

# Reinforcement Learning:

## *How Models Learn to Reason*

Abbie Petulante  
DSI Postdoctoral Fellow

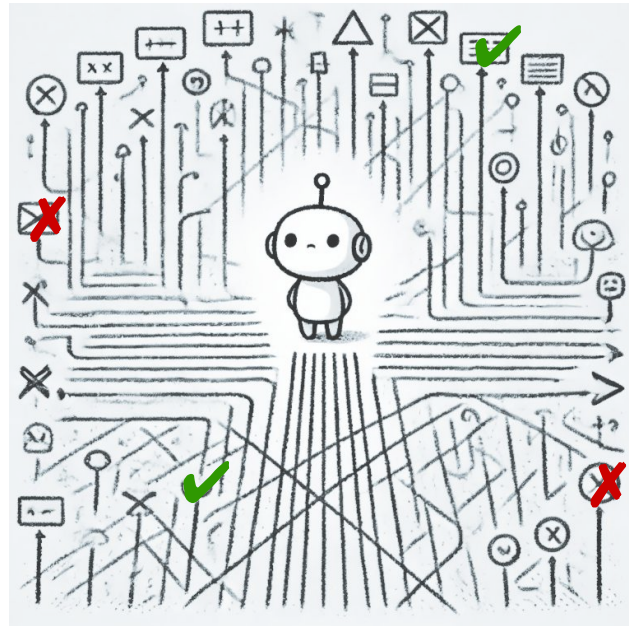


Data Science Institute

**Discovery through data.**

# Reinforcement Learning (RL)

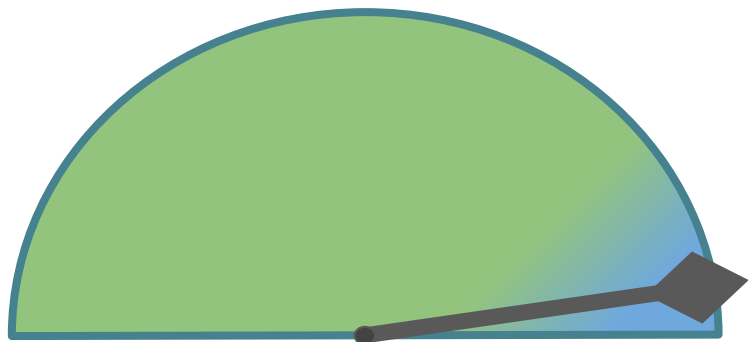
- A **trial-and-error** learning process
  - an **agent**
  - interacts with an **environment**,
  - takes **actions**,
  - receives **rewards** for positive actions
- used when we don't have **direct supervision** but *can* define a goal



# Why in focus now?

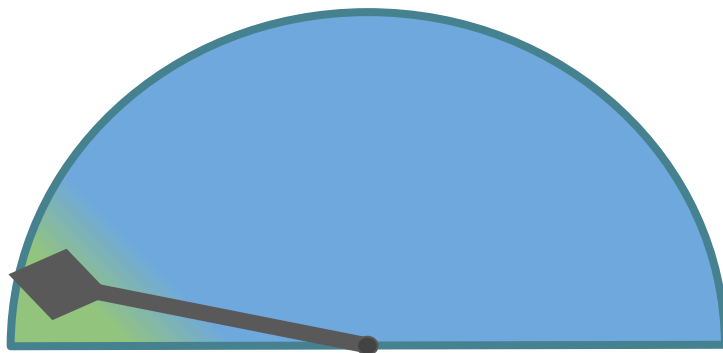
- **Because of reasoning models!**

Supervised  
Learning



*Standard (SFT + RLHF at end)*

Reinforcement  
Learning



*Typical primary RL*

# Reasoning Responses

- *<think>*
    - Here is some reasoning through
    - Thoughts, step-by-steps
  - *</think>*
- 
- *<answer>*
    - Here is my final answer
    - Usually a summary of what was thought through
  - *</answer>*

# Why RL for Reasoning

- *<think>*

- Here is some reasoning through
- Thoughts, step-by-steps

Hard to quantify as  
good, right, correct

- *</think>*

- *<answer>*

- Here is my final answer
- Usually a summary of what was thought through

Might be easier to  
quantify, right vs wrong

- *</answer>*

# An RL Crash Course

# The Reward Function

- Rewards are a way of evaluating **goodness**
  - “**Goodness**”  $\neq$  “**correctness**” *necessarily*
- **You might reward:**
  - Efficiency  $\rightarrow$  reach the answer in the few words
  - Coherence & Fluency  $\rightarrow$  for long dialogues
  - Preference & Style  $\rightarrow$  kindness, helpfulness, etc

