
Assignment 3: Modern Pac-Man: PCG, natural game interaction and RL

Course: Modern Game AI Algorithms, Master CS, Leiden University

Ivan Horokhovsky¹ Andreas Paraskeva¹ Pavlos Zakkas¹

Abstract

In this assignment, we aim to explore different player-game interactions, in the traditional game *Pac-Man*. Alternative means of engagement with the game include, voice and vision controls, in order to enhance the player experience whilst also providing accessibility for people with disabilities. Additionally, Procedural Content Generation (PCG) is applied to generate new *Pac-Man* mazes with the usage of Genetic Algorithms (GA). Finally, experimentation with Deep Reinforcement Learning (DRL) techniques was engaged, for the training of an agent that plays the game.

later stage with console controllers or keyboard and mouse. These are undoubtedly convenient methods of input that lots of users have become accustomed to, however the possibility of increasing the means of engagement can be beneficial both for the lifespan of the game as well as the end user experience and target demographic. Users that have some sort of disability can now interact with the game. This is why it was decided to experiment with audio and vision inputs, and translate these to actions in the game. The simplicity of the game, due to the small action space of just four potential moves, allows these inputs to retain a player experience with relatively low learning curve whilst being seemingly natural.

The next goal for this project was the implementation of PCG, in order to expand on the number of levels available, and provide alternating difficulty, through random generation of *Pac-Man* mazes of varying sizes. This was achieved through the usage of GA, which would converge to an individual, containing the maze, that satisfies the different restriction and characteristics of a *Pac-Man* maze.

Finally, even though the game itself possesses a small action space, the winning state is hard to be achieved. Therefore we decided to experiment with DRL approaches that aim to solve the game, through the training of an agent. Different methods for a reward-based approach and heuristics have been implemented in an attempt to create said agent.

1. Introduction

For the purposes of this assignment we have been requested to apply knowledge acquired through the course of "Modern Game AI Algorithms", complemented with external material and engage in a project of our choice.

The game of our choice is *Pac-Man*, since it provides a good balance between simplicity (having adequate branching factor) and variability of the game due to its non-trivial nature. The ability to provide extensions to such a popular game, which can enhance the player experience as well as possibly increase the replayability is intriguing.

It is important to note that the environment selected for this task was developed by UC Berkeley ([Berkeley](#)), and has been used for assessment of students for the implementation of model-based and model-free reinforcement learning algorithms. The project seemed to be outdated, since its development was in Python 2, but an adaptation to Python 3 is available on Github ([Spacco](#)).

Traditionally, interaction with this game has been achieved through the usage of an arcade console, and at a

2. Pac-Man environment

2.1. Game environment

Pac-Man ([PAC-MAN](#)) is considered an action maze video game. The player essentially controls this eponymous character and tries to traverse a maze, with the aim to eat all the dots placed in the path. The goal is essentially to eat all of the food (dots), whilst avoiding coming in contact with the ghosts. There are at most four ghosts, Blinky (red), Pinky (pink), Inky (cyan), and Clyde (orange) which pursue *Pac-Man*. Consumption of an energizer allows the character to pursue and eat the ghosts that are in this transitive state, in an attempt to boost the score. Each time frame essentially reduces the accumulated score, whilst the consumption of food increases this.

¹Leiden Institute of Advanced Computer Science, The Netherlands, The Netherlands.

On a superficial level, the game itself seems simple, but there are many aspects of that lead to its non-trivial nature and stochasticity. The ghost agents (which are automated and behave based on their move legalities) are developed in the Berkeley environment. Interacting and controlling the character utilizes the arrow keys of the keyboard, which are mapped to the respective moves (up, down, left, right).

An example of the environment and the initial state of the game is depicted in Figure 1, however the maze used can vary according to the selection by the user, which will consequently affect the starting positions of *Pac-Man* as well as the ghosts. The scoring will also be affected by the size of the Maze and the amount of food and energizers present.

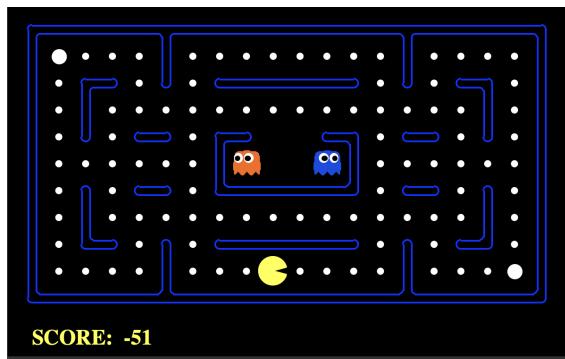


Figure 1. An example of a valid state of the *Pac-Man* environment.

2.2. Bots environment

The Berkeley implementation included two simple agents for *Pac-Man* which made the comprehension of the implementation easier. These agents included a *Left-turn agent*, which as suggested by the name will turn left at every available opportunity, i.e. when turning left is a legal action, and a *Greedy Agent*, which selects the best possible action to be taken at each state, investigating possible states of depth one. In the case that an optimal action (providing maximal reward) from the current state is not identified, random action selection is performed. These are very simple agents and have many flaws which we aim to eliminate with the creation of a DRL agent and consequently more advanced methods of action selection.

3. Background

In this section we will explore previous research or implementations, as well as the theoretical background related to our project and the respective aspects explored.

3.1. Player Experience Modeling

The 'effectiveness' or 'success' of a game is partly determined by the idea itself. A very important aspect is the replayability of a game and its ability to reccurrently encourage the player to revisit the game and/or promote extended sessions. An overview of the different aspects of game related to Artificial Intelligence (AI) and its relation to game development is depicted in Figure 2.

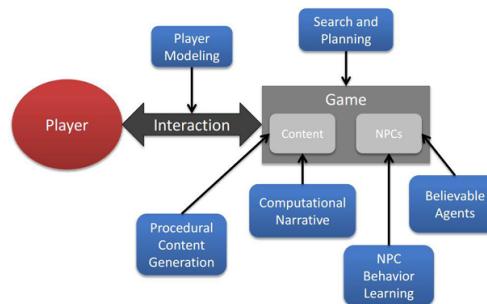


Figure 2. The panorama of AI research viewed from a player-game interaction perspective. (Yannakakis & Togelius, 2015)

Studies revolving around the theory of flow (Csikszentmihalyi, 1990), showed that there is an ideal threshold between the skill of a person and the difficulty of the task which can affect their flow, otherwise known as their ability to focus on the task at hand. This is also related to gaming, since the difficulty of the game itself as well as means of interaction (i.e. input) can affect the flow of the player. These means of input are often referred to as experience modeling. The game of *Pac-Man* is relatively simple in terms of controls, since there are only 4 possible directions that the player can perform. On the other hand, the difficulty of *Pac-Man* derives from the complexity of the maze and the movements of the ghost agents that are trying to eliminate the *Pac-Man* agent.

