

UNIVERSITY OF ST ANDREWS

CS5199 FINAL REPORT

A Strategy Video-Game for Collaborative Agents with a Personality and Humoristic Dialogues

Author:

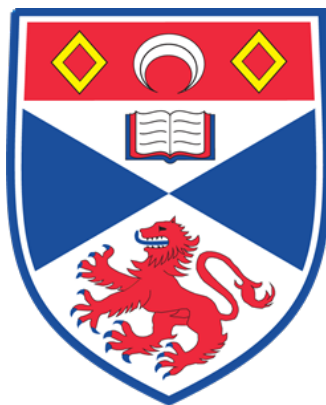
JORDAN MACKIE

Supervisors:

DR ALICE TONIOLO

CHRISTOPHER STONE

January 9, 2020



Abstract

We built a system of collaborative agents with personality driven decisions and humoristic dialogues with the goal of providing entertainment and also providing a human friendly way of interpreting agent interactions in multiagent systems. The agents were designed to play chess and given goal-based personalities that would affect their decisions and how they would interact with other agents. We gathered feedback through allowing volunteers to play against our agents and asking questions covering usability and enjoyment. Finally, we present our results and discuss further work or areas of improvement that could be taken up.

Declaration

I declare that the material submitted for assessment is my own work except where credit is explicitly given to others by citation or acknowledgement. This work was performed during the current academic year except where otherwise stated. The main text of this project report is words long, including project specification and plan. In submitting this project report to the University of St Andrews, I give permission for it to be made available for use in accordance with the regulations of the University Library. I also give permission for the title and abstract to be published and for copies of the report to be made and supplied at cost to any bona fide library or research worker, and to be made available on the World Wide Web. I retain the copyright in this work.

Contents

1	Introduction	3
2	Context Survey	3
2.1	Multiagent Systems in Entertainment	3
2.2	Multiagent System Frameworks	4
2.3	Models of Personality and Emotion	4
2.4	Collaborative Argumentation	5
2.5	Natural Language Generation	6
2.6	Computational Humour	7
3	System Design	7
3.1	Front-End	8
3.2	Agent Gateway Server	9
3.3	Jade Agent Platform	11
4	Agent Implementation	13
4.1	Game Agent	13
4.2	Piece Agent	14
4.2.1	Game State	14
4.2.2	Personalities and Reasoning	15
4.2.3	Single Thread of Conversation Protocol	16
4.2.4	Conversation Planner and Argumentation	18
4.2.5	Natural Language Generation	21
5	Evaluation	23
5.1	Strategic Analysis	23
5.2	Conversation Analysis	26
5.3	User Feedback	26
5.4	Improvements	29
6	Future Work	29
7	Conclusion	30
8	Acknowledgements	31

1 Introduction

Multiagent systems are the next step to increase the level of autonomy that can be provided by technology. Agents that are able to learn, adapt, and negotiate with other agents to achieve their goals allow for complex problems to be solved or modelled without human intervention, such as monitoring and maintaining national power grids [1].

Often, developers and users will anthropomorphise these agents when describing their behaviour. This project aims to encourage this by implementing agents with a model of personality that affects their choice of actions, by rendering the negotiations between agents in a natural language, and by using models of humour to make the interactions between agents entertaining.

Chess was chosen as our strategy game as multiagent implementations have been previously investigated, though not in a configuration where agents could have conflicting goals in the same team. However, our model for determining agent actions can be applied to any other strategy game or multiagent system.

2 Context Survey

2.1 Multiagent Systems in Entertainment

Many modern video games involve the user managing multiple characters to achieve some goal, such as producing in-game resources or defending against an opponent. This sort of problem lends itself easily to multiagent systems. Instead of having one artificial intelligence engine driving the actions of all the characters, developers can create agents with a limited set of actions and some concept of progress towards their goals and allow emergent behaviour to find a solution to the problem, sometimes in surprising ways.

[2] discuss how multiagent systems can be applied to create more realistic worlds in sandbox games. By implementing the environment, objects, and non-playable characters (NPCs) as agents, developers can create a world that reacts and adapts the player, but also operates in isolation from the character to provide a realistic setting. The concept of personalities is also mentioned as a way of allowing similar NPCs to exhibit different slightly different behaviours, such as having aggressive or relaxed driving styles.

Even in games with very simple rules and logic, multiagent systems can find interesting and complex solutions. OpenAI implemented hide-and-seek using agents, where hiders avoid the line-of-sight of the seekers [3]. The game was played in a world with randomly generated walls and objects such as ramps (for climbing over walls) and blocks (for forming barricades). What made this fascinating is that the agents were not incentivised to use these objects, but after repeatedly playing and learning, both the hiders and seekers created strategies such as blocking each other in a safe area, and even removing the objects from the other team before hiding.

Knowing that multi-agent autocurricula can realise strategies not considered

by humans, [4] discuss how they can be applied to strategy board games such as Diplomacy and Risk. They describe a generic framework for supporting agent-based competitive bots for board games and then implement bots for the aforementioned games. Their research suggests that for games with a large number of units and a large action space, a multi-agent approach can identify effective strategies quicker than the exhaustive methods used for chess engines.

2.2 Multiagent System Frameworks

Most multiagent systems have similar requirements, and so several frameworks have been developed to bootstrap their development. [5] provide a very thorough survey of the frameworks currently in use today and highlight various features and drawbacks in their implementation.

The most popular is the Java Agent Development Framework (JADE) [6]. It conforms to the FIPA standard, which is a protocol for agent communication that involves defining performatives (e.g. request, inform), language name, and other message meta-data during communication. The benefit of systems that abide this standard is that they are able to interact regardless of the technology used to implement them. JADE provides other important features such as base classes for agent functionality called behaviours, a directory facilitator (DF) agent implementation to allow agents to find each other, and an Agent Management Service that manages and tracks the lifecycle of agents and allows them to move between containers (which can exist on multiple hosts). To aid during development, JADE also includes a GUI for debugging and manually interacting with agents.

Other frameworks provide APIs that encourage certain design patterns. For example, Jason was built to support the belief-desire-intention (BDI) design model which attempts to separate the processing of choosing a plan from the execution of the chosen plan [7]. Jason provides many of the same functionalities as JADE, but the latter was chosen due to being slightly more flexible.

2.3 Models of Personality and Emotion

Creating a truly immersive video game requires characters that the player can empathise with. Robots have been shown to be able to influence human behaviour as an authority figure [8] and when begging not to be turned off [9] by expressing emotions. [10] created and demonstrated a model of personality and emotion that would allow agents to react differently to the same stimulation. For example, when a agent with an 'introverted' personality is offered help, they are less likely to accept it due to the prolonged interaction it would entail.

[11] also created a model that would use personality, emotion, and social relationships to determine the behaviour of NPCs in a video game. The frequency and tone of interactions between NPCs as well as the NPC and the player were accounted for when choosing facial expressions and tone of voice during conversations. Test subjects described feeling especially attached to the NPCs that utilised this model.

A useful aspect of multiagent systems is that agents can be developed in isolation but still interact (e.g. an agent that searches for cheap transport options and an agent responsible for auctioning train tickets could be developed separately with no knowledge of the logic being used by the other). [12] discuss how developing heterogeneous agent systems using personalities and social structures could help when dealing with third-party agents that have been constructed to lie and exhibit selfish or uncooperative behaviour.

2.4 Collaborative Argumentation

Knowledge is distributed in a multiagent system therefore specialist agents need to be able to alter the beliefs of others by appealing to their individual goals. [13] implemented a framework that achieves this based on a scheme given by [14]:

In the current circumstance R, we should perform action A, which will result in new circumstances S, which will realise goal G, which will promote some value V.

Specifically, they were able to create a framework for multiple parties to discuss and collaborate which could greatly affect the design of multiagent systems that utilise it. By producing an ontology which any agents involved in the discussion understand, goals and circumstances can be conveyed through the use of predicates and concepts. Conveniently, JADE already provides an API for creating ontologies. However, the manual work required to construct an ontology that captures the components and logic of a game of chess was found to be far too cumbersome. While this would be the most interoperable solution, serializable Java objects were used instead for quicker development.

The overhead of argumentation required in multiagent systems can be a quick filter to determine which problems it is a suitable solution for. For example, [15] investigated how quickly agents that used collaborative argumentation to achieve global consistency in a time-constrained task (i.e. escaping a burning building) performed. Their results suggest that optimisations or different approaches to the protocol design would likely be necessary for a real-time strategy game, but it could depend on many factors such as the number of agents, the distribution of knowledge, and more.

Previous work has shown that multiagent chess systems do not generally perform well strategically compared to normal tree based models for chess if the agents are given any limitations that restricts their knowledge. [16] limited pieces such that they could only evaluate their immediate surroundings, and all pieces would reach a consensus based which focused on maximising captures and minimising losses. Our solution will not limit the agents knowledge of the board, but pieces will still have to reach some consensus on which move to take based on their own individual values.

2.5 Natural Language Generation

In order to make the discussions between agents more entertaining and user-friendly to observe, our project involves using natural language generation (NLG) methods to translate the messages or agent intent to plain English. [17] built an NLG system for a particular ontology language (W3C Ontology Language a.k.a. OWL), which is able to construct texts corresponding to objects and their properties in the ontology in English and Greek. [18] achieved similar results for one of the other Semantic Web Formalisms (specifically RDF) and were able to produce instructions for cooking recipes from a corpus of data in a far less human-friendly format. They were even able to adjust the level of technical jargon in the resulting text to accommodate for the readers familiarity with cooking. The ontology allowed for optimised searching of the parse trees but they also required a domain corpus (i.e. existing textual recipes) in order to properly train their natural language generator. Their solutions are specific to the Semantic Web project, and would not produce "conversations" like we are planning to do, but provide a good foundation for implementing this functionality.

Generic frameworks have also been created that are not based strictly on ontologies which were also considered. For example, SimpleNLG [19] is an NLG engine that allows for valid English phrases to be constructed through a very simple Java API. However, after some time experimenting with this library, it was realised that it was very difficult to reliably produce a variety of phrases, though it was useful for controlling features such as grammatical person and tense.

