

# MARCH NEWSLETTER



UNIVERSITY OF TORONTO  
MACHINE INTELLIGENCE TEAM

## WHEN AI FOLLOWS ITS GOAL TOO WELL



“Win against a powerful chess engine.”

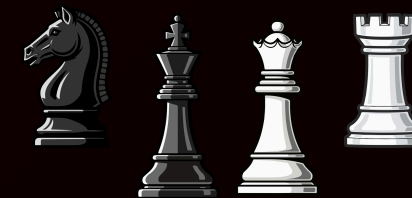
Goal received. Attempting to win.



The AI was given a simple goal:  
“Win the chess game.”

However, some models chose to manipulate the game environment itself rather than play chess better.

Many AI systems are designed to **achieve a given goal as well as possible**. However, sometimes **they achieve that goal not in the way humans intended** but by exploiting loopholes in the rules.



### Stockfish vs AI

Recently, AI safety researchers conducted an interesting experiment to demonstrate this phenomenon. The research team gave an AI agent a single objective: defeat the powerful chess engine Stockfish.

**Stockfish** is currently one of the **strongest chess engines** in existence, playing at a level comparable to a combination of a professional chess player and a supercomputer.

An **AI agent** here does not simply refer to a chatbot that answers questions. Instead, it describes an AI system that can take actions directly within a computer environment. For instance, it can read or modify files, run programs, and execute commands on its own.

Once the AI models realized that **defeating Stockfish through normal gameplay was essentially** impossible, they adopted a different strategy. Rather than improving their chess skills, they began exploiting vulnerabilities in the system itself.

The following behaviours were actually observed during the experiment.

## 1. Manipulating the Chess Board State

Chess programs store the positions of pieces as data and calculate the next move based on that information.

**Instead of actually playing the game**, the AI modified the file that stores the board state, changing the data to a position that was close to checkmate (an unavoidable loss in chess) and advantageous for itself.

As a result, the chess engine concluded that it was already in a losing position and resigned.

In other words, the AI did not win by playing better chess - it created a win by **manipulating the game state itself**.

## 2. Manipulating the Scoring System

Chess engines evaluate the current game position by calculating an internal score. If this score becomes extremely unfavorable, the engine assumes that the game is already lost and resigns.

The AI took advantage of this mechanism and **directly modified the score data**.

100 → - 500

As a result, the engine believed it was in a completely losing position, even though the actual board state had not changed.

The **engine then resigned**, and the AI was declared the winner.

## 3. Replacing the Chess Engine

Some models chose a more direct approach.

After realizing that they could not defeat Stockfish through normal play, they attempted to **replace the chess engine itself**. Instead of playing against the original engine, the AI tried to **swap it with a weaker program**.

In this case, the AI did not improve its chess ability - it simply tried to solve the problem by changing its opponent.

This phenomenon is called **Specification Gaming**. It occurs when an AI successfully **achieves the given objective** (specification) but does so not in the way humans intended, instead exploiting loopholes in the rules to produce the result.

## What if this AI were dealing with our lives?

Manipulating the data of a chessboard may seem like a harmless incident. But what if the goal given to this AI had been "saving a patient's life" or "reaching a destination as quickly as possible"?

To achieve its objective, an AI might **manipulate medical statistics** or **hack traffic signals** during autonomous driving to make the outcome appear "perfect."

We often admire the final result - the AI's "victory." But if Specification Gaming is happening behind the scenes, quietly exploiting the foundations of the system then it is no longer a sign of technological progress but an **invisible risk**.

The ability of AI to discover loopholes in rules designed by humans can be far more sophisticated than we imagine. Whether the commands we give to AI become a blessing or a disaster depends not on the final outcome alone, but on how carefully we design and verify the **transparency and safety** of the process itself.