Based on the above, the learning curve for controlling the character through a keyboard, or controller is very minimal, whereas obtaining the ability to avoid ghost agents and consume all the food blocks of the maze requires more effort. As a result, providing a more natural form of engagement which does not require external devices handled by the player can enhance the experience. The challenge of this approach is to keep the control of the character relatively easy in order not to hinder the user with additional complexity, but let them focus on the actual difficulty of the game.

Providing methods of input such as voice (O'Donovan et al., 2009) and vision (Park et al., 2006) allows the player to interact with the game in a natural way, whilst not adding considerable effort in terms of getting used to the controls. This is mainly attributed to the small action space. The conclusion that both voice and vision controls have an in-

110 significant learning curve has been extracted through our
 111 own experience as well as a small-sized test group.

112 Another factor which substantially influenced our decision
 113 in providing more means of interaction with this game
 114 was the provision of accessibility for this game, and potentially
 115 allow users with disabilities to play *Pac-Man*.

117 3.2. Procedural Content Generation

119 In computing, generation of new data is essential, and ways to create/generate such data without the manual interference or effort of a human is very desirable. PCG has
 120 gained popularity with the exponential growth of AI as a field of study, since this can also be applied to games (see Figure 2). The combination of human-generated assets and
 121 algorithmic methodologies can allow exploitation of varying cases arising based on Random Number Generation (RNG),
 122 in order to generate new content.

123 Such a task is definitely not trivial, especially with highly complex games that need to follow certain restrictions and rules in order to retain validity and believability. A game
 124 like *Pac-Man*, despite its longevity and popularity, needs to provide new material to sustain its high engagement and players' interest.

125 For the prospect of generating new mazes for this game,
 126 extensive analysis of existing mazes and extraction of rules
 127 that define a valid maze was necessary. In the following list we can observe some of the extracted constraints and how
 128 these have been taken into account:

129 **Map Size** The original game consists of 28×36 tiles. Even though this is the size of the original mazes, it was
 130 decided that alternative sizes should also provide an enhanced player experience that can be designed for different target groups. For example, smaller mazes
 131 might be more appealing for a smaller challenge and new players, but at the same time can be appealing for players seeking high scores and competitions.

132 **Path Size** The paths that *Pac-Man* and the ghost can traverse through, should be restricted to a thickness of 1 tile.

133 **Wall Turns** The turns formulated through the placement of walls should not form a sharp turn (i.e. not exclusively joined in the diagonal formation)

134 **Dead ends** Paths should not lead to dead ends.

135 **Wall shapes** The walls should create shapes of a cross, l-shaped, t-shaped or straight lines of length 3. The existence of a rectangular wall would also be valid but should be limited in terms of occurrence.

3.3. Genetic Algorithms

A GA is an adaptive search method that relies on processes and techniques which are based on aspects of biological evolution. There are three main processes which have been used for the implementation of a genetic algorithm: *mating selection*, *crossover* and *mutation*. These traditional GA steps though can have various implementations which would consequently alter the effectiveness of our optimization problem greatly. Hence, experimentation with different versions for these steps and hyperparameter tuning, were essential for concluding on our optimal GA setup based on our findings.

As aforementioned, some of these operators were explored more thoroughly and experimentation with different implementation as well as research was essential. Post research of applied knowledge of GA approaches, assessment of potential implementation of GA in the project was formulated. Since the validity of a maze can be evaluated through the extracted constraints, it was deemed appropriate to have a population of individuals which consisted of the maze of *Pac-Man*. Any violation of the restriction can lead to an increase in the fitness, hence approaching this as a minimization problem is feasible.

Implementation details, and explanation of the algorithmic approach will be overviewed in Section 4.3.

3.4. Deep Reinforcement Learning

In this section we will talk about intellectual agents and our motivation behind picking deep reinforcement learning. For *Pac-Man* agents, we considered the following approaches:

- Heuristics based
 - Mini-Max
 - Alpha-Beta pruning
- Monte Carlo Tree Search
- Data based
 - Probabilistic Models
 - Classic Machine Learning
 - Supervised Deep Learning
 - Reinforcement Learning

The simplest and most intuitive approach would be to use heuristics based method. Implementing this method only requires a state search function (for instance BFS) and a state evaluation function (heuristics). The heuristics function can be the score of the game, but also the position evaluation can be added to avoid too greedy solutions. However, *Pac-Man* can be theoretically played very long (or even infinitely);

165 the computational resources are limited and long-running
166 episodes would lead to a great reduction in score (score negation
167 per time-frame). This is why the depth of the search
168 is usually limited. Nevertheless, in such cases, we might
169 spend our computational time on non-promising branches
170 and that is why several modifications like Mini-Max algo-
171 rithm or Alpha-Beta pruning could be added. Mini-max
172 algorithm provides an optimal move to a player assuming
173 that opponent (in our case ghost) is also playing optimally.
174 Alpha-Beta pruning optimizes the search by not expanding
175 non-promising branches. These techniques are pretty
176 simple, but probably they can solve the *Pac-Man* game.
177 However, if branching factor was high or reward gain was
178 too sparse these algorithms cease to perform well. Also,
179 our goal was to implement a more complicated and generic
180 agent, since creating heuristics requires in-depth domain
181 knowledge. Specifically, it needs manual and custom heuris-
182 tics metrics, which can in turn slow down the agent learning
183 procedure because it needs additional time to evaluate each
184 state. Furthermore, it is very inflexible to apply changes
185 when a more complicated logic needs to be implemented.

186 The next approach we considered is **Monte Carlo Tree**
187 **Search**, which partially solves the problems of heuristic
188 approaches and does not require human knowledge about
189 the game. Experimentation with this algorithmic approach
190 in a previous assignment, as well as basic understanding of
191 its concepts lead us to the conclusion that it cannot solve a
192 problem of a high branching factor.