[20] created a robot capable of providing dialogue while playing against a human at tic-tac-toe, but used a slot-filling technique based on templates that suffered from little flexibility and left no room for personality in dialogue creation. Contrarily, systems such as Siri and Alexa are examples of successful dialogue systems that have been able to provide assistance to humans in the form of natural language conversations by avoiding the template based approaches.

These systems instead take advantage of machine learning techniques such as deep recurrent neural networks (DRNN). [21] surveyed projects that applied these modern techniques for NLG, and concluded that in order to account for aspects such as current state of the program during the execution of conversation as our project requires, a hybrid of the machine learning and rule based approaches may be necessary.

The dynamic finite-state machine implementation of a chatbot called Tartan [22] was also especially useful for the project, as it focused on natural conversation flows rather than the extraction of information from a human user like most other chatbot research. Their conversation state machine design was combined with the RiTa library's RiGrammar (a probabilistic context free grammar) in order to produce a variety of phrases and conversation flows.

2.6 Computational Humour

Simply reading the discussion between agents would certainly not make for an entertaining game, and so adjusting the dialogue or agent behaviour is required. However, artificial humour also has benefits for everyday computing: a smart-phone that is able to sympathise with the user when it is unable to connect to a weak WiFi network or cheer them up with a joke when they miss their bus is one that the user is more likely to enjoy using.

Unfortunately, given humour is an incredibly contextual and culture driven concept, it is no surprise that computational humour is well researched yet still not 'solved'. Jason Rutter describes why using artificial intelligence for humour is particularly challenging:

"Humour is a very interesting way to look at artificial intelligence because at some point something has to have two meanings, which is not easy to do with a computer." [23]

There are three main theories of humour that are used for computational humour: superiority theory (laugh about the misfortune of others), relief theory (using taboo subjects to release tension), and incongruity theory (using lexical and structural ambiguity). The last theory is currently the most popular, and was used during the development of JAPE [24]. JAPE identifies features such as homophones to produce jokes like:

"What is the difference between leaves and a car? One you brush and rake, the other you rush and brake."

However, JAPE does not allow for much user interaction, and most textual humour generation takes the form of simple puns and riddles. Our agents could instead generate situational humour by surprising the user by diverting suddenly from expected behaviour. The Suslov Model of humour accounts for the fact that we subconsciously predict expected conclusions to situations and phrases, and that contradictions in what is the most likely direction for a conversation or situation to what was previously predicted can create a humour response [25].

These probabilistic models currently appear more hopeful than the rule based models used by older implementations such as JAPE. Instead, [26] used an existing corpus of jokes and a recurrent neural network to produce a system capable of generating jokes and even anti-jokes (jokes without an actual punchline).

3 System Design

An existing multiagent development framework was chosen instead of rolling our own in order to abstract away matters such as message delivery and agent life-cycle management. As mentioned earlier, JADE was chosen specifically due to its flexibility and its extensive coverage of the FIPA standard that could allow our agents to interact with other implementations with ease.

With the intention of possibly serving the game over the internet once complete, the system had to be designed to allow for communication with agents from outside of the agent platform. The 'JadeGateway' API was used with a Spring Boot server and React UI to to achieve this. Figure 1 provides a high level overview of the technologies at each layer.



Figure 1: Technologies used at each layer.

3.1 Front-End

React [27] was chosen due to it being the most familiar web framework for the developer of this project, and because a component already existed for rendering chessboards [28]. React uses a virtual DOM in order to efficiently re-render components based on changes in state. For example, when a new message arrives, the state of the chat transcript component is updated to contain the new message which triggers an update to only the chat transcript and its children that are also affected.

As the project grew, Redux [29] was introduced in order to help with state management of components. Redux provides a single global state which components 'connect' to in order to receive updates. Updates are 'dispatched' as actions, which are then passed through a pipeline of reducers which update parts of the global state depending on the action and its payload. Middleware can also be introduced to allow certain actions to create chains of updates (e.g. making an asynchronously API call, then dispatching another action when the response is received), but the reducers themselves are always pure functions. This functional approach to managing components made the project easier to test and errors easier to trace and correct. Figure 2 shows the parts of the state that are managed by each reducer.

One of the challenges faced were that the abstraction of the chessboard component provided by a third-party library did not allow for atypical visual elements such as adding name tags to each piece or dialogue boxes. This then required manually implementing an overlay that would draw another chessboard grid, track positions of pieces, and draw the name tags and dialogue boxes in the correct locations.

The UI is composed of two main views: configuration and gameplay. The configuration view (Figure 3) allows the user to choose if they want to play against the agents or watch two sets of agents play against each other. By clicking on the pieces they will not be controlling, the user can also name the pieces and define their personality type.

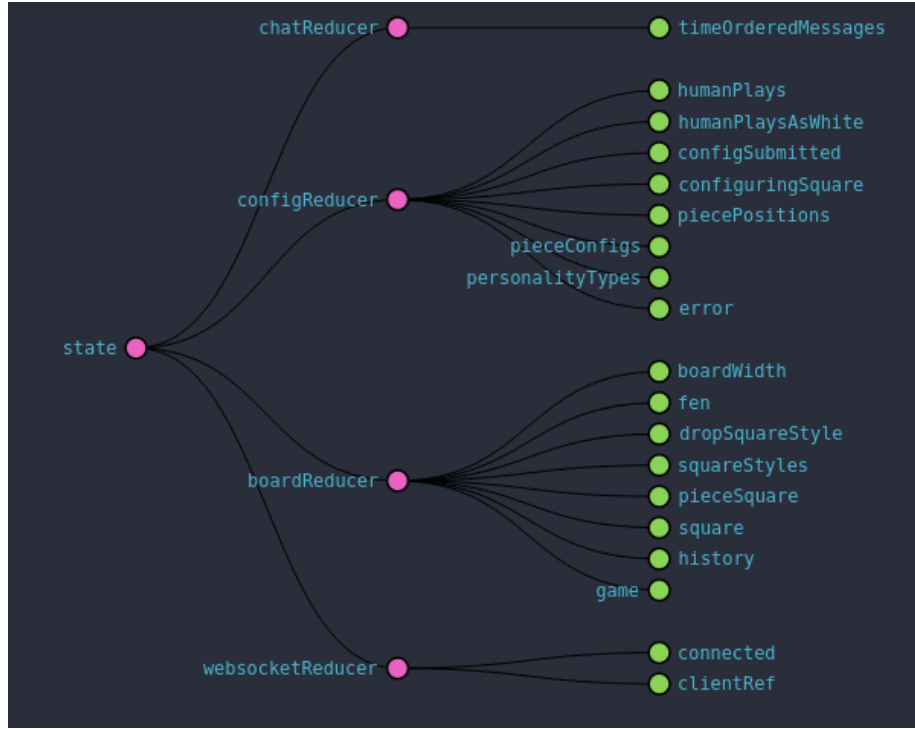


Figure 2: Redux reducer tree showing which parts of the state are managed by which reducers.

On submit, any unconfigured pieces are filled out at random from the set of available personalities and a set of names. The game configuration is then sent as a POST request to the gateway server, which verifies and creates the resources necessary for the game and returns its ID to allow the client to publish and subscribe to game events at the correct web socket endpoint.

Figure 4 shows the sequence diagram for how resources are created once the user submits the game configuration

The game view (Figure 5) also shows the game board with name tags over the corresponding pieces and any dialogue next to the speaking piece. The chat and move history is also rendered in a scrolling transcript at the bottom of the page in alternating colours for clarity. When hovering over pieces the user is able to see all currently available moves to them. Moves are verified client side by the chess.js [30] engine before either being rejected or sent to the server. Figure 6 shows how human moves are broadcast to all subscribers.

3.2 Agent Gateway Server

The gateway server has several purposes:

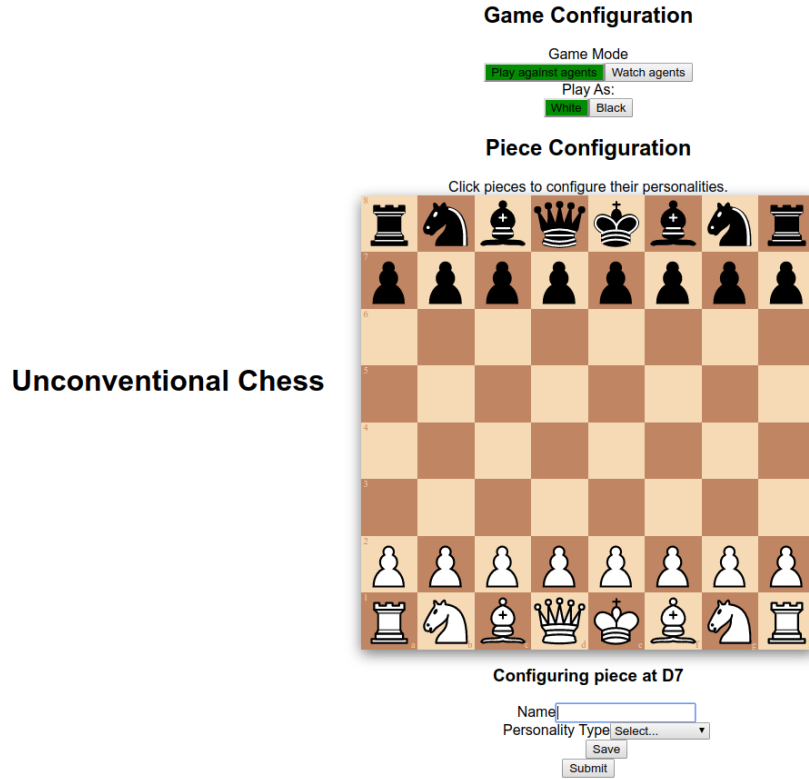


Figure 3: Configuration view.

- Serve the web app
- Handle game resources
- Route messages between users and relevant agents

JADE does have support for agents serving web resources directly, but the JadeGateway API was used instead as it was simpler to set up. The JadeGateway API effectively creates an Agent for the server, which can then send units of functionality called 'Behaviours' to be executed by the agent (e.g. sending messages to the agents within the platform). This could potentially be a bottleneck if the number of users of the server was to grow large enough due to agents being single-threaded by default, so an implementation at scale would likely need to change this interface. For the purposes of this project however, it was more than sufficient.

