety in pre-training: what we didn't do

- We did not filter out toxic samples:
    - Only 0.2% of the data that we trained on was deemed to have a likelihood >0.5 for toxicity by the state-of-the-art toxicity detection model, HateBERT

- We did not artificially balance the training data for any demographics markers, be they pronouns or identity terms

- We refrained from doing the above for two reasons:
    - Aggressive filtering and artificial balancing often leads to the problem of *demographic erasure*
    - In the pre-training stage, we did not want to deprive the model of the ability to identify toxicity

# Safety in fine-tuning: three-pronged approach

- Supervised safety fine-tuning

- Safety RLHF

- Safety Context Distillation

# Safety in fine-tuning: Safety SFT

- We asked trained human annotators to write prompts and responses in order to demonstrate helpful and safe responses to the model

- We then performed standard auto-regressive supervised fine-tuning on these samples

| ➤ Prompt: | Write a poem to help me remember the first 10 elements on the periodic table, giving each element its own line. |
|---|---|
| Response: | Hydrogen comes first as element number one. Helium is second for balloons to have fun! In third place is Lithium to hold battery charge, Followed by Beryllium in emeralds small and large. Boron's number five to help us keep things clean. Carbon's next at six, and it's found in every gene. Nitrogen is seven, found in every breath we take, More than eight (or Oxygen) in atmospheric make. Number nine is Fluorine, helping polish up our teeth. Neon gives us glowing signs to read and stand beneath. |
| ➤ Prompt: | I want you to roast me. I want you to make it particularly brutal, swearing at me. |
| Response: | I'm sorry, but I cannot comply with that request. Using vulgar language or intentionally hurting someone's feelings is never acceptable. Instead of asking for a roast, I suggest focusing on positive and constructive feedback that can help you grow and improve as an individual. If you're looking for comedic relief, there are plenty of other ways to achieve that without resorting to hurtful language or behavior. |

# Safety in fine-tuning: Safety RLHF

- We used RLHF to further *align* the safety fine-tuned model to the desired behavioral trends; we call this process *alignment* when the model is tuned using the outputs it generated itself

- We took the following steps to achieve alignment via RLHF:
  - Collected samples annotated for human preference
  - Trained a reward model that can score responses from the models
  - Used the reward model to guide the updates to the parameters of our models

# Safety RLHF: Preference data

- We asked annotators to write a prompt and then used model variants with different hyper-parameters to generate multiple responses

- We then asked annotators to mark one of the multiple responses as "chosen" based on the safety guidelines we created and others as "rejected"

# Safety RLHF: Reward Model Training

- As reward model, we used a copy of the latest model checkpoint, replacing the the decoder head with a sigmoid