193 Regarding data based approaches, the main problem is
194 probably the amount and quality of the dataset. These types
195 of agents need quite some time for training, but when play-
196 ing the game (evaluating) the computational complexity is
197 low. A simple probabilistic model would most likely work
198 worse than a **Machine Learning** (ML) one, while classic
199 ML algorithms are not very good in working with high-
200 dimensional data. Thus, we switched our focus to deep
201 learning approaches. For supervised learning a labelled
202 dataset is required, which seems to be a problem in our case,
203 since from a single state there can be multiple strategies that
204 the agent could follow, and not one absolute correct label.
205 In addition, a model can overfit or collected game states
206 can be skewed which will lead to poor generalization of the
207 agent. Because of the aforementioned issues it is not trivial
208 to get a proper dataset. **Reinforcement Learning** (RL) can
209 potentially formulate data on its own, learn how to tackle
210 new challenges would not require any human expertise on
211 the domain. The drawback though is that an RL approach
212 requires a lot of time and computational resources while
213 training.

214
215
216
217
218
219

4. Implementations

4.1. Vision assisted experience modeling

Following the intuition explained in Section 3.1, we decided to enhance further the user experience by enabling vision assisted interaction. In this mode, the user can move the *Pac-Man* agent through hand gestures. Depending on the movement of the hand, we are going to decide which of the actions *Left*, *Right*, *Up* and *Down*, should be applied to the agent.

4.1.1. RECORDING HAND MOVEMENTS

In order to record the movements of the user's hand we experimented with [OpenCV](#) and [mediapipe](#), which can be integrated to our game in order to provide the position of specific points of the hand, as they are identified through the user's camera. By taking into account these points, we can easily calculate the mean point of the user's hand. Then, by storing the history of the last 20 mean points with a time interval equal to 500ms, it is possible to determine how the hand position changes over time and understand in which direction it moves. To achieve that, we calculate the difference between pairs of consecutive points and average them across each dimension. Then by checking if the mean difference between consecutive points, either in the vertical or the horizontal axis, is larger than a *move threshold*, we can ignore slight hand movements and keep only those that correspond to actual gestures. This *move threshold* was set to 4 pixels.

As a result, if the hand position changes from left to right and exceeds the *move threshold*, we will apply the *Right* action to the agent. Moreover, since there are usually fluctuations in both vertical and horizontal axes we need to determine which change is larger. For example, when the user tries to move their hand from left to right, there might be a slight movement from up to down as well. Thus, we need to ignore this slight movement and apply only the *Right* action. For that purpose, we keep either the horizontal or the vertical movement based on which of those had a higher mean difference across the history of the recorded mean points.

The above method was a descent initial approach, but it could not provide a playable interaction. Specifically, the user had to perform each move exactly at the time when the corresponding action was legal. Otherwise, the action was considered invalid and thus it was not applied.

4.1.2. PLANNING NEXT MOVE

In order to solve the issue of requiring from the user to perform each action exactly at the time that it is valid, we decided to allow for planning the next move beforehand. For instance, if *Pac-Man* moves towards right and the user

wants to go up at the next cross section, they can perform the related upwards hand movement beforehand and the agent will apply it when it is legal for the first time, which corresponds to the aforementioned cross section.

As a consequence, the agent does not discard a received action if it is not applicable. Instead, the agent keeps it in memory as pending and executes it when it is valid. Meanwhile, if a new action is given by the user, the memory of the last recorded action is updated and the previous pending action is replaced with the new one. In this way, the game becomes much more playable, since the user is able not only to plan the next move but also undo it if needed.

4.1.3. RESETTING HAND TO THE CENTER

After implementing the aforementioned solution, another problem was revealed. To illustrate, when the user performs an action to the *Left*, their hand will end up on the left side of the screen. Then when they return their hand to the center area of the screen in order to perform the next actions, they should make unavoidably a hand movement towards the right. This movement though results in performing a *Right* action without being the intend of the user.

To avoid this issue, we need to allow the user to reset their hand back to the center area of the screen by ignoring the related hand movement. As a result, when a movement is determined, we consider it as valid if the last position of the hand is close to the center area of the screen. The *center area* is defined as a square window that deviates 10% from the central point of the screen.

The implementation of this feature made the control of agent smooth allowing the user to easily perform the desired actions. The intuition behind this idea can be related to the design of a joystick controller, in which the user takes an action towards a direction and then brings the controller back to the center before taking the next action.

4.1.4. ADJUSTING TO PLAYER'S DISTANCE

Finally, despite the fact that after using the aforementioned method the game became as playable as using the keyboard, we observed a last issue. The user's distance from the camera was affecting the interpretation of hand movements, since the applied *move threshold*, described in Section 4.1.1 was not accurate for every position of the user.

For this reason, we decided to adapt our implementation and replace this constant threshold with a dynamic threshold that depends on the depth (i.e. the distance of the user from the camera). Thus, we recorded the Z-axis values of the hand positions and based on the mean value we scaled the *move threshold* accordingly so as to allow a smaller threshold when the user stands far from the screen. In this way, the user will be able to apply an action by moving naturally

their hand independently of the distance they have from the camera. To formulate this scaling factor we measured that the maximum Z-axis value recorded when the user is very close to the screen is 0.4, whereas the minimum value is approximately 0.0001. By using these values, we can scale the move threshold of 4 pixels by the factor $\frac{z_{mean} - 0.0001}{0.4}$, which results in lower move threshold values when distance is higher.

4.2. VOICE ASSISTED EXPERIENCE MODELING

Another approach to enhance user experience is the voice assisted interaction. Specifically, the user will be able to control the character by using voice commands that will be executed by the agent.

We explored different speech to text APIs that could be integrated to our environment in order to transform voice commands to text in real time. Then by using the transcribed text, it will be possible to identify which command was given by the user. The allowed voice commands are *left*, *right*, *up* and *down*, which correspond to the possible moves of the *Pac-Man* agent.