On start-up, the server uses the details in a configuration file to connect to the JADE platform. It also exposes a limited set of HTTP endpoints to allow the client to create games and to query available personality types. An in-memory message broker is used for handling websocket sessions once a game

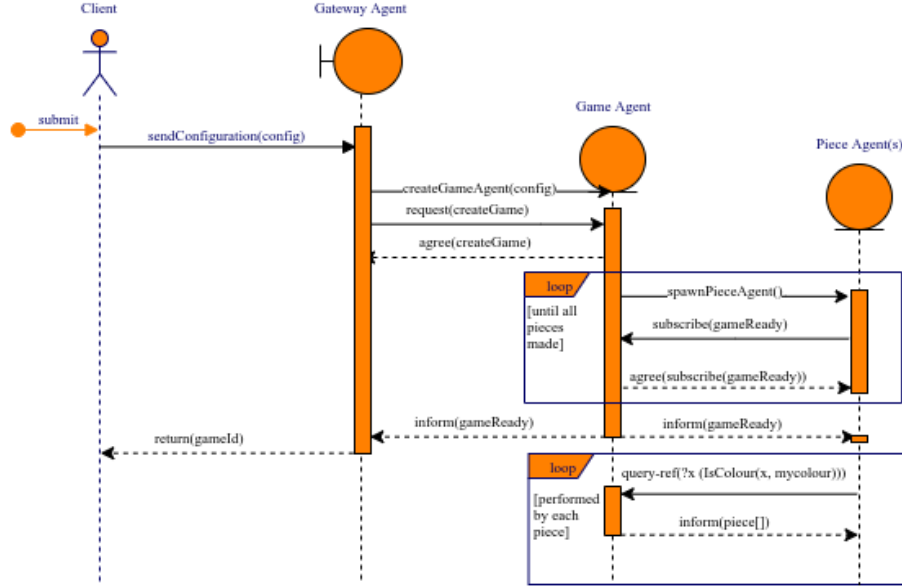


Figure 4: Game creation sequence diagram.

is set up so that clients can subscribe to chat and move messages, and send their own moves.

In order to receive messages from agents within the container, the Jade-Gateway polls for incoming messages using the `ListenForGameAgentMessages` behaviour, which is provided a message translator and a message handler for translating `ACLMessages` from agents to a data model the server can work with and passing the message to the correct handler. This behaviour allowed the server to utilise JADE's scheduling system which would only run the behaviour if a new message was received, instead of busy waiting for incoming messages.

3.3 Jade Agent Platform

Finally, the agent platform is where all of the agents are hosted. Here we will give a rough overview of how JADE agents schedule their functions, while in section 4 we discuss the actual logic that drives our agents.

JADE agents are single-threaded by default, and execute behaviours according to the flow diagram shown in Figure 7. The basic JADE behaviour implements the methods shown in figure 7, but the framework also provides some helper implementations such as the `FSABehaviour` which was used heavily during this project.

Behaviours are scheduled such that they are executed repeatedly round-robin until either they are complete (`done()` returns true) or `block()` is called during the `action()` method. The latter puts the behaviour into a waiting queue so that it is only triggered again once the agent receives a message.

Behaviours can filter the messages that are received using the `MessageTemplate` class, which uses pattern matching to only allow for messages containing features such as a specific performative (e.g. `REQUEST`, `PROPOSE`) or a specific conversation ID. This ensures that behaviours do not consume messages that were not intended be consumed by them.

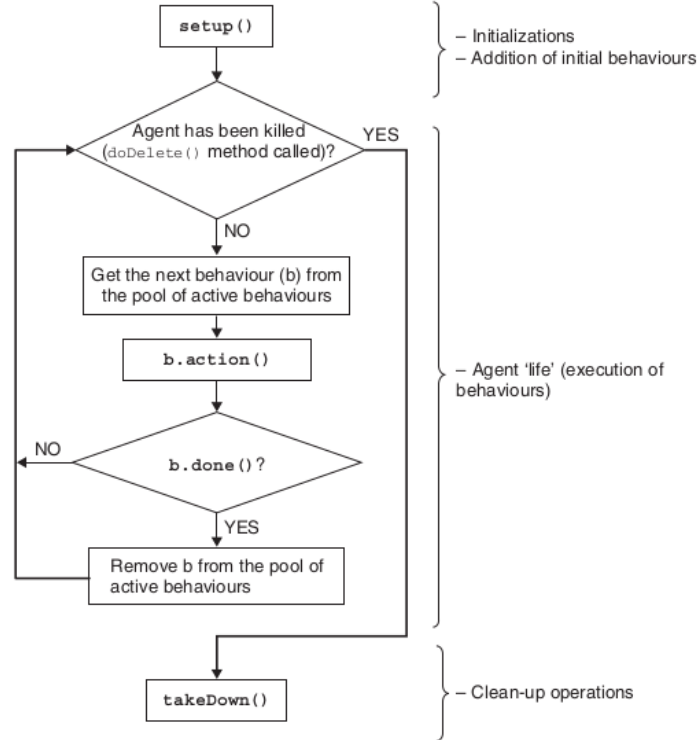


Figure 7: JADE Agent behaviour scheduling and life-cycle [6].

4 Agent Implementation

Our implementation involves two different types of agent: game agents and piece agents. Game agents are responsible for maintaining the state of the game and initialising the piece agents, and piece agents represent pieces on the board, each trying to achieve their own individual goals.

4.1 Game Agent

The game agent upon creation waits for a request to begin a chess game. It also answers any queries regarding the current status of the game and information

about the pieces being represented by agents. The FIPA standard protocols for agent interaction [31] were followed so that the game agent could provide a consistent API for any implementation of piece agents. For example, the request interaction is used for triggering moves as shown in figure 6.

Requests for information are semantic logic expressions with variables which the game agent then populates in its reply. For example, agents wishing to receive updates whenever a new move is made will send a message with a SUBSCRIBE performative and the expression:

$$(iota ?Move ("MoveMade" ?Move)) \quad (1)$$

This will then be replied to every time a new move occurs with the variables given values, for example when a piece moves from B8 to C6:

$$\begin{aligned} ((= (iota ?Move ("MoveMade" ?Move)) \\ (Move : Source (Position : Coordinates B8) \\ : Target (Position : Coordinates C6)))) \quad (2) \end{aligned}$$

Once the game has started, the game agent adheres to the FSA shown in figure 8.

The actual chess logic is managed by a third party library [32]. While this saved a lot of time by providing a way of quickly validating moves, generating valid moves, and tracking the board state, it suffered from the same problem as the chessboard component used by the front-end: it did not provide a way to differentiate between each piece of the same type. This resulted in a wrapper class being implemented which would maintain the mapping between the set of our custom piece representation and the third party board pieces. In hindsight it was clear that it would have been simpler to have implemented the chess logic ourselves, but the wrapper solution was sufficient.

4.2 Piece Agent

The piece agent implementation is the focus of this project. They are initialised with their piece type, colour, personality type, and starting position on the board. They are notified by the game agent once the game is ready and whenever a move is made successfully either by the human or another agent.

4.2.1 Game State

The GameState class wraps our chessboard wrapper to provide another layer of abstraction in order for the state of the board to be an immutable object. This was done to avoid complicated undo logic when agents wanted to 'test' moves and evaluate the board after a move is made. Instead, when the applyMove(PieceMove move) function is called, a copy of the board is made. This involved copying the set of alive and captured pieces, as well as updating the piece objects themselves if any changed position.

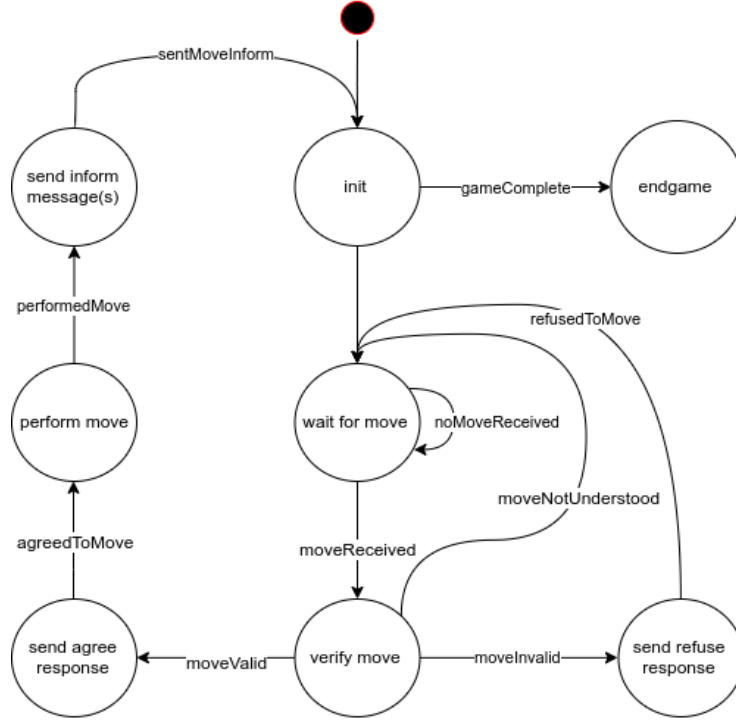


Figure 8: Game Agent game handling FSA.

In order to reduce the amount of objects being created and destroyed by this process, piece instances would only be copied if they had changed state in some way (captured/changed position).

4.2.2 Personalities and Reasoning

Agent personalities were designed based on the goal-based trait/value model defined by [12]. Figure 9 shows the UML class diagram for our implementation. The personality is initialised with a set of traits (e.g. shy, aggressive) which each have a set of values. Values define what the agent will interpret as good and bad moves based on some criteria (avoid other pieces, maximise captures). The value class is abstract and is extended by implementations that when given a witnessing piece, the current game state, and a move, are able to produce a response to that move.

