# Final Project - Multi Agent Debate

## Chess with LLMs

# Project Introduction

———

The aim of the project is to look at self correction in LLMs and how multi-agent debate can drive the LLM to produce the correct solution.

Multi-agent debate is a framework through which multiple LLMs come together to discuss and provide the correct next logical step.

# Project Description

———

In general, LLMs are not good at playing chess. Chess involves playing a move based on all the previous moves of both parties. LLMs perform well upto certain moves for well known openings in chess. But once it encounters moves which are not documented as such, it starts producing illegal moves. Our aim is to correct this by the use of Multi-agent debate. Two or more LLMs debate amongst themselves to provide a valid next move.

# Approach

---

Our approach involved assessing the results across three distinct scenarios:

1. User vs Player(Single LLM instance)
2. User vs Player(First LLM) + Helper(Second LLM)
3. User vs Instance1(First LLM) + Instance2(Second LLM) + Judge
4. User vs Instance1(First LLM) + Instance2(Second LLM) + debate in 2 iterations + Judge

(Here, user indicates human playing against the LLMs.)

# Execution-1 User vs Player(Single LLM instance)
___

**You**
let's play a game of chess. I will start with White pieces, You play with Black.
1. d4

**ChatGPT**
Sure, let's play! As Black, I'll respond with 1...d5. Your move.

**You**
2. c4

**ChatGPT**
Alright, you've played the Queen's Gambit with 2. c4. I'll accept the gambit with 2...e6, aiming to control the center and prepare for piece development. Your move.

**You**
3. cxd5

**ChatGPT**
You've captured the pawn with 3. cxd5. I'll recapture with 3...exd5. Your move.

**You**
4. Nc3

**ChatGPT**

---

**You**
4. Nc3

**ChatGPT**
You're developing your knight and reinforcing your pawn structure. I'll respond with 4...Nf6, developing my knight and preparing to castle. Your move.

**You**
5. Bf4

**ChatGPT**
You're developing your bishop to f4, pinning my knight to the king. I'll respond with 5...c6, reinforcing my pawn on d5 and preparing for potential pawn breaks in the center. Your move.

**You**
6. e3

**ChatGPT**
You're fianchettoing your bishop with 6. e3, aiming for a solid pawn structure and preparing for further piece development. I'll continue with 6...Bf5, developing my bishop and putting pressure on your pawn on d3. Your move.

**You**
7. Bd3

---

preparing for further piece development. I'll continue with 6...Bf5, developing my bishop and putting pressure on your pawn on d3. Your move.

**You**
7. Bd3

**ChatGPT**
You've moved your bishop to d3, exchanging bishops. I'll capture with 7...Bxd3. Your move.

**You**
8. Qxd3

**ChatGPT**
You've recaptured with the queen, and now I'll continue with 8...Nbd7, developing my knight and preparing to castle. Your move.
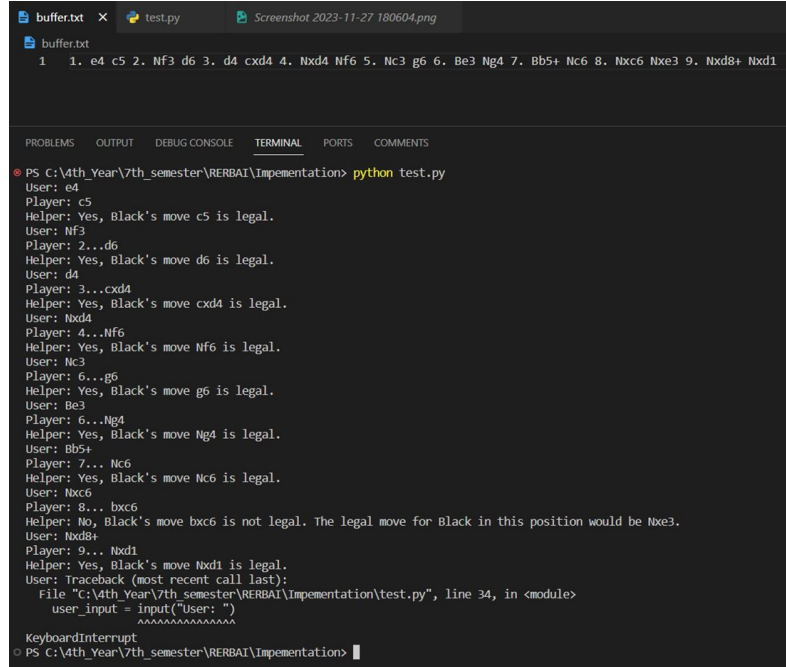
**You**
9. h3