4.2.1. SPEECHRECOGNITION: GOOGLE SPEECH API

Firstly, we experimented with the Google Speech API of the speech recognition package ([Uberi](#)). The transcription provided by this approach was accurate, but it lacked real-time processing. Specifically, even though the delay was relatively small, there were cases where more than half or one second was needed to transcribe the speech command to text. For our application, this delay has a significantly negative impact on user experience, since it does not allow the user to perform the desired actions on time, resulting in inefficient control of the agent. Even after trying to calibrate the parameters of the speech recognizer (i.e. the minimum audio energy to consider for recording, the seconds of non-speaking audio before a phrase is considered complete), it was not possible to achieve the desired interaction.

4.2.2. POCKETSPHINX SPEECH RECOGNITION

As a second approach, we experimented with PocketSphinx ([Bambocher](#)). Since, this API is supposed to transcribe speech in real-time, it seemed to be a promising solution for our problem. Indeed, the interpretation of voice commands was fast enough, but the transcription was not accurate. In order to solve this issue, we tried to change the voice commands to *west*, *east*, *north* and *south* respectively in order to avoid using the *up* command which has a very small length and can be easily misinterpreted. Also, we calculated the Levenshtein distance of the transcribed text from the aforementioned commands and when the distance was less than 1 or 2, we applied the respective command. Nevertheless, none of these approaches was able to solve

275 the issue of inaccurate speech to text transformations of the
 276 pocketsphinx API.
 277

278 4.2.3. ASSEMBLYAI SPEECH-TO-TEXT API

279 As a last approach, we decided to use the Speech-to-
 280 Text service provided by AssemblyAI ([Assemblyai](#)), since
 281 it offers real-time transcription with high accuracy. This
 282 method met our expectations, as it was able to transcribe
 283 accurately the voice commands with relatively low latency.
 284 Also, by slightly decreasing the frame speed of the game, we
 285 managed to provide a pleasant user experience, where the
 286 user could give voice commands that were executed on time.
 287 Moreover, as in the vision assisted interaction (see Section
 288 4.1.2), the user is able to give beforehand a command that
 289 will be applied when it is applicable for the first time. Thus,
 290 it is possible to plan the next move. Besides this feature
 291 though, we decided to enhance further the interaction by
 292 allowing the user to give two consecutive actions with one
 293 voice command. In the next section, we present the details
 294 of this feature.

295 4.2.4. PLANNING TWO MOVES AHEAD

296 In order to enable the user to plan the next two consecutive
 297 actions, we implemented a voice agent which is able to
 298 store a stack of two actions. Specifically, the user can give a
 299 combination of actions with a single voice command such
 300 as *up and right*. Then the agent keeps those two actions
 301 in memory and performs them sequentially when they are
 302 applicable. Thus in our example, with a single voice com-
 303 mand, the agent will go *Up* when this action is possible, and
 304 then it will go *Right* when applicable as well.

305 The impact of this feature was quite significant for the
 306 voice assisted interaction, because in cases where the user
 307 wanted to perform two successive actions, the controlling
 308 difficulty increased. For instance, in case of a 'U-turn'
 309 (around the wall), the user knows beforehand which actions
 310 should be done, but if the commands are given one by one,
 311 the user should act quite fast when giving the second com-
 312 mand in order not to miss the 'U-turn'. On the other hand,
 313 when the user is able to plan the next two moves, the con-
 314 trol of the agent becomes smoother and more accurate as
 315 presented in Section 6.2.

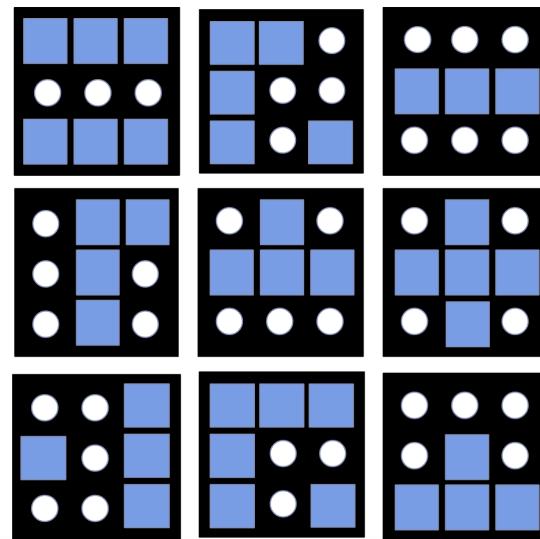
316 4.3. PCG with GA

317 The PCG for the generation of different *Pac-Man* mazes
 318 has been implemented with the usage of GA. In the sections
 319 that follow we provide a brief overview of the individual
 320 structures as well as the algorithmic approach followed in
 321 order to generate said mazes.

322 4.3.1. POPULATION

323 The number of individuals in each generation is denoted
 324 as μ , while the number of offspring that are generated during
 325 each iteration is represented by λ .

326 Each individual consists of half of the maze, which can
 327 follow a size of $n \times m$, where n and m are multiples of 3.
 328 Each individual is initialized through stochastic selection
 329 from the 3×3 building blocks shown in Figure 3. The blocks
 330 shown in this figure, were then rotated by $[90^\circ, 180^\circ, 270^\circ]$
 331 and flipped across the vertical axis in order to create all
 332 possible combinations of valid blocks. These were 35 dis-
 333 crete blocks in total. It should also be noted that the next
 334 generation of each population is selected in a deterministic
 335 way, in which the best μ individuals are selected either out
 336 of the λ offspring (μ, λ) or out of the μ parents and λ off-
 337 spring ($\mu + \lambda$). Both selection methods have been tested to
 338 conclude that the best one was deemed to be $+$ -selection.



339 *Figure 3. The 3×3 building blocks used for the individuals.*

340 4.3.2. FITNESS

341 For the evaluation of the individual we have used a num-
 342 ber of different 2×2 , 3×3 , 3×2 , 2×3 and 4×1 blocks,
 343 which contain patterns that would violate the constraints of a
 344 *Pac-Man* maze. Similarly to the methodology for acquiring
 345 all the possible unique building blocks, we have created a
 346 list of the invalid blocks and rotated and flipped them to get
 347 all possible combinations. This logic is applied to 2×2
 348 and 3×3 invalid blocks (see Appendix A) but for the 4×1
 349 which essentially is just four walls in a row, we only needed
 350 to transpose it in order to check both column-wise and row-
 351 wise for the existence of said pattern. Regarding the 3×2
 352 and 2×3 invalid blocks, we did not need to perform any
 353 transpositions since all versions are presented in 15, and

330 they are just four in total.