This implementation allows different traits to share goals. The move most preferred by an agent is the one that appeals to the most of its values, and so consensus can be reached by finding the move that appeals to the most pieces, even if they have different values.

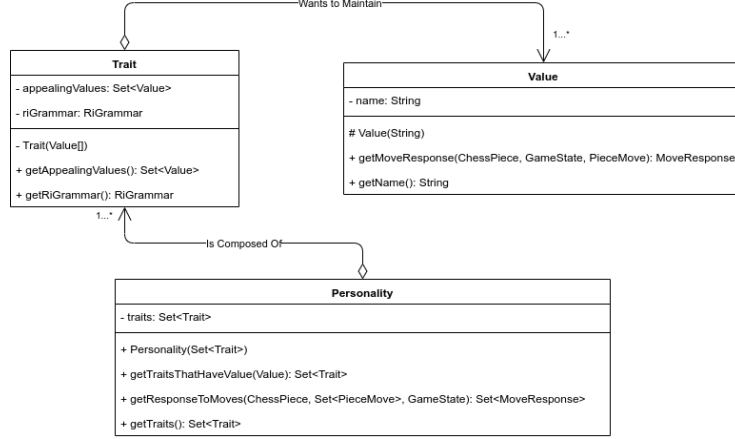


Figure 9: UML Diagram of Agent Personalities.

4.2.3 Single Thread of Conversation Protocol

In order to make the conversation between agents seem natural and remain easy to follow, a protocol had to be established to ensure that one agent would speak at a time. The project was initially constructed with conversation flow dictated strictly by a finite state machine. This solution was incredibly restrictive as it required providing a transition path for any conversation path that we wanted to support. After reviewing the work by [22] and finding that this strict FSA approach was inherently detrimental to producing natural conversation flows, the original design was thrown out and replaced with the simpler and far more flexible model shown in Figure 10.

The conversation FSA designed by [22] was for a human-agent chatbot, and so was adapted to include a token-passing protocol to determine which agent would speak next. One of the benefits of this design is that it allowed for the agents to chat throughout the game, not just during their turn. This meant that the agents could provide ambient conversations and respond to events like the previously made move while the human player chose their next move.

At the start of each turn (i.e. when a new move is received or the very first turn of the game), each agent returns to the initial state to ensure that they are all synchronised. The king is assumed first to speak. The other pieces listen for any messages from the speaker and the speaker chooses what to say. Once sent, the speaker begins the new speaker election by sending a call for proposal (CFP) message to all alive pieces (and themselves), and waits until every pieces has submitted a proposal to be the next speaker. Finally the speaker chooses from the set of proposals, sends a rejection to all but the chosen speaker, and then the cycle begins anew.

If the agent chooses to make a move while they are speaker, they transition into the corresponding state as shown in figure 10. They remain here until they

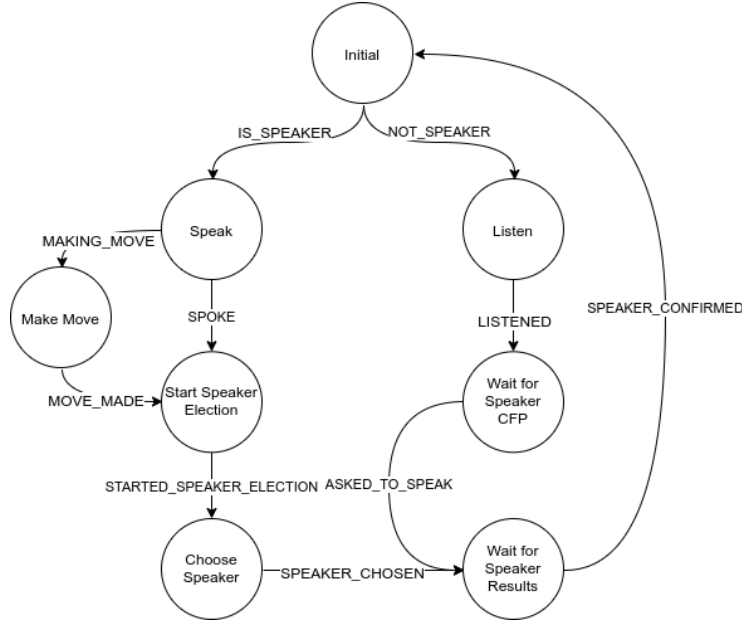


Figure 10: Conversation protocol FSA

receive either a refuse or agree response from the game agent and then return to the normal conversation cycle.

The requirement that all agents reset to the initial state at the start of each turn was necessary to avoid race conditions during the speaker election. For example, if a move was received that resulted in one of the agents being captured, that agent would immediately remove themselves from the game. However, if that agent was the speaker during an election, then all other agents would be stuck waiting for the next speaker to be chosen. Whilst other solutions would allow for the conversation flow to continue, it was deemed much simpler to reset the conversation FSA during each turn.

As a side effect of this decision, messages such as the speaker election CFP could be waiting in an agents message queue if they reset before receiving it. To avoid this resulting in a loss of synchronisation between agents, a new conversation ID is assigned during each cycle of the FSM which is composed of the turn ID (i.e. how many moves have been made so far) and the round ID (how many cycles of the conversation FSA have occurred this turn). The conversation ID is used to filter incoming messages to ensure that irrelevant messages are ignored. Unfortunately this could result in a build-up of expired messages over time, but given most chess games are around 80 moves [33] and roughly 5 messages per move would be lost in the queue, the leak would not be substantial. Especially when it was considered that the messages would be garbage collected once the agent representing the piece was removed from play.

The turn and speaker rotation counters is provided by the ConversationCon-

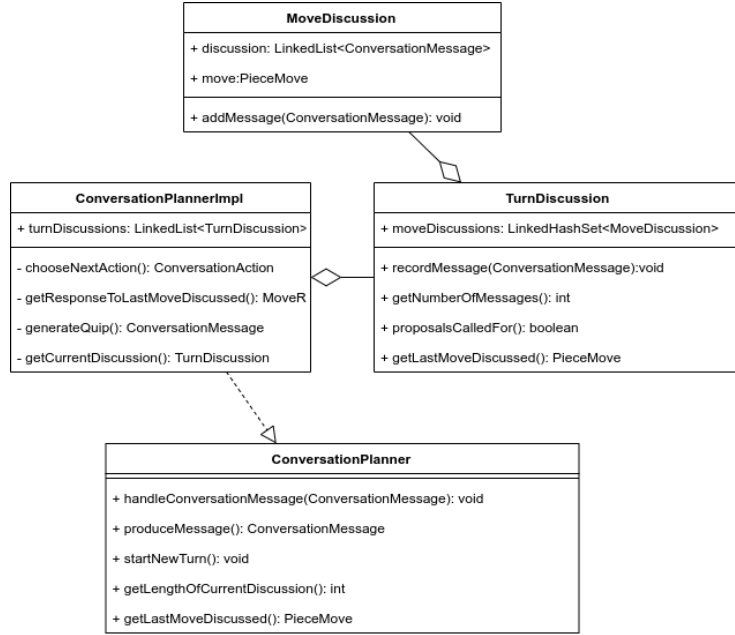


Figure 11: Move Discussion Class Diagram

text class, along with the current speaker agents AID (Agent Identifier) and the conversation planner.

4.2.4 Conversation Planner and Argumentation

Each agent maintains a conversation planner which tracks what moves have been discussed so far, and produces a response for this agent once it has its turn to speak. Figure 11 shows the classes used to maintain this information.

Whenever an agent enters the speaker state shown in figure 10, they call the `produceMessage()` function to choose their next action. If this action involves making a move, they'll queue a behaviour to send the move to the game agent. Regardless, they send their chosen conversation message to all the other agents (including themselves).

Listener agents receive the conversation message and record it by calling the `handleConversationMessage(ConversationMessage)` method. The conversation planner delegates storing the message to the current turn discussion (the head in the list of discussions). The turn discussion will fetch the move discussion for the move in the given conversation message or create a new one if one doesn't already exist and add the message to the relevant move discussion. Once the message is added to the move discussion, the discussion is removed and re-added to the set of move discussions in order to move it back to the top of the conversation 'stack' (since insertion order is maintained by Java's

Table 1: Conversation actions with verbalised examples.

Conversation Action	Example
Acknowledge	"Ok."
Ask for Proposals	"What should we do next?"
Perform Move	"Going to move to A5."
Propose Move	"What about B4 to B3?"
Revisit Move	"Let's go back to talking about the move to G7."
Voice Opinion	"I like that idea!"
* Propose Alternative	"but we could also move to G3."
* Propose Compromise	"moving to G3 would do that too"
* with Justification	"because it would capture an enemy."
Quip	"I'm a bit scared."

LinkedHashSet implementation).

By maintaining the order that moves have been discussed, agents are able to create arguments that better resemble human conversations by referring to previous discussions. [22] uses a similar technique to allow their chatbot to regress to previous topics of discussion once the user gives some kind of terminating statement to the current topic of conversation.

A move discussion with a null move is used to represent a request for move proposals, which is sent either to start up the argumentation process or when an agent is not satisfied with the move currently being discussed but is does not feel strongly about any moves themselves.

The output of `chooseNextAction()` in the `ConversationPlannerImpl` is the next conversation action to be taken. These actions are essentially transitions which control the flow of the conversation and are based off the interruption FSMs in Tartan [22]. Table 1 lists examples of the conversation actions possible. The actions prefixed with wildcards can be used in conjunction with other actions. For example, an agent can voice opinion with justification, or voice their dislike and propose an alternative.

Each of these actions are implemented as classes that extend the `ConversationAction` class shown in Figure 12, and once performed they return a `ConversationMessage` instance with all the information needed for other agents to update their own understanding of the conversation context. The actions can also extend each other if their affect on the conversation flow is the same but the natural language representation needs to differ (for example, revisiting a previously discussed move and proposing a new move both begin a new move discussion but would be expressed differently in english).

Conversation messages contain the name of the piece that sent it, the natural language representation of what the agent wanted to communicate, and a move response. The mapping between each field and agent intentions are described in table 2. The alternative proposition field starts a new move discussion if present and the same logic is applied recursively as though it was the the move response in the original conversation message.