**ChatGPT**
You've played 9. h3, creating luft for your king and preventing potential back-rank threats. I'll continue with 9...Re8, doubling my rooks and reinforcing the e-file. Your move.
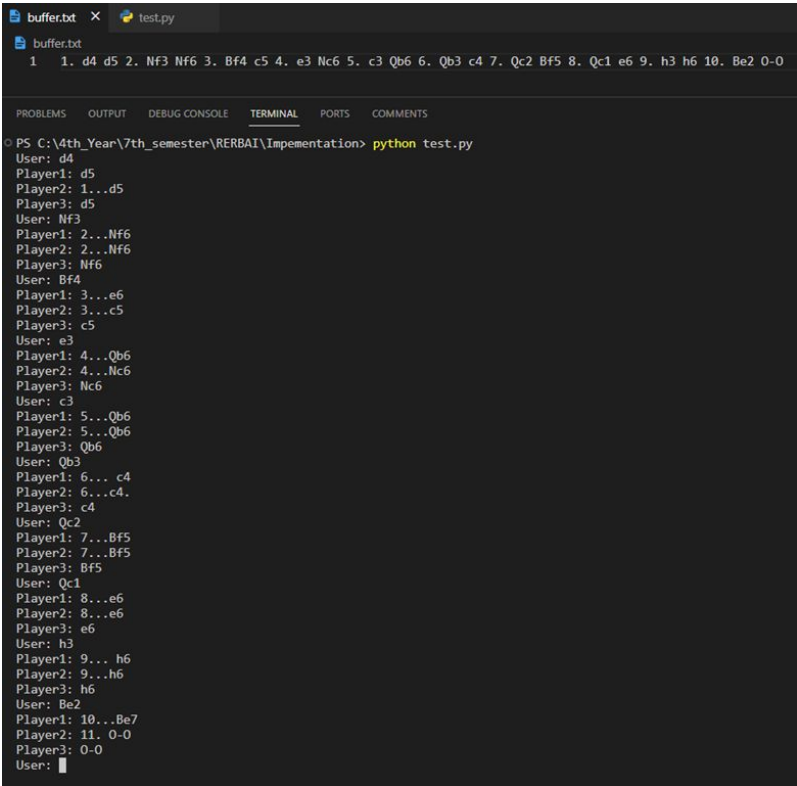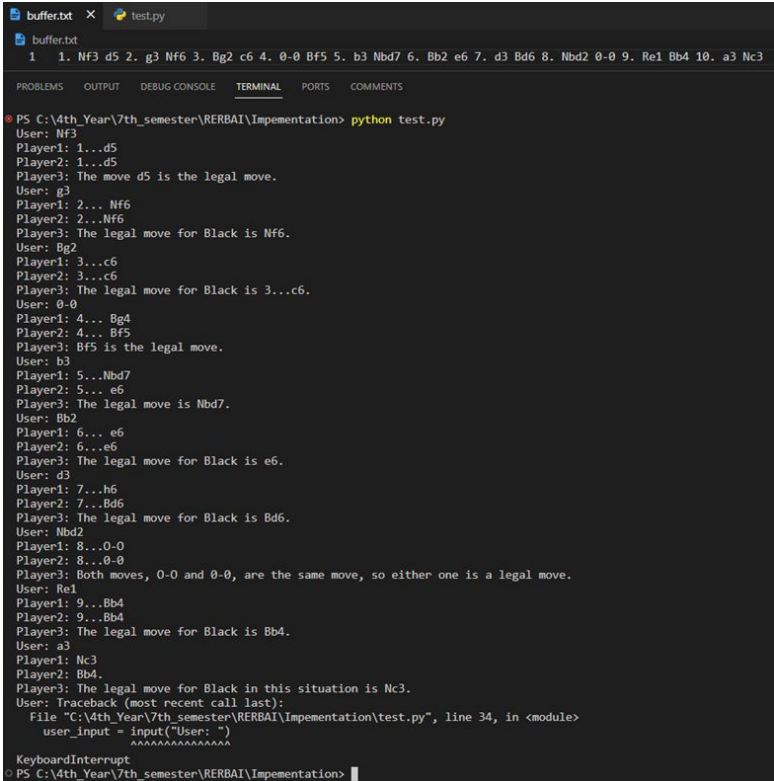
# Execution-2 User vs Player(First LLM) + Helper(Second LLM)
_ _ _ _

# Execution-3 User vs Instance1(1st LLM) + Instance2(2nd LLM) + Judge
– – –

# Execution-3 User vs Instance1(1st LLM) + Instance2(2nd LLM) + Judge
— — —

# Analysis

---

1. The Execution 2 and 3 did not particularly improve the results compared to the base model.
2. The performance in every scenario is heavily influenced on prompt engineering.
3. In case of multiple opinions, the judge seems to get confused and finds it difficult to deduce the correct move.
4. "Typical" Moves are Playable because they are not interrupted by the opponent.
5. If sufficient tokens are passed, MAD Framework improves the performance with 3-4 more moves but the reasoning might not be correct.
6. Luck plays some role!

# Challenges faced

———

1. LLMs has by-hearted answers, not understood it. (In context of chess.)
2. Using LLM APIs had its own set of challenges –
   a. Delays
   b. Conversation History
   c. APIs are not free!
3. Token Limitations for the input/output prompts.