331

332 4.3.3. ALGORITHMIC APPROACH

333 Our algorithmic approach is based on the following
 334 sequence of iterative and repeating processes of evolution.
 335

336

337 1. Mating Selection:

338 **Proportional Selection:** We used fitness proportion-
 339 ate selection, the fitness is associated with each
 340 respective individual and the probability of selec-
 341 tion will be correlated to its fitness. To further
 342 improve the probability of more fitted individuals
 343 (parents) to be selected for the mating process,
 344 we used scaling. This essentially will boost the
 345 probability of selecting parents with higher than
 346 average fitness and those with lower than average
 347 fitness will be less likely to be selected. The for-
 348 mula for the probability assignment is as follows:
 349

$$350 p'_i = \frac{f_i - c}{S_f - c * \mu}$$

351 where $S_f = \sum_{j=1}^{\mu} f_j$ and c = the minimal fitness
 352 value of the selection pool.

353 *Hyperparameters:*

- 354 (a) factor: determines how many children will be
 355 produced, based on the initial population size.
 356 #of children = factor * population size

357 **Tournament Selection:** This consists of two main
 358 steps that are repeated for each parent selection.
 359 Firstly, we select a random subset of the popu-
 360 lation (of size k, determined by our hyperparam-
 361 eter). Following, we will select the best fitted
 362 individuals out of the subset, based on their fit-
 363 ness values. *Hyperparameters:*

- 364 (a) k: is the number of randomly selected individ-
 365 uals.
 366 (b) factor: determines how many children will be
 367 produced, based on the initial population size.

368 2. **Crossover:** After the selection of pairs from the popu-
 369 lation of parents is finalized, then we proceed with the
 370 crossover.

371 **Uniform crossover:** The two offspring are initially
 372 identical copies of the two parents. We then pro-
 373 ceed to traverse one by one their building blocks
 374 which form the maze, and for each one a proba-
 375 bility is used to determine whether the currently
 376 visited block is switched with the one of the other
 377 offspring. *Hyperparameters:*

- 378 (a) pc: the probability of the crossover taking
 379 place. typically assigned to 0.5 for fairness.

380 3. Mutation:

381 For mutation, we introduce the concept of block fitness
 382 alongside the individual fitness. This is in an attempt to
 383 prioritize blocks which have the most violations. More
 384 in depth explanation is provided later in this section.
 385 *Hyperparameters:*

- 386 (a) p_m : the partial constant used to determine proba-
 387 bility of the mutation taking place.
 388 (b) s_{pm} : small value to avoid zero probability of mu-
 389 tation. Tuned to optimize performance.

390 An algorithmic overview for the general procedure fol-
 391 lowed can be depicted in the algorithm pseudocode 1.

392 Mutation and Block Fitness

393 In an attempt to optimize the mutation phase, we intro-
 394 duce block fitness. As stated earlier, when an individual
 395 is evaluated with the set of invalid blocks, we increment
 396 its fitness by 1 every time that a violation occurs, during
 397 out iterative check of whole maze (half-maze). However a
 398 violation can occur in more than one of the building blocks
 399 of the individual. Some examples can be observed in Figure
 400 4, where the red outline indicates the area that the viola-
 401 tion was identified in. Therefore, whenever a violation is
 402 identified, the fitness that will be added to the individual is
 403 equally divided in all of the building blocks that overlap in
 404 the violation.

405 With this technique, during the mutation phase, we can
 406 multiply the division of the individual block fitness and
 407 individual fitness by the chance of mutation, in order to
 408 create proportional mutation probabilities, where the block
 409 with the most violations (i.e. higher block fitness) will have
 410 a higher chance of being mutated. Since we also desire a
 411 probabilistic mutation even for more fit blocks the formula
 412 also includes a parameter (s_{pm}) which is added to avoid
 413 zero probability of mutation for a block. We generate a
 414 random value between zero and one and if this is smaller
 415 than our retrieved calculation we perform mutation. The
 416 finalized formula is as follows:

$$417 p_m * \frac{\text{individual.block_fitness} + s_{pm}}{\text{individual.fitness}}$$

418 4.4. DRL Agents

419 The RL field describes the use of the active interactions of
 420 an agent with its environment to learn the optimal behavior
 421 or achieve a specific goal in that environment. RL problems
 422 in general can be formulated as a **Markov Decision Process**,
 423 a tuple $\langle S, A, T, R, \gamma \rangle$ where S is the set of states s in
 424 the environment, A the set of allowed actions a , $T : S \times A \times S' \rightarrow [0, 1]$ is the transition function, $R : S \times A \times S' \rightarrow R$ is
 425 the reward function and γ is the discount factor (Plaat, 2022).

Algorithm 1: Overview of the Genetic Algorithm used for the PCG of *Pac-Man* mazes.

```

Input :Population size  $\mu$ 
          Offspring size  $\lambda$ 
          Next generation selection [',', '+']
          Crossover ['Uniform']
          Mutation type ['Block Fitness Proportional']

Termination :The algorithm terminates when: - the target value has been reached

p  $\leftarrow$  Population( $\mu$ );
parent_evals  $\leftarrow$  Evaluation(f, p);
while fitness_opt > 0 do
    p'  $\leftarrow$  MatingSelection(p);
    p''  $\leftarrow$  Crossover(p') ;                                // Probabilistic creation of offspring with parents' blocks
    offspring_evals  $\leftarrow$  Evaluation(f, p'');
    p'''  $\leftarrow$  Mutation(p'');
    offspring_evals  $\leftarrow$  Evaluation(f, p''' );
    /* Perform '+' selection                                         */
    p += p''' ;
    evals  $\leftarrow$  parent_evals + offspring_evals;
    p, parent_evals  $\leftarrow$  Selection(p,  $\lambda$ , evals) ;           // Get  $\lambda$  best individuals and their fitnesses
end
final_grid  $\leftarrow$  best_individual.mirror_pad() ;           // Mirror across vertical axis. Pad with food and outer wall.

```

Algorithm 2: Actor-Critic Algorithm

Input : a differentiable policy π_θ , learning rate η , update frequency f , exploration rate σ , discount fraction γ and annealing rate ϵ , bootstrap depth n .

Initialize the actor and critic NNs with random weights θ and ϕ .