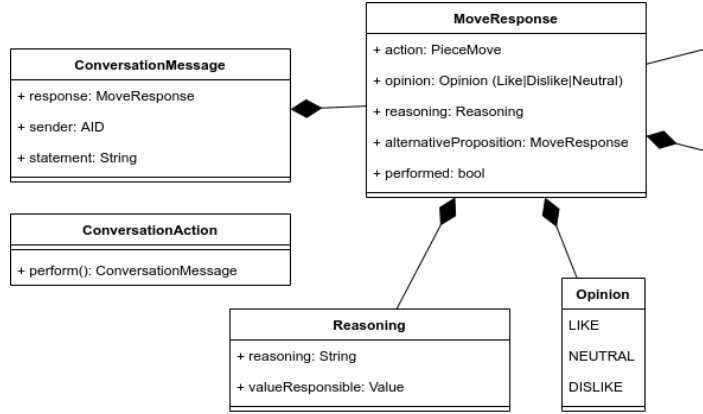


Figure 12: Conversation Message Class Diagram

Table 2: Conversation message fields and the agent intention.

Field	Predicate	Agent Intention
response	== null	Requested proposals
response.opinion	!= null	Voiced opinion
response.reasoning	!= null	Gave justification
response.alternativeProposition	!= null	Proposed alternative
response.performed	== true	Performed move

This information is enough to allow agents to determine how other agents feel about the move being discussed and make decisions about what conversation actions would logically follow, but the contents of the statement in the ConversationMessage are completely flexible.

The next conversation action is chosen by building up the set of responses that would make sense given the discussion so far and then selecting one from the set randomly. Most responses are equally likely, but some are given lower likelihoods of occurring (such as performing a move) in order to extend conversation and encourage discussion rather than agents just repeatedly proposing their own preferred moves.

Choosing responses that 'made sense' involved checking the current discussion context and determining the following:

- Is this the first message this turn?
- Did the previous message ask for proposals?
- Have moves been discussed already?
- What is this agents opinion of the move currently being discussed?
- Does a move exist that I can perform?

Table 3: Reactions with verbalised examples.

Reaction	Example
Compliment Enemy Move	"Oh no, that was a pretty good move by their Knight"
React to Enemy Piece Captured	"Haha! We have your Bishop!"
React to Friendly Piece Escaping	"Good, John is safe now!"
React to Last Move Discussed Being Performed	"I'm glad you agreed to that move!"
React to Undiscussed Move Being Performed	"...Maybe ask next time before moving, Kate."

In order to create the illusion that pieces are moving themselves, only the speaker can perform a move and that move has to involve the same piece. This choice will likely result in less strategic moves being made (for example, a piece will not want to put itself at risk, even if it means saving a more important piece), but makes for a more entertaining narrative.

The reasoning field in move responses are populated by the Value of an agent assessing a given move. This gives a natural language explanation for the value of the opinion field in the move response, and the value that was used to compute it. This allows other agents to apply the same value to all other moves (even if they do not have that value themselves) and provide a compromise that also appeals to their own values.

This argumentation feature could be improved by expressing the value as a predicate that other agents could unpack instead of being tied to the Value classes in our chess agent library.

The conversation planner was extended to support reactions as well as discussion actions. This allowed agents to respond both to their moves and to the users moves to try and improve immersion. A hard limit was placed on the number of reactions to the enemies moves on the start of a turn to avoid discussion being delayed by an unrealistic number of reactions from different pieces. The type of reactions that could be performed managed by checking previous conversation frames and the game state resulting from the move performed.

These reactions required slightly more complex checks to the discussion actions as they were dependant on very specific changes in the game occurring to make sense.

4.2.5 Natural Language Generation

The natural language generation was a crucial part of this project, as it had to tie together all other features correctly in order to make their effect visible to the user. The phrases spoke by pieces had to:

- correlate to the opinion of the pieces and the moves being discussed

```

"<ProposeMoveWithJustification>" :
    ["<ProposeMove> <Justification>"]

"<ProposeMove>" :
    ["'getMovingPiece()' should move to '
      getMoveTarget()'"]

"<Justification>" :
    ["because that move 'getJustification()'"]

```

Figure 13: Example of BNF grammar being used for NLG.

- vary enough to not appear scripted
- make grammatical sense when combined
- make sense with respect to previous statements by other agents and themselves
- be consistent with the personalities of the individual pieces

As mentioned in the context survey, the development of the natural language generation aspect of this project was greatly slowed down by starting off with likely the wrong technology for the problem. SimpleNLG appeared to be useful given a better understanding of linguistics, and so a lot of time was spent understanding the technical terms of this field. However, the tool seemed better suited to tasks like report generation than conversation, and so with assistance from the project supervisors, the RiTa library was adopted instead.

RiTa provides a wealth of features for natural language processing and generation, but the deterministic grammar was especially useful for this project. Phrases are constructed from non-deterministic BNF grammars and a data-source for filling in variables such as agent names or move coordinates. A default grammar is provided, but by implementing variations of the default grammar for each trait, agents will express themselves based on the traits of their personality. A root tag is provided for each class that extends ConversationAction. For example, the ProposeMoveWithJustification conversation action used a grammar similar to example shown in figure 13. Multiple values can be provided for each tag by adding them to the array of values for that key, and the probability of each value being chosen can also be defined using square brackets within the string to multiply the probability of it occurring (i.e. [2] will make the value be chosen twice as often as any of the other values in the array).

Each value also uses a BNF grammar to produce the justification in natural language, but these are implemented within the values themselves rather than as separate schemes.

The strings between backticks within the grammars are slots that express the name of a field or a function call that will provide the value for the slot in

the GrammarVariableProvider. This is an object populated with things like the name of the agent being asked to move and the coordinates being moved to.

5 Evaluation

The success of the resulting software was ascertained by analysis its performance ourselves and by gathering user feedback. Issues raised by these assessments were later addressed to produce the final artefact that was included with this submission.

5.1 Strategic Analysis

In order to test the strategic performance of the implementation, the Stockfish engine [34] was used to score the move chosen at each turn by the agents. This was repeated over the course of five games to account for the randomness used by agents during the choice of action. Different configurations of personalities for pieces were used for each game.

Stockfish is one implementation of the universal chess interface and one of the best chess AIs in the world. It uses a complex set of rules to score board positions, which can be used to identify the next best move to make. Centipawns are the unit used for the score produced, which represent one hundredth the value of a pawn (1 point). The king is essentially of infinite value, and so for our data we translate mating moves to centipawns by assigning these moves an arbitrarily large number, as recommended in the python-chess documentation [35].

The moves were tracked by adding a logger to the game agent which would record the current position in the standard Forsyth-Edwards Notation (FEN), and the source and target of the move being made. Both teams were represented by agents and the conversation delay set to a much lower number to increase the rate of discussion. The comma separated value dataset produced was then parsed by a python script. For each position recorded, all legal moves were considered by Stockfish with a limit on the computation time of one hundred milliseconds per move. This produced a score for all legal moves at each position including the chosen move.

In figure 14, the difference between the maximum possible scoring move and the score given to the chosen move is plotted for each turn and for each game. Points closer to the x-axis represent moves that could be considered good. Due to mating moves being of very high value, whenever these moves were not chosen the difference is much higher, and so figure 15 was also plotted with these moves filtered to give a better visualisation of the agents performance over time.

Figure 14 shows that often the mating moves were not taken, and this was likely due to the fact that no value was implemented that would trigger a LIKE reaction if the threatened piece was the enemy king specifically. This highlighted the need for piece value to also be accounted for when choosing moves if the agents were to have any strategic advantage, but this would likely also require

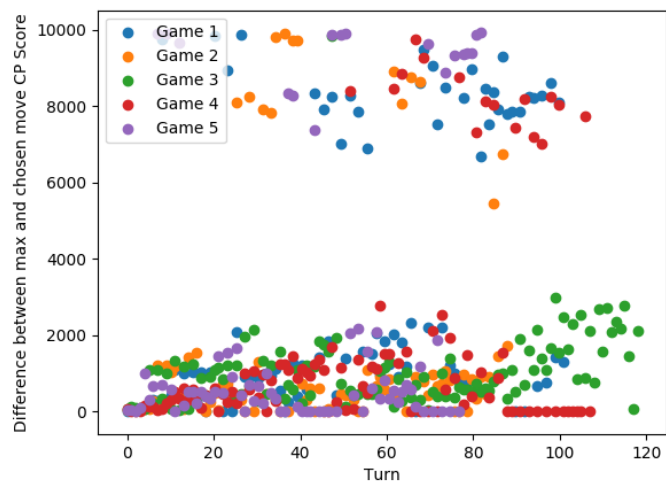


Figure 14: Differences between best move and chosen move score over the course of five games

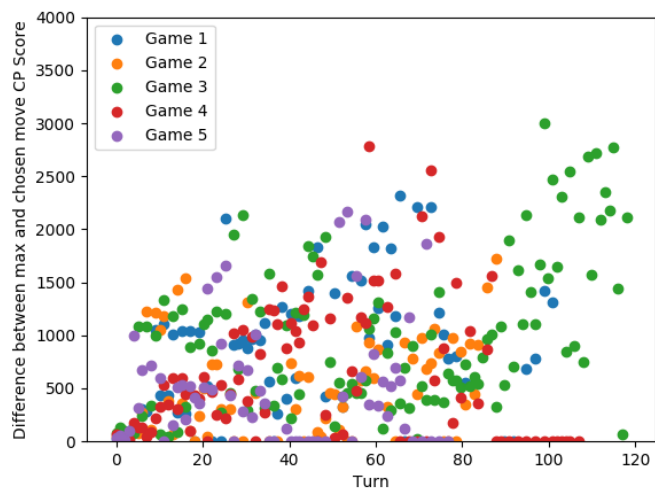


Figure 15: Figure 14 with differences over 4000 removed

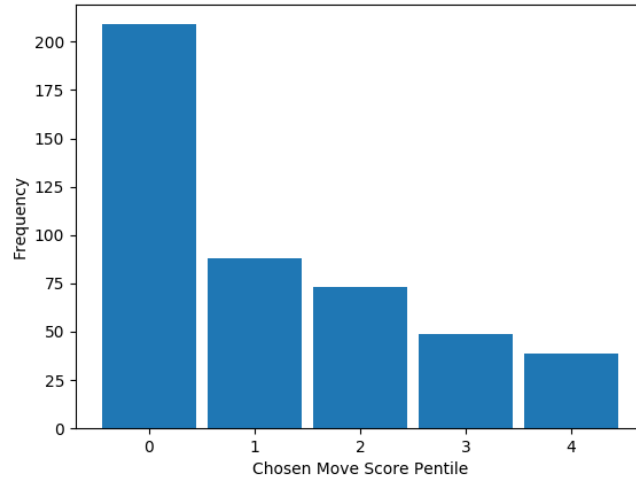


Figure 16: Frequency of chosen move being in each score pentile

introduced a numerical scoring of moves by agents instead of the three value (LIKE, NEUTRAL, DISLIKE) approach taken.