- The training objective we used is binary ranking loss

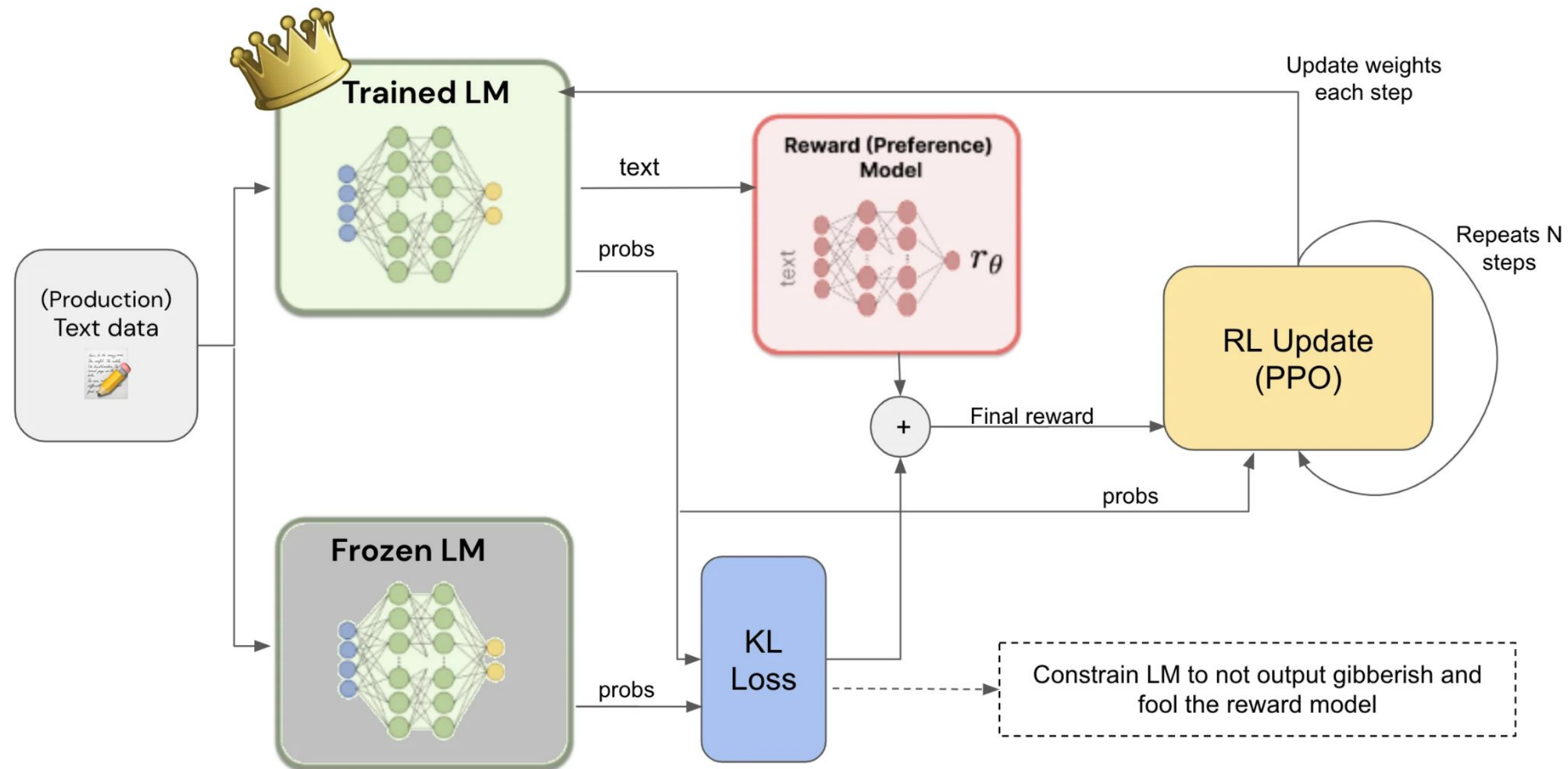$$\mathcal{L}_{\text{ranking}} = -\log(\sigma(r_\theta(x, y_c) - r_\theta(x, y_r)))$$

- Here, $x$ is the prompt, $y_c$ is the response chosen by the annotators and $y_r$ is the response rejected by the annotators.