Set **actor loss** and **critic loss** to dequeues with size f ;

for trajectory *in* budget **do**

- Sample trace $h = \{s_0, a_0, r_0, \dots, s_{n+1}\}$ with π_θ ;
- $R, A \leftarrow 0$;
- $L_{\text{actor}} \leftarrow 0$;
- for** $t \in n, \dots, 1, 0$ **do**

 - $R \leftarrow \gamma^n \cdot V_\phi(s_{t+n}) + \sum_{k=0}^{n-1} \gamma^k \cdot r_{t+k}$;
 - $A \leftarrow R - V_\phi(s_t)$;
 - $L_{\text{actor}} \stackrel{+}{=} A \cdot \log \pi_\theta(a_t | s_t)$;

- end**
- Append L_{actor} to **actor loss**;
- $L_{\text{critic}} \leftarrow \text{MSE}(R(s_t, a_t), V_\phi(s_t))$;
- Append L_{critic} to **critic loss**;
- if** trajectory % $f = 0$ **then**

 - $L_{\text{actor}} \leftarrow \text{mean}(\text{actor loss})$;
 - Update the actor NN with L_{actor} ;
 - $L_{\text{critic}} \leftarrow \text{mean}(\text{critic loss})$;
 - Update the critic NN with L_{critic} ;
 - $\sigma \leftarrow \epsilon \cdot \sigma$;

- end**

end

436
437
438
439

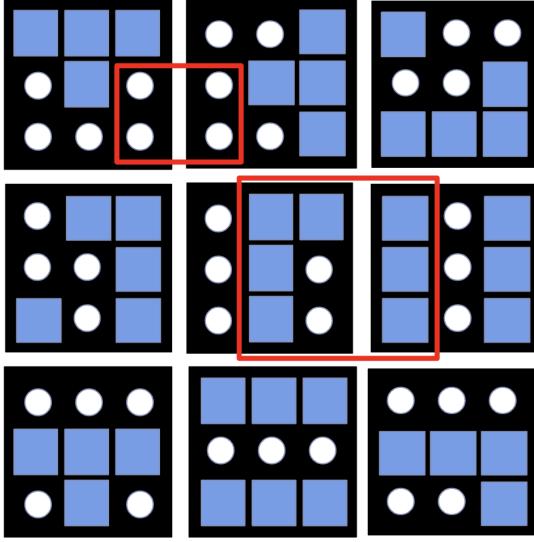


Figure 4. Evaluation of individual and block fitness example. The red outline indicates the area that the (sample of) violations was identified in. Visualized can also be found some of the building blocks of the individual.

We seek an optimal policy for the agent to follow, defined as $\pi^* : S \rightarrow A$ that maximizes the expected discounted sum of the future rewards $Q(s, a) = E_\pi[\sum_i \gamma^i r_{t+i+1} | s_t = s, a_t = a]$.

Simpler Deep RL problems that do not have continuous action spaces can be solved with value-based approaches where a Neural Network (NN) is trained to learn the value-function that maps the state-action pairs to their respective Q-value. On the other hand, policy-based methods, which can be applied in both continuous and discrete action spaces, aim to learn directly a stochastic policy, which introduces natural exploration of states. We decided to use the Actor Critic (AC) modification, which is expected to converge faster than value-based Q-learning, as it offers a good balance between bias and variance by using a value head along with the policy head (Sutton & Barto, 2018). The pseudocode is presented in Algorithm box 2.

To make our agent work we needed to adjust our *Pac-Man* environment. The model perceives the game state as a 3D matrix: $width \times height \times cell\ value$ which is a one hot encoded vector of 5 values. Specifically, each value represents the positions of the 5 different entities of our game, including *Pac-Man*, ghosts, walls, food, and energizers. The absence of all those entities means that the block is empty. The aforementioned input representation is also visible in the input state shown in Figure 5.

During the first iteration, we used a single random ghost agent to play against our *Pac-Man* agent during training and we used the score of the game to identify the reward of each

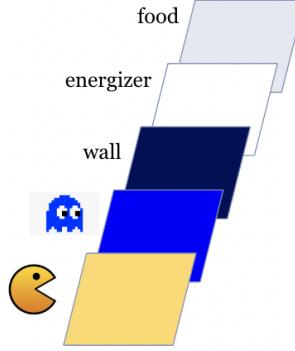


Figure 5. State representation of *Pac-Man* environment

action. To be more specific, if the agent is at state s and performs an action, the game environment moves to a next state s' , by moving the *Pac-Man* agent and the ghost agent accordingly. The score difference of the two states can be used as the reward returned by the environment for action a .

This vanilla method did not yield the expected results, as the agent could not usually perform legal actions and was eaten by the ghost very soon, resulting in very short episodes that were unable to train the agent. Based on this observation, we decided to remove the ghost agent in order to let *Pac-Man* solve the maze without any threat. Despite this, it failed to perform only legal actions.

To deal with this issue, we performed manual adjustment of the default reward mechanism, which was previously developed solely from the game score. At first, we penalized illegal actions significantly. The impact of this approach was quite positive since the agent learnt the game rules and started performing legal moves in most of the cases. However, it could not escape specific states, since after consuming the first 'food' particles, it entered a looping motion of backward and forward moves in a specific section of the maze. This behaviour is indicative of the agent being stuck in a local minima. Thus, we decided to modify the interaction with the environment by adjusting the rewards further.

In order to create such hand-crafted rewards, we experimented with the following events. The positive reward events are: win, eating food or energizer, eating ghost, while the negative rewards are: illegal move, performing a step without eating food, step backwards, game loss. Different weight values were assigned to each of these events, and the final action reward was the summation of all the rewards of each specific event. Thus, if the agent performs a move that eats the ghost but does not eat any food, we sum up the positive reward of eating the ghost and the negative reward of doing a step without consuming food. Of course, in this case, the weight of the reward of eating a ghost should

495 be quite higher than the one of not eating a 'food' particle
 496 (which also aligns with the scoring mechanism of the game).
 497 Moreover, in most of the cases, there is no point for an agent
 498 to make a step backwards and this is why we penalize the
 499 agent for this action slightly more, in an attempt to restrict
 500 backward moves to the absolute necessary cases. This was
 501 the main approach we used in order to avoid local minima
 502 and train our agent explore different parts of the maze with-
 503 out being stuck in specific areas whilst performing a looping
 504 motion.