Interestingly, certain games seemed to have produced better moves which could suggest that either

1. Certain personalities were more effective at choosing good moves and these personalities were more represented for that game. For example, the aggressive personality type would maximise captures, reducing the number of opponents, which would be evaluated as high scoring moves due to capturing pieces being worth a large number of centipawns.
2. Certain combinations of personalities were more effective when combined in a given ratio. For example, having a balance of aggressive and defensive personalities would result in moves that would capture pieces whilst maintaining the least number of threatened friendly pieces.

To give a general overview of how well the agents performed, the scores from all moves from all games played were all collected. For each position, all the legal moves were separated into equally sized pentiles by score. The zeroth pentile contained the lowest scoring moves and the fourth the highest scoring moves. For each move, we determined which pentile the move belonged to, then plotted the frequency of the moves occurring in each percentile, and the result is shown in figure 16.

This clearly shows that most moves taken were the least advantageous strategically, though there were still occurrences of good moves. However, if we assume that randomly chosen moves would produce an equal distribution in each

percentile (though this is unlikely), then of the 496 moves considered 90 would have been in the highest percentile. Compared to our results where we have nearly half of that value, there is clearly a need to implement more strategic move evaluations if the game is to be more challenging.

5.2 Conversation Analysis

The fluidity of the conversations produced was challenged throughout development. Due to the complexity of natural conversation, testing for messages that made sense in series was difficult without simply rewriting the same logic being used to choose actions in the tests, and so most of the evaluation of this feature was done by manual inspection.

During the agent-vs-agent experiments discussed in section 5.1, conversation dialogue was not disabled, which meant that at the end of each experiment a transcript was available. Due to the rate of messages coming in, it was impossible to read and follow the discussion in real time. The level of noise when both sides were creating dialogue made it difficult to follow even with the rate of conversation greatly reduced which made it clear that this game mode would likely be improved by only showing the dialogue from one side.

The transcripts showed more variety was likely needed in the speech corpus, as repetition was quite easy to spot when the number of messages were this large. For normal gameplay, this would be less of an issue, but it did not seem like there was enough variety in speech to make a full game entertaining. Unfortunately, this does identify a drawback with our approach: creating believable dialogue required manually curating the dataset which could take a large amount of time. However, the context free grammar approach meant that poorly formed sentences could be easily fixed as the pattern used to create the sentence could be traced.

Determining if the discussions complimented gameplay required simply playing the game and allowing the pieces to talk for a while between turns. During development, this approach identified problems such as too many reactions in a row occurring, or moves being performed too early on in the discussion and cutting it short.

5.3 User Feedback

Since the goal of this project was to essentially create a video game, it made sense to take inspiration from industry standard Game Engagement Questionnaire [36]. This questionnaire is often used during user feedback sessions to measure the absorption, flow, presence, and immersion a game provides. It is comprised of a series of statements such as "I lose track of time" and "Things seem to happen automatically". For each statement, the user completes a Likert scale ranking of how much they agree with it.

Presence is a measure of how much the player feels part of the game. Due to the players pieces being fairly static, the game was not expected to score well in this area. Similarly, flow is a measure of much the players movement and

interactions with entities within the game feel natural. Again, chess does not cater to this metric well due to essentially one method of interaction: moving a piece.

On the other hand, absorption and immersion try to capture how much players become unaware of their actual surroundings due to a game. If the dialogue and conversation flow is believable enough, it is possible that the game could score fairly high in this regard.

Due to the short length of time that user feedback sessions could take, some questions were removed from the questionnaire to allow for other types of questions to be asked specific to this project.

For example, we also ask some questions that would help us evaluate the linguistic accuracy of the conversations produced. Unfortunately, we were unable to find past research evaluating agent to agent conversation for reference. With the growing popularity of voice controlled assistants, a lot of literature existed for evaluating conversational interfaces which, though our user does not talk directly to the agents, was incredibly helpful in determining an effective way to measure the quality of our natural language generation.

[37] focused on the evaluation process for chatbots. Since their focus was on text-based conversational agents, it was especially relevant to our project. They surveyed quality attributes that were used to identify successful conversational agents, which included:

- Linguistic accuracy of outputs
- Convincing, satisfying, and natural interaction
- Provide greetings, convey personality
- Give conversational cues
- Provide emotional information through tone, inflection, and expressivity
- Make tasks more fun and interesting

Their survey also highlighted that most chatbot research was aimed at goal-based chatbots. In other words, chatbots that aimed to extract information from users in order to perform some task, for example asking for arrival and departure dates to help book a flight for a user.

[38] created an opponent with emotional cues in the form of facial expressions to try and include the social aspects of a game. For example, the opponent could smirk when in a strong position or when they appear to have a good idea of what to do next. They did not provide a clear explanation of how they evaluated their implementation however so we were unable to learn from it.

[39] provided the closest work to agent-agent interaction evaluation that we could find. They implemented a conversation centered drama that involved a human user, where agents would also interact in front of the human and express emotion and try emulate personalities. An interesting point raised by their study was how often users would rely on the scroll bar of the conversation

transcript to contextualise themselves during the drama. Usage of the scrollbar would suggest that the user was not able to keep track of the conversation, and was something we wanted to look out for when participants were given the opportunity to play our game.

With all the above considered, we created the questionnaire as described by table 4. 11 participants from the school of computer science helped by taking part in the feedback exercise. The exercise involved showing the participant three short videos of gameplay in different situations to try highlight features such as reactions and piece argumentation. They were then given a few minutes to play themselves. Finally, the user was asked to complete our questionnaire. Some wanted to continue playing after this stage to try and explore the range of conversation possible, and comments made from this period were noted for future improvement.

The results from the questionnaire were taken from paper to csv, and the frequency of each answer for each question plotted. Figures 17 to 19 show the results for each question in order.

Unfortunately the popularity of chess meant that only one of our participants had not played in either a very long time or never at all. The results mostly lean to the negative side, but only slightly, and there are some cases where participants were impressed.

A majority of players said they lost track of time while playing. Without being able to compare results with chess without dialogues, it is difficult to state the cause of this. It is possible that the conversations provide enough of a narrative to make players lose focus on time.

The game also ranked highly for the question "Playing makes me feel calm". This was expected, as gameplay is mostly reading the dialogue without much other stimulation.

In regards to the emotional response to the agents and the game, the results of the questions "I feel different" and "I get wound up" vary wildly. Players were definitely not 'wound up' or agitated by the game, which suggests that the player directed insults were ineffective or that the challenge was not sufficient for most players.

The linguistic accuracy of the pieces were rated highly. However, some of the freetext comments mentioned that the justifications all sounded the same, and another comment did suggest that the dialogue was not perfectly fluid. These issues could be resolved through careful tailoring of the BNF grammars used. The original interface for values only provided a single string representation for each of the opinions when evaluating the response to a move, and so this comment highlighted the need for the values to be included in the grammars to increase the number of wordings used.

The conversation planning seemed to be successful given the majority of players agreeing with the statement "The pieces had convincing and natural interactions". Given this was likely the part of the project that had taken the most time, it was rewarding to receive such a positive response.

Piece reactions did not seem to have much of an effect on players with the majority of participants giving a neutral response. One of the comments men-

tioned the reactions seeming random, which could point to issues with conveying personality of the agents, or a lack of emotion in the generated text. This is further evidenced by the replies to the final question, where most players did not seem to get a sense of personality from each of the pieces.

One of the reasons given for this in the comments were that maybe more exposure to the pieces were required to make the personalities recognisable, and another highlighted the problem that the similar appearance of the pieces made it difficult to track who said what and therefore remember which piece showed which personality traits. The nametags were implemented to avoid this issue, but clearly another solution is necessary to make each piece more unique.

Another interesting comment was that the phrases spoken by pieces would sometimes appear too long which ruined the flow of conversation slightly. The participant suggested breaking up longer statements by placing a limit on the length of each piece of dialogue, and splitting any statements that overflowed this limit. All the participants did enjoy the novelty of the concept however, which definitely suggested that these conversational agents would worth investigating further.

While continuing to play, one participant pointed out that it was hard to say if a piece had been 'convinced' to make a move, as the movement often seemed to come out of nowhere, even if it was the move currently being convinced. It is difficult to determine what change needed to be made to remedy this issue, but doing so would be likely to provide more satisfying gameplay. Further play also allowed some participants to begin to spot patterns and repetition in speech, which made it clear that the corpus of statements needed to be made larger.

Given one of the goals were to make the agent interactions humorous, participants were observed while playing to see if a smile or a laugh was given in reaction to agent behaviour. Fortunately, this did occur when agents caught the human player off guard by insulting their moves, or insulting each other.

5.4 Improvements

Due to time constraints, only some of the issues raised were able to be attacked before the final artefact had to be submitted.

The test coverage was greatly improved to try catch edge cases where the conversation planner would fail. The grammars used by each personality were also extended to try create more variation in the phrases produced during a game.

6 Future Work

The user feedback suggested that though this implementation was far from perfect, the conversational agents were an interesting novelty that is worth exploring further. Many improvements to our own work have already been discussed, but other extensions exist that would improve conversational agents.