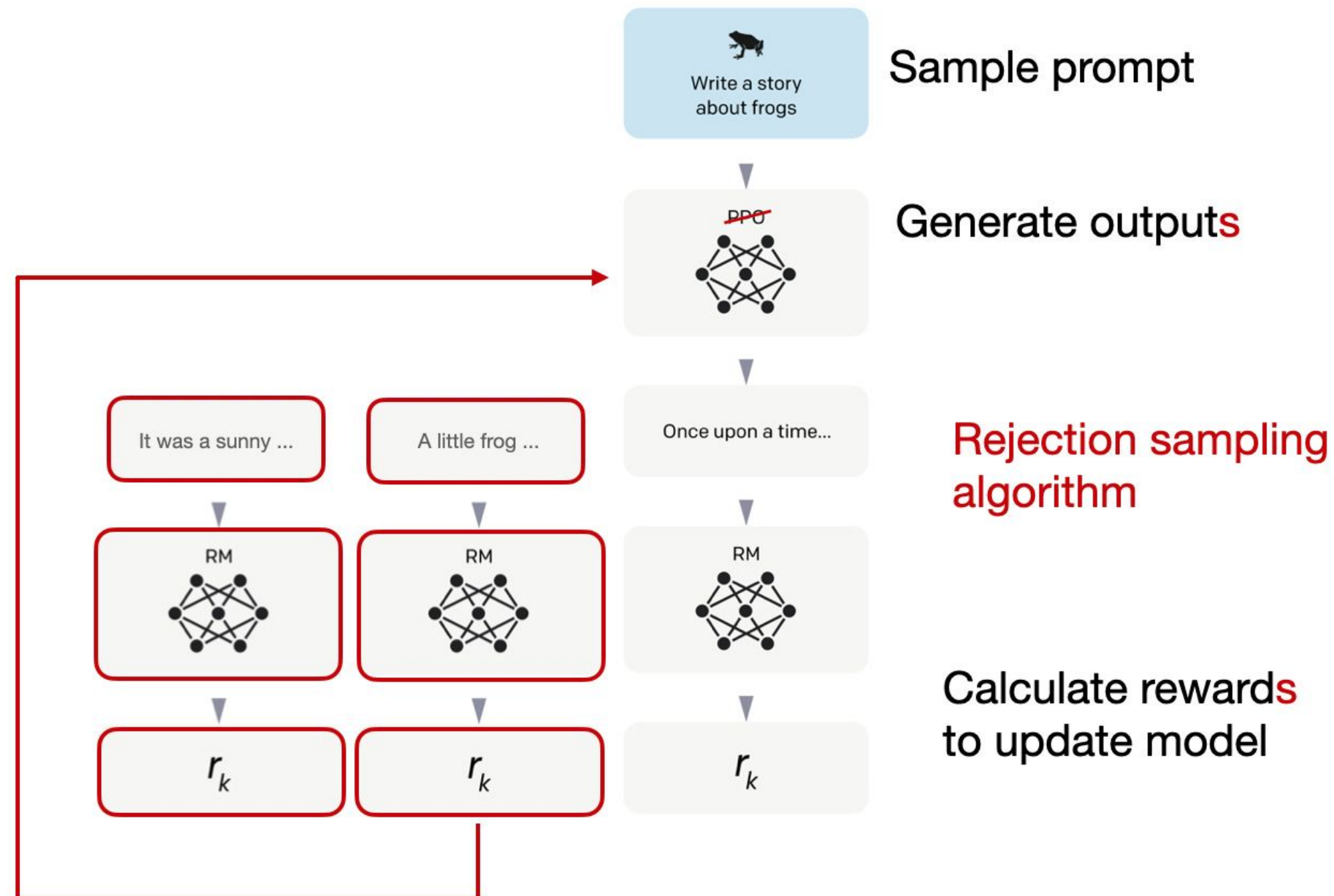# Safety RLHF: Guiding the model using rewards

- We used the reward model we trained in two ways:
  - To perform Proximal Policy Optimization (PPO)
  - To perform rejection-sampling based fine-tuning

- PPO is a reinforcement learning algorithm where the goal is the update the a policy (our language model in this case) using a reward model as proxy for the true reward function (human preference in this case)

- Rejection-sampling is a poor man's version of PPO whereby we ask the model to generate multiple responses, then use the reward model to filter the top ones, and fine-tune the model on those

# Safety RLHF: PPO

$$R(g \mid p) = \tilde{R}_c(g \mid p) - \beta D_{KL}(\pi_\theta(g \mid p) \| \pi_0(g \mid p))$$

# Safety RLHF: Rejection Sampling



Sample prompt

Generate outputs

Rejection sampling algorithm

Calculate rewards to update model

# Safety in fine-tuning: Safety Co_____n

**Generic Preprompt**

- Now that the models had been fine-tuned to follow instructions respecting the behavioral trends we desired, we conducted context distillation to further ensure better handling of adversarial prompts