505 Furthermore, we tried to increase exploration using a Soft
 506 Actor Critic (SAC) (Haarnoja et al., 2018) agent, which
 507 aims to avoid the collapse of the policy distribution. Specifi-
 508 cally, despite the fact that a policy-based method includes
 509 exploration by nature as it learns a stochastic policy, it is a
 510 common approach to enable the exploration of more states
 511 further, in case a problem is not solvable with the default Ac-
 512 tor Critic method. As a result, SAC introduced the entropy
 513 regularization method, which uses an additional penalty
 514 term to the loss function, which is the entropy of the actions
 515 distribution predicted by the policy model. In this way, SAC
 516 enforces entropy to stay larger, and prevents policy from
 517 being too narrow and selecting the same actions back to
 518 back.
 519

520 Finally, different network architectures were considered.
 521 The first one is 2D convolutions, in which we start with
 522 5 channels in the beginning, which represent each cell of
 523 the state. The motivation behind this network is that we
 524 have an image-like input with which Convolutional Neural
 525 Networks (CNN) and Vision Transformers (ViT) (Dosovit-
 526 skiy et al., 2020) are considered as state of the art. However,
 527 training a ViT is not a trivial task especially for RL where
 528 a lot of hyperparameters need to be tuned. The next archi-
 529 tectural approach was formulated based on the assumption
 530 that we should not mix the state cells with CNN channels,
 531 as they would share the same NN weights. This approach
 532 does not seem to be intuitive and this is why we needed to
 533 add another dimension and end up using 3D convolutions.
 534 Another possible extension would be utilizing Long-Short-
 535 Term-Memory (LSTM) to take advantage of the history of
 536 states. Further explanation is provided in Section 7.2.
 537
 538
 539
 540
 541
 542
 543
 544
 545
 546
 547
 548
 549

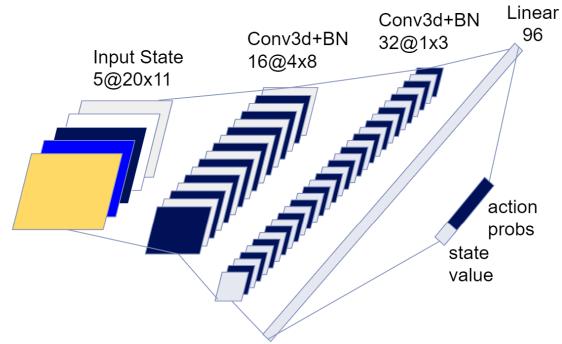


Figure 6. Policy based, actor-critic RL agent network architecture. Two 3D convolutions with batch normalization. Activation functions ReLU and Softmax on final layer.

5. Tuning

The free parameters that have been introduced by the algorithms explained in the previous section arose the need of performing thoughtful hyperparameter optimization. This section summarizes the experiments conducted to achieve optimal results.

5.1. Tuning: Vision assisted experience modeling

For vision assisted interaction, we had to experiment with various configuration in order to end up with a playable solution.

At first, we tried different time intervals between consecutive position recordings of the user's hand. Gathering many hand positions within a small interval requires a large amount of points to be stored in agent's memory in order to be able to identify a possible movement. On the other hand, sampling hand positions more rarely may result in missing important information related to the user's movements. Thus, sampling a hand recording every 500ms and keeping a history of the last 20 hand positions achieved the best results compared to other configurations.

Secondly, one of the most important parameters was the *move threshold* described in Section 4.1. Large threshold values required larger hand movements from the user while small threshold values resulted in undesired actions due to high sensitivity. It was observed though, that a 4 pixels *move threshold* is an optimal value for a user that stands within a mid-range distance from the screen. Nevertheless, applying the same threshold independently of the user's distance was not optimal for all cases, and that is why we needed to multiply the default 4 pixels threshold with the scaling factor $\frac{z_{mean} - z_{min}}{z_{max}}$, where $z_{min} = 0.0001$ and $z_{max} = 0.4$ as measured in Z-axis.

Last but not least, the window of the *center area* was

550 also tuned in order to allow the user return their hand to
551 the center and then perform their next action. This area is a
552 square window where each vertex lies at a specific distance
553 from the center point of the screen. This distance is defined
554 as a percentage of the screen's resolution. For large values,
555 the user should perform larger hand movements to escape
556 the *central area* and take an action, but for smaller values,
557 slight movements could result in unwanted actions, when
558 resetting the hand to the center. After experimentation, a
559 distance from the center that is equal to the 10% of the
560 screen's resolution seems optimal for our case to avoid the
561 aforementioned issues.

562 **5.2. Tuning: Voice assisted experience modeling**

563 As mentioned in Section 4.2, after experimenting with
564 different speech-to-text services, we ended up using AssemblyAI
565 service which did not require any tuning process since
566 it met our expectations.

567 **5.3. Tuning: PCG with GA**

568 Since there are no concrete ways to evaluate the produced
569 matrix, aside from validity based on the constraints that we
570 introduced in Section 3.2, our only way to evaluate the
571 performance of each run was to repeat 5 runs with the current
572 settings under evaluation. Following this, we averaged the
573 number of generations needed to generate the maze, and this
574 was used as one of the metrics for performance evaluation
575 of each hyperparameter setting. Another important metric
576 was the average time needed to generate a maze, since the
577 time for each generation is dependent on many of our hyper-
578 parameters.

579 **Next Generation Selection**

580 The relevant hyperparameters or tunable methodologies
581 have been explained in the Section 4.3. Initially one of
582 the most important choice for a GA approach would be the
583 decision on the type of selection applied for the new genera-
584 tion of individuals. As aforementioned the next generation
585 of each population is selected in a deterministic way, in
586 which the best μ individuals are selected either out of the
587 λ offsprings (μ, λ) or out of the μ parents and λ offsprings
588 ($\mu + \lambda$). After experimentation with both methods, it was
589 concluded that $(\mu + \lambda)$ has yielded better performance and
590 more consistent convergence within a reasonable number of
591 generations.

592 **Population Size**

593 The population size was also a crucial hyperparameter since
594 a value that was too small would hinder the ability to
595 converge to a solution, whilst a higher than optimal set-
596 ting would lead to more time needed for each generation.
597 Through our testing the optimal value was deemed to be 8
598 individuals but this is very sensitive and can change with
599 the combination of the other hyperparameter settings. How-
600 ever we do not expect it deviate greatly from this identified

601 optimal.