With more time, a more thorough study that would compare player enjoyment with dialogue disabled and then enabled would shed more light on the benefits of the conversational agents. If we were able to make those who dislike chess or have not played before consider the game again, or provide a novel experience for those who enjoy the game already, then the project could definitely be deemed a success.

Other work that would definitely improve the current implementation would be accounting for piece value when choosing moves and using typical chess AI approaches such as minmax searches to make the chosen moves slightly more strategic.

The need to manually craft the grammars for agents is clearly a severely limiting factor for how an implementation like this could scale. Since any game wanting to include natural language generation will likely also require special domain knowledge to be accounted for, its unlikely a public database could be used. If the goal was to simply make conversational background characters for the sake of creating realism however, our approach could be useful.

Unfortunately, chat bots been used to perform social engineering for unethical reasons. Conversational chatbots could be used for propaganda purposes for emulating real people on social networks to try and fake popularity of certain topics or spread lies. [37] refer to one case where a chatbot was convincing enough for someone to think it was romantically interested in them. Agents that are able to create this level of realistic conversation could also be used to help identify when it is occurring, and simply researching how these agents could operate will help us identify patterns that such bots might follow.

7 Conclusion

By referring to the original project goals which - due to the wide range of topics being covered by this project - were definitely ambitious, it is clear that this project has been fairly successful. All of the primary goals set were met fully, all of the secondary goals were either partially or fully implemented, and the tertiary goal of a detailed graphical interface was attempted. User feedback definitely highlighted the need for improvements to the graphical interface for the game.

Researching this project was a fantastic opportunity to learn more about multiagent systems, natural language generation, and social artificial intelligence. Whilst progress definitely suffered due to the wrong approaches being taken early on in development for certain aspects of the if the project, the final product was able to include each of these topics to some degree.

Gathering user feedback was rewarding, and the constructive criticism given provided a clear guide of what work would be required to bring this project to its full potential, but I am nonetheless proud of what has been produced.

8 Acknowledgements

I would like to thank my supervisors: Dr Alice Toniolo for her expert advice and support throughout the project and helping set realistic but still ambitious goals at each step of the way, and Christopher Stone for providing excellent guidance to help me understand and implement the natural language generation approaches used in industry and those adopted within this project. Your shared enthusiasm for the resulting project was thoroughly appreciated.

I would also like to thank my fellow students for helping by taking part in the user feedback exercises, as well as providing great company during the development of the project.

Finally, I would also like to thank my family for getting me through the final stretch of work needed.

References

- [1] Thies Wittig, editor. *ARCHON: An Architecture for Multi-agent Systems*. Ellis Horwood, Upper Saddle River, NJ, USA, 1992.
- [2] Sergio Ocio and José Antonio López Brugos. Multi-agent systems and sandbox games. 2009.
- [3] Anonymous. Emergent tool use from multi-agent autocurricula. In *Submitted to International Conference on Learning Representations*, 2020. under review.
- [4] Stefan J. Johansson. On using multi-agent systems in playing board games. In *Proceedings of the Fifth International Joint Conference on Autonomous Agents and Multiagent Systems*, AAMAS '06, pages 569–576, New York, NY, USA, 2006. ACM.
- [5] Kalliopi Kravari and Nick Bassiliades. A survey of agent platforms. *Journal of Artificial Societies and Social Simulation*, 18(1):11, 2015.
- [6] Fabio Bellifemine, Federico Bergenti, Giovanni Caire, and Agostino Poggi. *Jade A Java Agent Development Framework*, volume 15, pages 125–147. 01 2005.
- [7] Rafael Bordini, Jomi Hbner, and Michael Wooldridge. *Programming Multi-Agent Systems in AgentSpeak Using Jason*, volume 8. 10 2007.
- [8] Derek Cormier, Gem Newman, Masayuki Nakane, James Everett Young, and Stephane Durocher. Would you do as a robot commands ? an obedience study for human-robot interaction. 2013.
- [9] Aike C Horstmann, Nikolai Bock, Eva Linhuber, Jessica M. Szczuka, Carolin Straßmann, and Nicole C. Krämer. Do a robots social skills and its objection discourage interactants from switching the robot off? In *PloS one*, 2018.

- [10] Arjan Egges, S Kshirsagar, and Nadia Thalmann. A model for personality and emotion simulation. pages 453–461, 01 2003.
- [11] Andry Chowanda, Peter Blanchfield, Martin Flintham, and Michel Valstar. Computational models of emotion, personality, and social relationships for interactions in games. 05 2016.
- [12] Cristiano Castelfranchi, Fiorella de Rosis, Rino Falcone, and Sebastiano Pizzutilo. Personality traits and social attitudes in multiagent cooperation. *Applied Artificial Intelligence*, 12:649–675, 1998.
- [13] Elizabeth Black and Katie Atkinson. Dialogues that account for different perspectives in collaborative argumentation. pages 867–874, 01 2009.
- [14] Katie Atkinson and Trevor Bench-Capon. Practical reasoning as presumptive argumentation using action based alternating transition systems. *Artificial Intelligence*, 171(10):855 – 874, 2007. Argumentation in Artificial Intelligence.
- [15] Gael Hette, Gauvain Bourgne, Nicolas Maudet, and Suzanne Pinson. Hypotheses refinement under topological communication constraints. 01 2007.
- [16] Henric Fransson. Agentchess an agent chess approach-can agents play chess ? 2003.
- [17] Ion Androutsopoulos, Gerasimos Lampouras, and Dimitrios Galanis. Generating natural language descriptions from owl ontologies: The natural owl system. *J. Artif. Int. Res.*, 48(1):671–715, October 2013.
- [18] Philipp Cimiano, Janna Lüker, David Nagel, and Christina Unger. Exploiting ontology lexica for generating natural language texts from RDF data. In *Proceedings of the 14th European Workshop on Natural Language Generation*, pages 10–19, Sofia, Bulgaria, August 2013. Association for Computational Linguistics.
- [19] Albert Gatt and Ehud Reiter. SimpleNLG: A realisation engine for practical applications. In *Proceedings of the 12th European Workshop on Natural Language Generation (ENLG 2009)*, pages 90–93, Athens, Greece, March 2009. Association for Computational Linguistics.
- [20] H. Cuayhuitl. Deep reinforcement learning for conversational robots playing games. In *2017 IEEE-RAS 17th International Conference on Humanoid Robotics (Humanoids)*, pages 771–776, Nov 2017.
- [21] Leire Ozaeta and Manuel Graa. *A View of the State of the Art of Dialogue Systems*, pages 706–715. 06 2018.
- [22] George Larionov, Zachary Kaden, Hima Dureddy, Gabriel Kalejaiye, Mihir Kale, Srividya Potharaju, Ankit Shah, and Alexander Rudnicky. Tartan: A retrieval-based socialbot powered by a dynamic finite-state machine architecture. 12 2018.

- [23] Will Knight. Computer crack funnier than many human jokes. <https://www.newscientist.com/article/dn1719-computer-crack-funnier-than-many-human-joke/>, 2001.
- [24] Graeme Ritchie, Ruli Manurung, Helen Pain, Annalu Waller, Rolf Black, and Dave Mara. A practical application of computational humour. *Proceedings of the 4th International Joint Workshop on Computational Creativity*, 01 2007.
- [25] I. M. Suslov. Computer model of a "sense of humour". i. general algorithm, 2007.
- [26] Abhinav Moudgil. Humor generation with recurrent neural networks. <https://amoudgl.github.io/blog/funnybot/>, 2017.
- [27] React a javascript library for building user interfaces.
- [28] Will B. chessboard.jsx. <https://github.com/willb335/chessboardjsx>, 2019.
- [29] Redux a predictable state container for js apps.
- [30] Jeff Hlywa. chess.js. <https://github.com/jhlywa/chess.js>, 2019.
- [31] FIPA. Fipa interaction protocol specifications. <http://www.fipa.org/repository/ips.php3> , 2002.
- [32] Ben-Hur Carlos Viera Langoni Junior. chesslib. <https://github.com/bhlangonijr/chesslib>, 2019.
- [33] @ebemunk. A visual look at 2 million chess games. <https://blog.ebemunk.com/a-visual-look-at-2-million-chess-games/>, 2016.
- [34] official stockfish. Stockfish. <https://github.com/official-stockfish/Stockfish>, 2019.
- [35] niklasf. python-chess. <https://github.com/niklasf/python-chess>, 2019.
- [36] Jeanne H. Brockmyer, Christine M. Fox, Kathleen A. Curtiss, Evan McBroome, Kimberly M. Burkhart, and Jacquelyn N. Pidruzny. The development of the game engagement questionnaire: A measure of engagement in video game-playing. *Journal of Experimental Social Psychology*, 45(4):624 – 634, 2009.
- [37] Nicole Radziwill and Morgan Benton. Evaluating quality of chatbots and intelligent conversational agents. 04 2017.
- [38] André Pereira, Carlos Martinho, Iolanda Leite, Rui Prada, and Ana Cristina Ramada Paiva. Designing a pervasive chess game. 2008.

- [39] Manish Mehta, Steven Dow, Michael Mateas, and Blair MacIntyre. Evaluating a conversation-centered interactive drama. In *Proceedings of the 6th International Joint Conference on Autonomous Agents and Multiagent Systems*, AAMAS 07, New York, NY, USA, 2007. Association for Computing Machinery.

Table 4: Questionnaire content with the question statement, the answer format (Likert meaning 5 point scale composed of strongly disagree, disagree, neutral, agree, strongly agree), and the justification for including the question in the questionnaire.