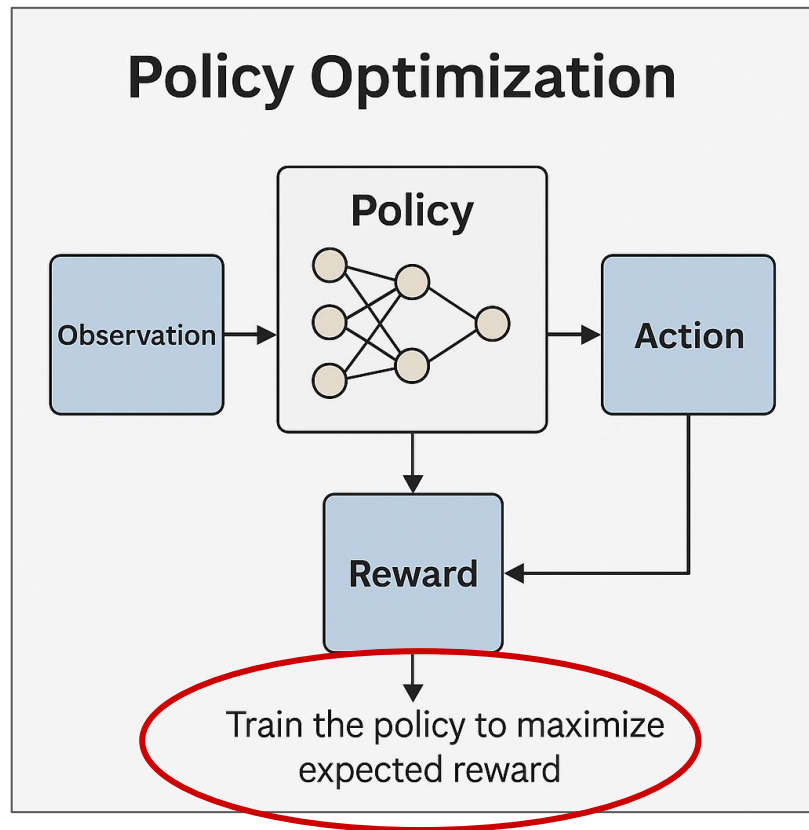
# Policy Optimization

- A “policy” is a strategy to decide what action to take in a given situation
  - For LLMs, *the model* - it's current weights and state!
- **Policy Optimization** refers to *how* you update the model to *get more rewards* next time.



# Policy Optimization

- Policy optimization strategies ask:
  - How do I determine what responses are better?
  - How do I change my model to make more of those?
- Differ in how model updates are chosen and implemented



# GRPO

- **GRPO** (Group Relative Policy Optimization) is a policy optimization strategy designed for reasoning

## Key Features:

- Compares groups of possible responses
- Encourages the model to prefer responses that are better than others *in the group*

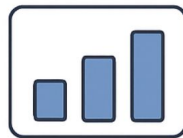
# GRPO for Reasoning

- By focusing on relative performance within groups, GRPO helps models develop better reasoning strategies
  - Comparison of different “approaches” to a problem
  - Aligns with how humans favor responses (seeing many options)

For each prompt, generate a group of responses



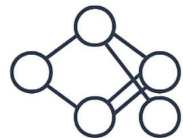
Evaluate the relative quality of the responses



