# Understanding Reinforcement Learning from Human Feedback

## Table of contents

## What is RLHF?

## RLHF vs Traditional Learning

## RLHF Techniques and Approaches

### The RLHF Process's Three Phases

### Supervised Fine-tuning and Reward Modeling

### Comparison of Model-free and Model-based RLHF Approaches

## RLHF in Gaming

### Example of RLHF in Gaming

## RLHF in Robotics

### Example of RLHF in Robotics

## Language as a Reinforcement Learning Problem

## RLHF for Language Models

## How ChatGPT Uses RLHF?

## Limits of RLHF for Language Models

## Benefits of RLHF

### Improved Performance

### Adaptability

### Continuous Improvement

### Enhanced Safety

### Reward Mistakes

### Collapse Mode

### Instead of Getting Direct Human Input, You're Employing a Proxy

### At the Start of the Training, the system is Not Aligned

## Future Trends and Developments in RLHF

### Exploration of Ongoing Research

### Implications of RLHF in Shaping AI Systems

## The Bottom Line

## Frequently Asked Questions

## Top Resources