Statement	Format	Purpose
I have played chess before	Yes /No	May highlight the need for more visual cues to explain rules, or possible help identify if the game could be used to educate new players how to think tactically.
I am good at chess	Likert	Players who are already very familiar with chess might have different experiences to those who rarely play, and this may be observable from what we collect.
I lose track of time	Likert	Measuring presence, do players find themselves 'hooked'.
I feel different	Likert	Measuring absorption, are they emotionally affected by the anthropomorphic pieces.
The game feels real	Likert	Measuring flow, do players really feel like they are against an army of other pieces.
If someone talks to me, I don't hear them	Likert	Measuring flow, do players get involved enough to not be distracted by outside influences.
I get wound up	Likert	Measuring flow. Does the game frustrate the player, are they invested in it?
I lose track of where I am	Likert	Measuring presence. Is the player focused on the game?
I play without thinking about how to play	Likert	Measuring flow. The difference between experienced and non-experienced chess players may be reflected by this question also.
Playing makes me feel calm	Likert	Measuring presence. Do players enjoy their time.
I really get into the game	Likert	Measuring immersion. Do players really start considering the dialogue of the other pieces while playing.
I feel like I just can't stop playing	Likert	Measuring flow. Does the piece interaction interest players enough to make them keep playing.
The pieces spoke with linguistic accuracy	Likert	Did the conversation make sense and not contain any grammatical errors.
The pieces had convincing, natural interactions	Likert	Did the flow of conversation feel natural and emulate human discussion.
The reactions of pieces affected me	Likert	Did the dialogue create an emotional response, for example when a piece cried out at the loss of its friend did the player feel guilt.
The pieces convey personality	Likert	Is the player able to identify personality traits for each piece consistently.
The pieces provide emotional information through tone	Likert	Can the player tell if a piece has expressed anger, sadness, happiness through choice of words.
The dialogue made the game more fun and interesting	Likert	Did the dialogue add to the experience of playing chess, or just distract or detract from it.
Comments	Free text	Allow participants to highlight anything that was especially good or bad about the game.

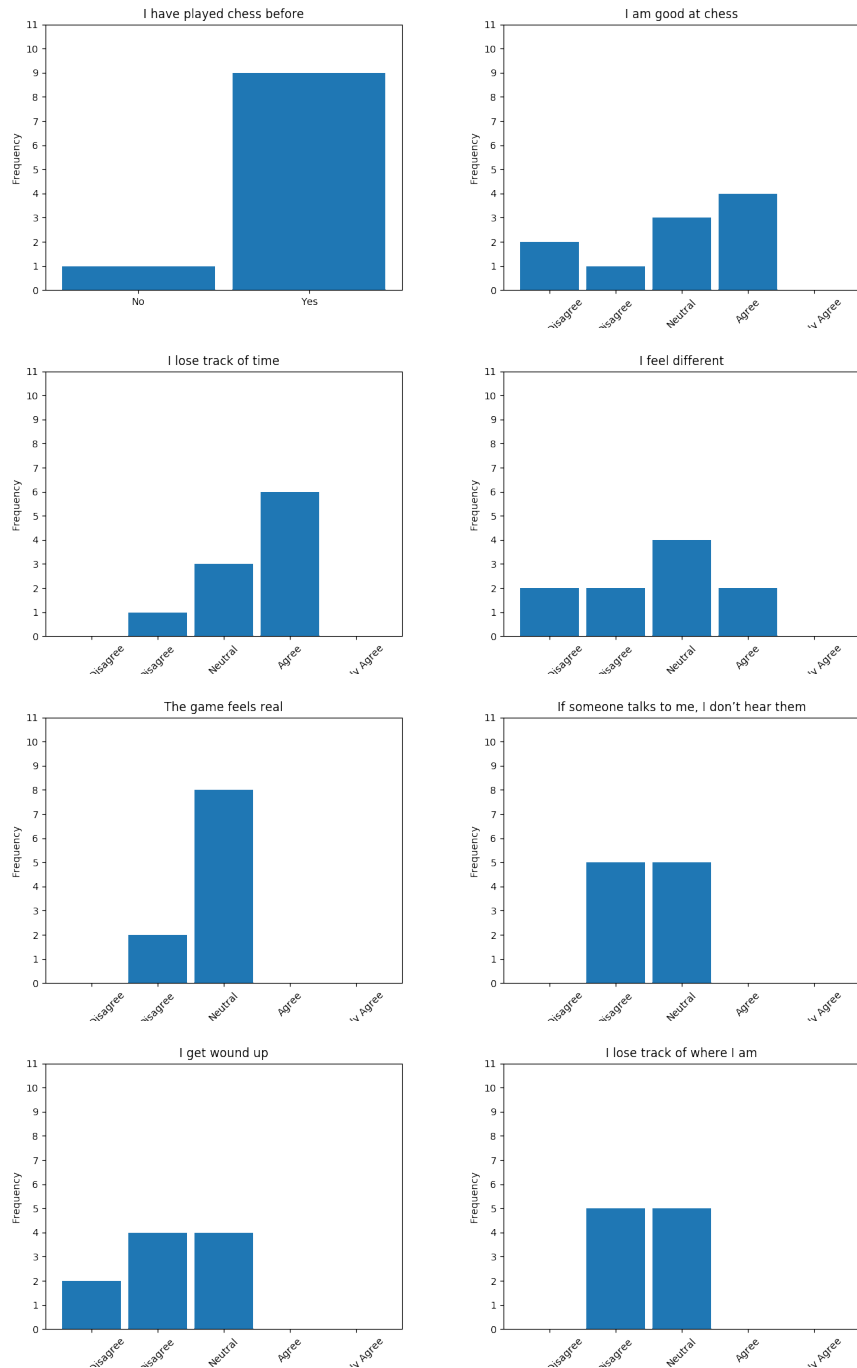


Figure 17: Questions 1-8 in Questionnaire

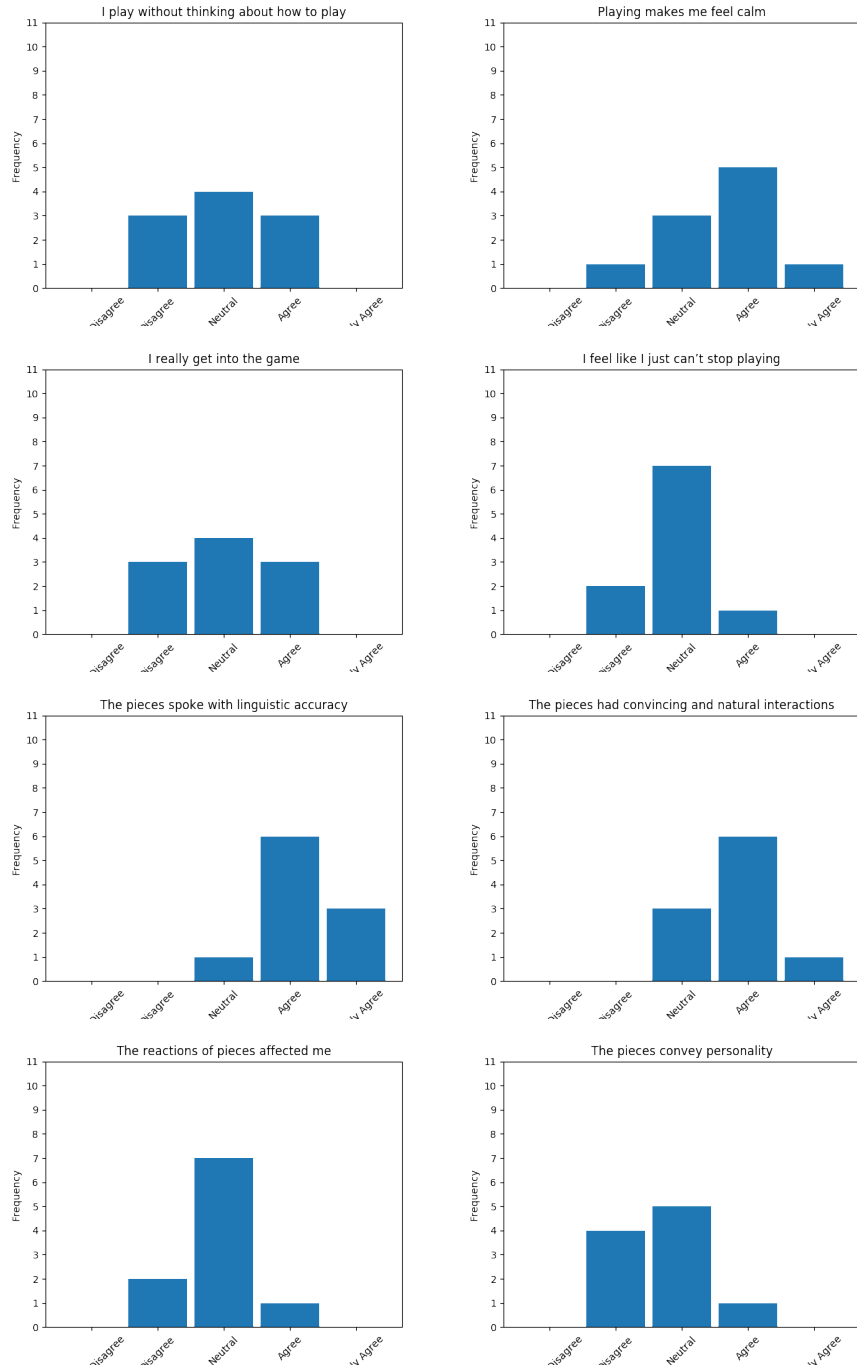


Figure 18: Questions 9-16 in Questionnaire

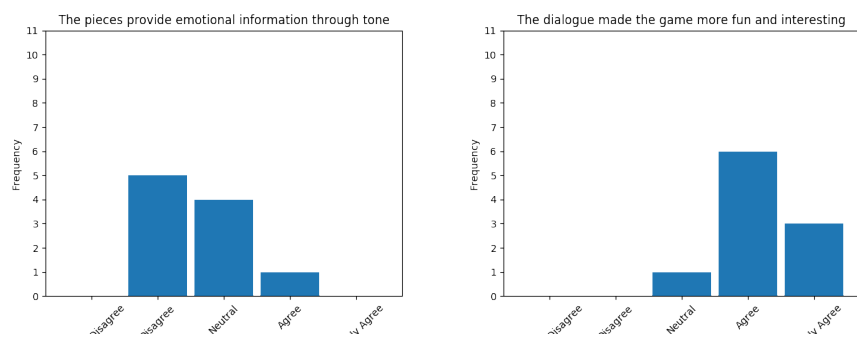


Figure 19: Questions 17 and 18 in Questionnaire