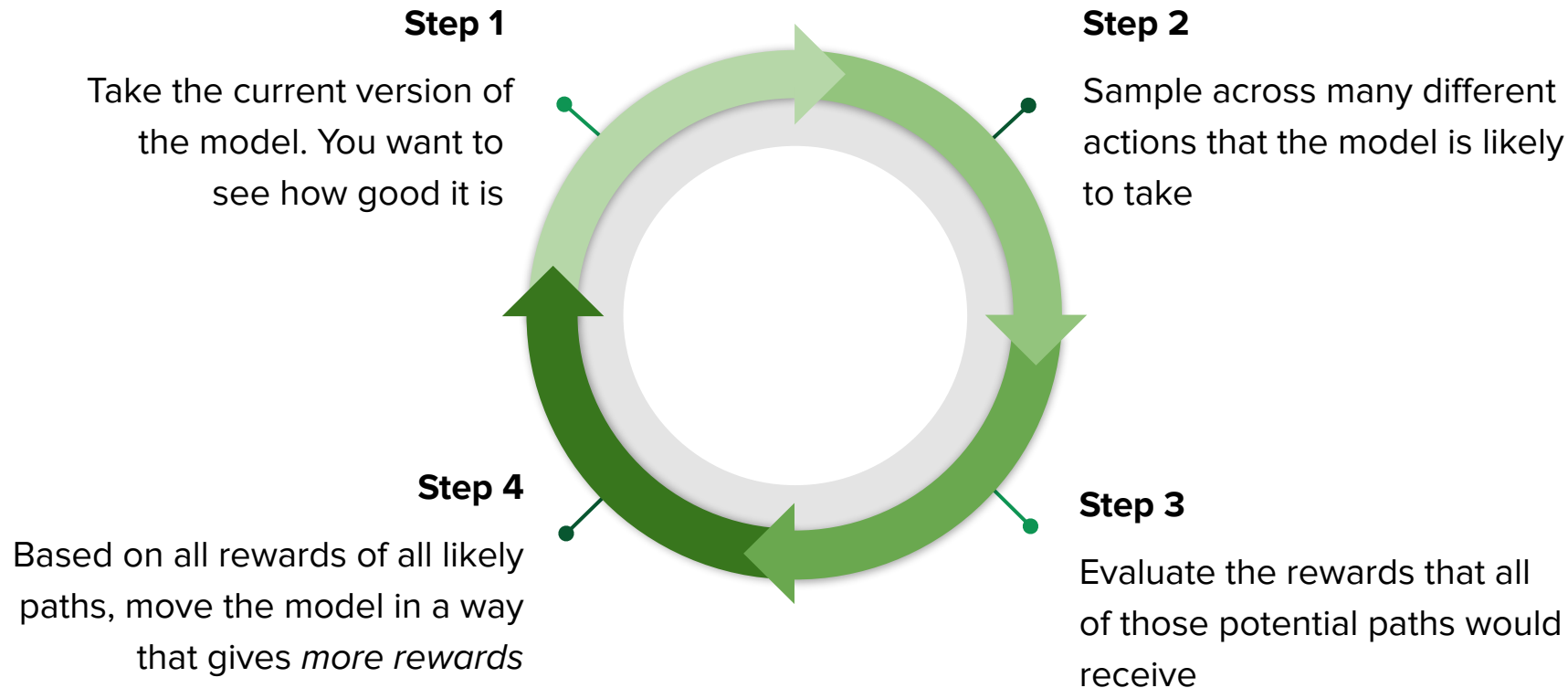
Assign higher rewards to better responses



Update the policy to improve response quality



# The RL Training Loop



Imagine You're an LLM...

# “Traditional” LLM Learning

- The model is trained to predict what comes next:

“I want to go to \_\_\_\_\_”

# “Traditional” LLM Learning

- The model is trained to predict what comes next:

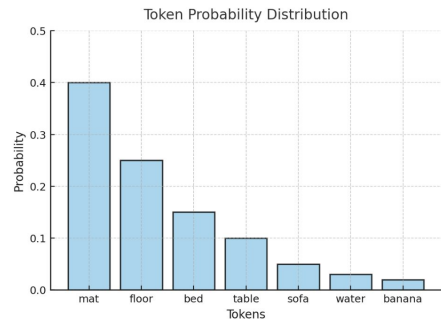
“I want to go to \_\_\_\_\_”

*bed?*

*sleep?*

*the?*

*see?*



Probably, pretty high probabilities for all

# RL

- The model is trained to predict what comes next:

“I want to go to \_\_\_\_\_”

*bed?*

All of these choices  
are *reasonable*...but  
which are *better*

*sleep?*

*the?*

*see?*

**Which do people  
prefer?**



# RL

- The model is trained to predict what comes next:

“I want to go to \_\_\_\_\_”

~~bed.~~

~~sleep.~~

*The movies later today because Ghostbusters just came out and I'm so excited to see it.*