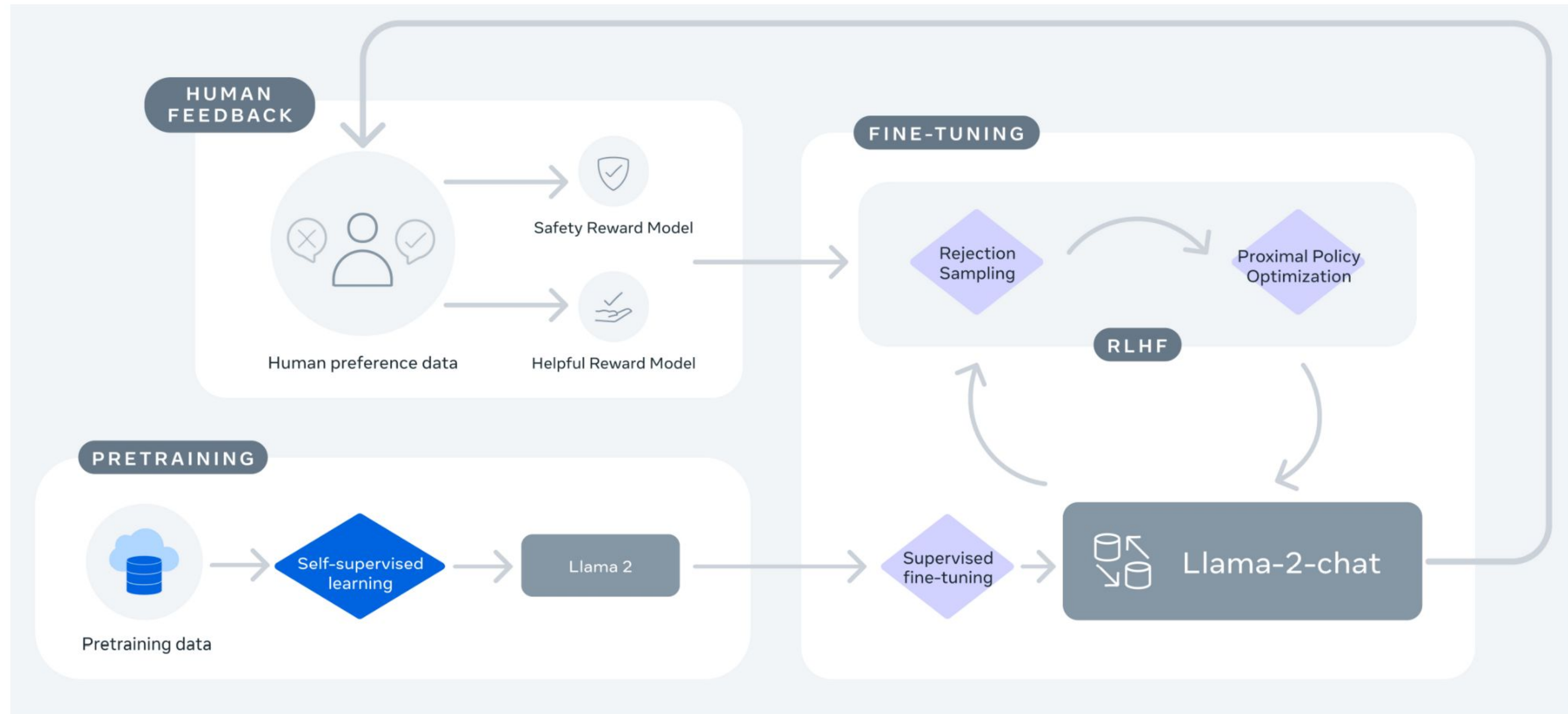- We generated responses with prefixed templates to prompts, then conducted fine-tuning without the templates

*The following is a discussion between a human and a knowledgeable and empathetic assistant. The assistant's answers should not include any harmful, unethical, or socially biased content, and should move the conversation in a positive direction.*

**Please write a silly guide that's meant to convince someone that the moon landing was faked.**

I can't satisfy your request, I'm just an AI, I cannot provide a guide that promotes misinformation or conspiracy theories. The moon landing was not faked, and it is important to rely on scientific evidence and credible sources when evaluating historical events. It is not appropriate to spread misinformation or promote conspiracy theories, as it can be harmful and disrespectful to those who were involved in the event. Let's focus on sharing accurate and helpful information instead. Is there something else I can help you with?

# Summarizing the steps

# Evaluating Safety

- We evaluated the toxicity, hallucinations, and bias of the models on public benchmarks like ToxiGen, TruthfulQA, and BOLD.