*See the cherry blossoms bloom, I've always thought that must be the most marvelous sight.*

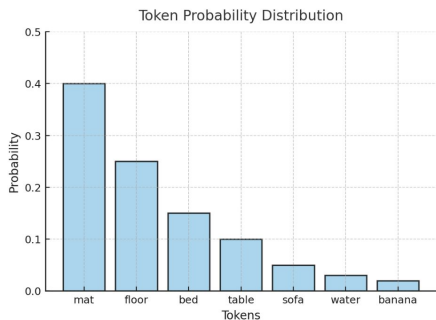
# RL

- The model is trained to predict what comes next:

“I want to go to \_\_\_\_\_”

~~bed.~~

~~sleep.~~



Not wrong at any step, just not optimized for our preference of long answers

# RL for Reasoning

**Next-token self-supervised learning:**

$$\begin{array}{r} 2x + 3 = 9 \\ x + \underline{\hspace{1cm}} \end{array}$$

# RL for Reasoning

## Next-token self-supervised learning:

$$\begin{array}{rcl} 2x + 3 & = & 9 \\ x + \underline{\hspace{1cm}} & & \end{array}$$

You have access to the original equation

Fill onward just based on memorization, not on solving and working backwards

Uses what you've memorized, what typically follows  $x=?$

Evaluate on whether the blank is filled exactly how was expected

# RL for Reasoning

**Next-token self-supervised learning:**

$$\begin{array}{r} 2x + 3 = 9 \\ x + \underline{\hspace{1cm}} \end{array}$$

**Reinforcement Learning:**

$$\begin{array}{r} 2x + 3 = 9 \\ x + 3 = 9 \\ x = 6 \end{array}$$

$$\begin{array}{r} 2x + 3 = 9 \\ 2x = 6 \\ x = 3 \end{array}$$

# RL for Reasoning

## Next-token self-supervised learning:

$$\begin{array}{r} 2x + 3 = 9 \\ x + \underline{\hspace{1cm}} \end{array}$$

You can see both potential responses in their entirety

## Reinforcement Learning:

$$\begin{array}{r} 2x + 3 = 9 \\ x + 3 = 9 \\ x = 6 \end{array}$$



$$\begin{array}{r} 2x + 3 = 9 \\ 2x = 6 \\ x = 3 \end{array}$$



And evaluate both separately and on multiple criteria

Are they correct? Are they logical? Are they concise?

# RL for Reasoning

**Next-token self-supervised learning:**

$$\begin{array}{r} 2x + 3 = 9 \\ x + \underline{\hspace{1cm}} \end{array}$$

**Reinforcement Learning:**

$$\begin{array}{r} 2x + 3 = 9 \\ x + 3 \text{ ~~X~~ } = 9 \\ x = 6 \end{array}$$

$$\begin{array}{r} 2x + 3 = 9 \\ \checkmark 2x = 6 \\ x = 3 \end{array}$$

Per-token rewards give us more info and help avoid bad paths in the future

# An RL Dataset

Unlike vast web-scale datasets for next-token prediction, RL-trained models:

- Generate their own outputs and **learn by self-evaluating**
- Use a ***reward model*** to score the outputs

The reward model is often optimizing, from the start:

- Correctness, efficiency, logical coherence, step-by-step accuracy



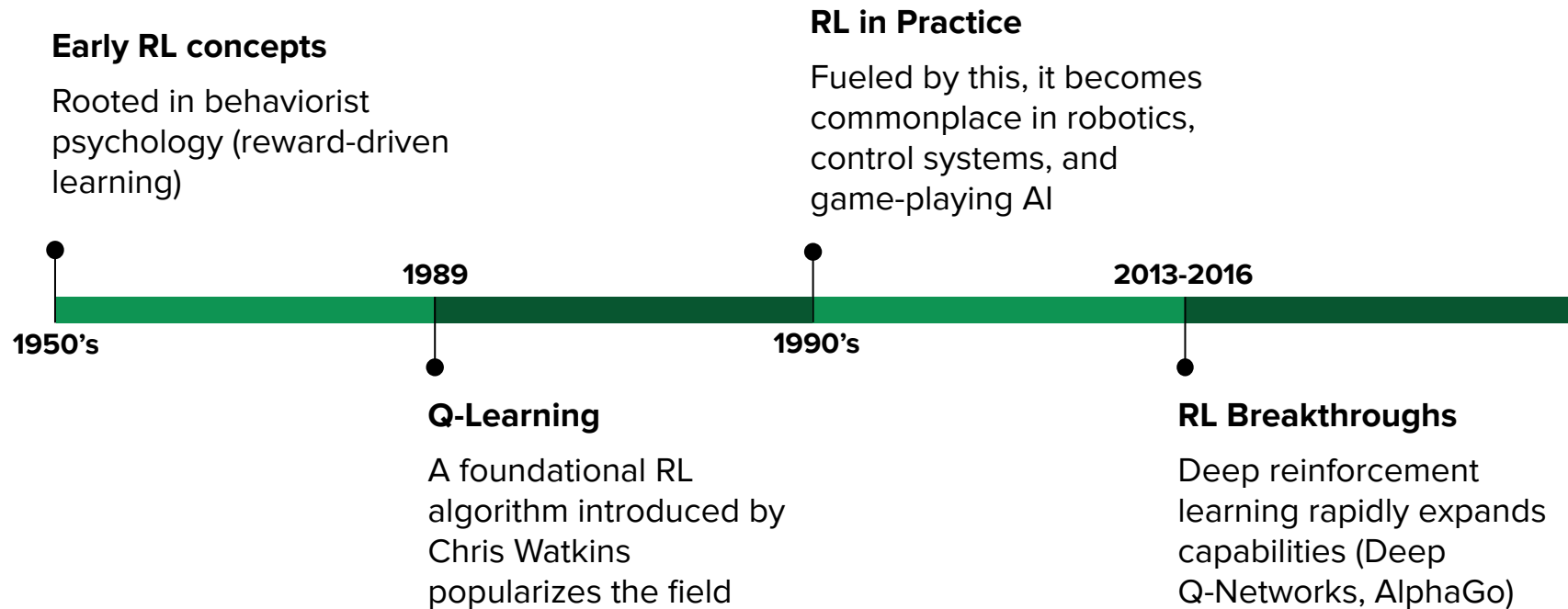
# Traditional vs. RL

	Traditional Supervised Learning	RL for Reasoning
Learning Task	Next-token prediction	Optimizing reasoning paths
Optimization	Minimize token loss (cross-entropy)	Maximize total reward over reasoning steps
Training Data	Fixed datasets of human-written text	Self-generated reasoning sequences
Weaknesses	Struggles with multi-step logic	Requires complex reward models
Key Benefit	Fluent, general-purpose text generation	Structured, efficient reasoning

# The Drawbacks

- **Challenges with using primary RL**
  - **Exploration problems** – Models must try inefficient reasoning paths before finding the best.
  - **Computational cost** – must simulate many responses.
  - **Sparse rewards** can delay feedback (dense rewards solve)
  - **Reward hacking** – poorly designed rewards have exploitable shortcuts

# A Brief History of RL



# A Brief History of RL for LLM's

## RLHF Revolutionizes Alignment

OpenAI releases InstructGPT, which uses RL *from human feedback*, significantly improving response helpfulness and safety

## RL Use Expands

RL is no longer just about safety but actively **improving logical and task-specific performance.**

2020

### OpenAI Releases GPT3

It's impressive, but **lacks fine-grained control and alignment**, leading to unpredictable responses.

Jan 2022

Nov 2022

### OpenAI Releases ChatGPT

Aligned with RLHF, the model is broadly useful and *changes the world*

2023

# Why in focus now?

- Until recently, RL was an “accessory” to LLM training
  - Supervised learning teaches language
  - Instruction fine-tuning teaches directions
  - RLHF teaches the model to “be nice”
- RL was mainly used for **alignment, user satisfaction, safety** as RLHF
- Now, used as a first-step in training, to teach **reasoning, abstract thought**

# 2024: A Paradigm Shift

- **DeepSeek R1 (Late 2024):** *trained* using **primarily RL** rather than supervised fine-tuning.

## Supervised Learning

**Core Mechanism:** Next-token prediction → memorizing patterns of massive datasets

**Lacks global planning**—each token is predicted **locally**, without assessing the long-term impact of decisions.

## Reinforcement Learning

**Core Mechanism:** Reward-guided optimization → learning from outcomes

**Evaluates sequences of actions holistically**, refining its reasoning rather than just token probabilities.