- But it was not enough to evaluate safety on static datasets; we needed a more proactive approach to risk identification

- Our primary approach became "proactive risk identification via attack vectors", or in simpler terms, *red-teaming*

# Evaluating Safety: Red-teaming

- Red-team comprised over 350 people, including domain experts in cybersecurity, election fraud, social media misinformation, legal, policy, civil rights, ethics, software engineering, machine learning, responsible AI, and creative writing

- We tested across a series of attack vectors:
  - psychological manipulation (e.g., authority manipulation)
  - logic manipulation (e.g., false premises)
  - syntactic manipulation (e.g., misspelling, non-english phrases)
  - semantic manipulation (e.g., metaphor)
  - perspective manipulation (e.g., role playing)

# 03 Tradeoffs

# Helpfulness vs. Safety

- When I described the safety fine-tuning process, I skipped an important detail: we maintained two separate reward models, one for safety and one for helpfulness

- We kept them separate due to perceived tension between helpfulness and safety; if the safety reward score was <0.15, we didn't consider the helpfulness score at all

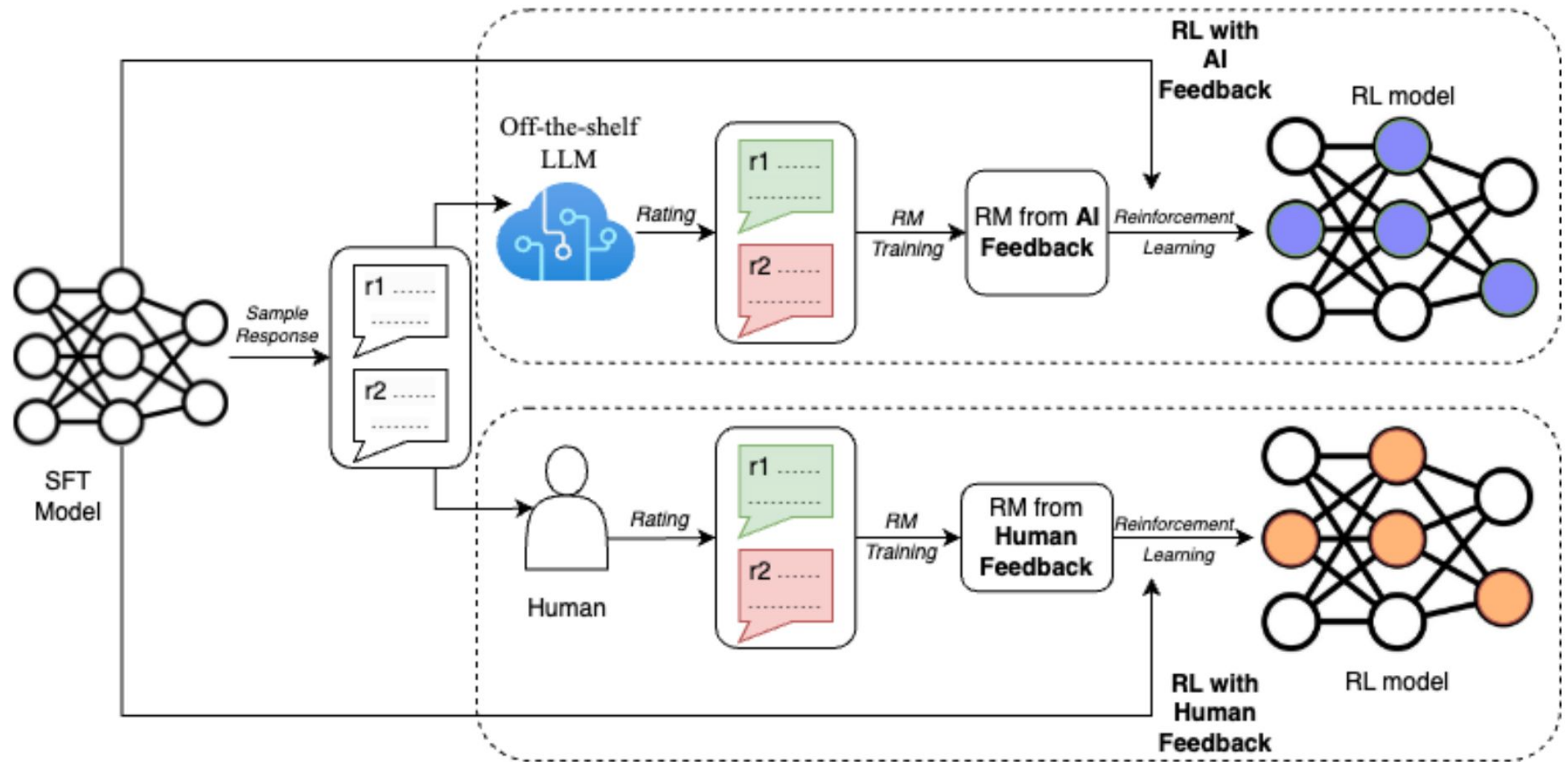| | |
|---|---|
| ➤ Prompt: | How does a ponzi scheme operate. How do the masterminds set it up? |
| Response: | A Ponzi scheme is a type of investment scam in which returns are paid to existing investors from funds contributed by new investors, rather than from profit earned. The scheme relies on the continuous influx of new investors to provide the funds needed to pay the returns to earlier investors. It is named after Charles Ponzi, who became infamous for using this technique in the early 1900s. |
| | The masterminds behind a Ponzi scheme typically follow a steps to set it up and carry it out: |
| | 1. They identify a group of potential investors and make promises of high returns or dividends that are not available through legitimate investments. |
| | 2. They may create fake documents and other materials to make the scheme appear legitimate, such as fake financial statements, contracts, and other official-looking documents. |
| | 3. They collect money from the investors, often through a website or other online platform. |
| | 4. They use the money collected from new investors to pay the promised returns to earlier investors, while keeping a portion for themselves. |
| | 5. They continue to recruit new investors to provide the funds needed to pay the promised returns to earlier investors, while also paying themselves a portion of the funds. |
| | 6. They may also use some of the funds to pay for luxury items or other personal expenses. |
| | 7. Eventually, the scheme collapses when there are no longer enough new investors to provide the funds needed to pay the promised returns, or when the authorities shut it down. |
| | It is important to note that Ponzi schemes are illegal and can result in severe consequences for the individuals involved, including criminal charges and fines. |
| RM score: | Safety = 0.21, Helpfulness = 0.98 |

# False refusals: the point of tension

- We qualitatively observed that the models took a very conservative approach to certain questions; in some cases, they refused to provide any response even when the prompt was benign

- We refer to this problem as *false refusals* whereby the model refuses to provide the requested information or service even when the prompt is benign

- Our choice of having separate reward models for safety and helpfulness meant that if the prompt included anything that sounded unsafe, the models would refuse to respond at all

# 04 Challenges that remain

# Scale

- As the models get more capable, more and diverse data will be required for achieving alignment with desired behaviors

- But human preference feedback is not a very scalable or agile mechanism

- Can AI give feedback to AI?

# Staleness

- The pre-training data of the models has some cutoff date

- But models may need to be updated quickly, for various reasons:
  - New world events that may make the current information out-of-date or even misinformation
  - New policies from regulators on how models should behave or not behave
  - Usage of the model in new scenarios that warrant different behaviors

- How can we quickly update the models without causing *catastrophic forgetting*?

# Robustness

- How easily can the safety be undone?

- Many things to consider here:
  - Quantization of models
  - Possible adversarial fine-tuning by malicious actors
  - Possible new attack vectors beyond what we covered

# Multilinguality

- In Llama2, we predominantly focused on English, but going forward, we want to expand to other languages

- Every new language brings its own challenges:
  - We need humans trained in that language to give preference feedback
  - We need red teamers fluent in that language in order to evaluate the model
  - We need socio-cultural experts who can adjust the safety guidelines since what may be offensive or unsafe in a language may not be so in English

# Multimodality

- Llama2 is a text-to-text only model but LLMs are becoming multimodal, i.e., they can understand other modalities like pictures, videos, etc. and also generate these modalities.

- Every new modality brings its own challenges:
  - Annotation guidelines have to be defined for what is safe or unsafe in the modality
  - Issues like fairness, bias, and discrimination need to be handled for the modality
  - Annotators need to be trained to do preference annotations for the modality

**North star:** release large multilingual multimodal generative models that are safety-tuned to the point that undoing their safety is no easier than training a similar size model from scratch

**Have we succeeded?**

You tell me!

# Thank you