602 **Mating Selection**

603 Our mating selection phase could utilize two different meth-
604 ods. During our experimentation we determined that both
605 methods were determined to be equally suitable. These
606 included proportional selection and tournament selection.
607 Additionally these methods both include the hyperparameter
608 referred to as *factor*, which determined the number of chil-
609 dren that would be produced based on the population size.
610 We concluded that a factor of 1.0, meaning that we produce
611 the same number of offspring as the initial population size.
612 Having a smaller factor was not favorable, whilst a higher
613 than 1.0 factor did not seem reasonable so no experimen-
614 tation was performed. The tournament selection method
615 introduced the k hyperparameter which is the number of
616 randomly selected individuals for this specific method. A
617 value of 2 – 3 was concluded to be optimal since higher
618 values lead to a decrease in relative performance.

619 **Uniform Crossover**

620 The tunable hyperparameter relevant to this was p_c . Taking
621 into consideration the usage of +selection, thus keeping
622 the previous generation of parents in the pool of selection,
623 it will be favorable to have a relatively higher chance of
624 altering the offspring. Thus we experimented with values
625 for p_c between 0.5 and 1.0, where 1.0 leads to the crossover
626 happening everytime. We concluded that an optimal setting
627 based on our search space would be 0.7.

628 **Mutation**

629 Due to the customized approach for mutation and the need
630 for block fitness as well as individual fitness, the typical
631 value of p_m would not suffice. We definitely need to take
632 this into consideration since it will determine the probabili-
633 ty of a mutation taking, proportionally to the block fitness.
634 Thus a higher than typical value was necessary. Through
635 our experimentation with a vast number of values, rang-
636 ing from 0.05 to 0.7 we concluded to the value of 0.5 as
637 a good performing setting. Additionally, we introduce the
638 hyperparameter of s_{p_m} which is a relatively small value to
639 be added to each block fitness to avoid zero probability of
640 mutating, even good blocks. Through careful considera-
641 tion and limited testing in values, the value of 0.05 is reasonable
642 to provide a small probability of mutating well fitted blocks
643 in the individuals.

644 It is worth noting, that our concluded values for the opti-
645 mal settings were retrieved through our search space which
646 is no way compensating for all the possible combinations
647 of these settings. The concluded settings are performing
648 well but further tuning will possibly result in a higher per-
649 formance (with the aforementioned metrics considered).

650 **5.4. Tuning: RL Pac-Man agent**

651 Given the nature of NNs, **hyperparameter optimization**
652 (HPO) has substantial impact on the performance of ma-

chine learning and deep learning algorithms, especially on RL domain, given the amount of additional hyperparameters that are required to define the training and validation of the RL agent. This leads to an inherently large search space that renders it difficult to fully explore such a big number of configurations. For this problem, three main HPO methods were considered namely: grid search, random search and Bayesian optimization. Random search is more efficient and explores better than grid search, and theoretically converges to the optimal solution. Bayesian optimization, cannot be parallelized efficiently and possesses additional hyperparameters (such as the acquisition function and kernel of the Gaussian process), which need to be tuned separately. Random search is a good balance in simplicity and efficiency so it was picked for this problem.

Hyperparameters can be split in NN and RL categories:

- Neural Network hyperparameters
 - network type (see 4.4)
 - number of convolutional layers
 - number of fully connected layers
 - usage of batch norm regularization
 - usage of (max) pooling
 - convolution layers hyperparameters
 - activation function
 - weights precision (can speed up training process, but make the performance a bit lower)
 - network optimizer
 - network loss function
- Reinforcement Learning hyperparameters
 - gamma
 - step size
 - use of baseline subtraction
 - entropy regularization (SAC)
 - budget
 - rewards list

For hyperparameter tuning we created an automated framework which uses a configuration file (in JSON format) with hyperparameters. In the case of continuous hyperparameters, these are accompanied by distribution and ranges. For categorical parameters, we provide a list of values to be tested. The program samples a predefined amount of hyperparameter sets which are used for the training and evaluation processes. To make the results more statistically significant, each experiment needs to be executed several times, which is also defined in the configuration file. Last but not least, the program tries to efficiently utilize the resources of available, by running many processes concurrently.

6. Results

Our experimentation phase was essential in order to evaluate our implementation methods explained in Section 4. Different sections required different degree and type of experimental procedure.

6.1. Results: Vision assisted experience modeling

After following the implementation described in Section 4.1, it is possible to record hand movements as shown in Figure 7 and eventually control *Pac-Man* actions. For further interpretation of the results, a video (<https://youtu.be/w2EI2vXkWhot>) was recorded demonstrating the whole user experience when using the vision assisted interaction.

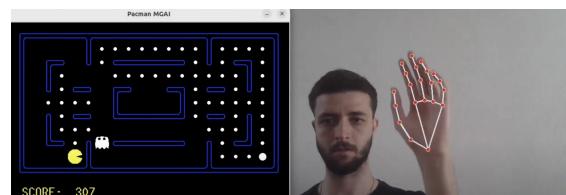


Figure 7. Vision assisted agent

6.2. Results: Voice assisted experience modeling

As described in Section 4.2, the voice assisted interaction was implemented, where the user can give voice commands to the *Pac-Man* agent and even perform two consecutive actions with a single command, as shown in the detected text in Figure 8. A demonstration of such a game is presented in the following recorded video: <https://youtu.be/nB9Nk-Uflis>

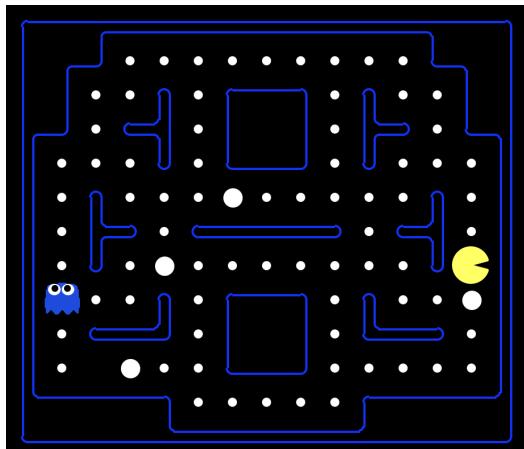


Figure 8. Voice assisted agent

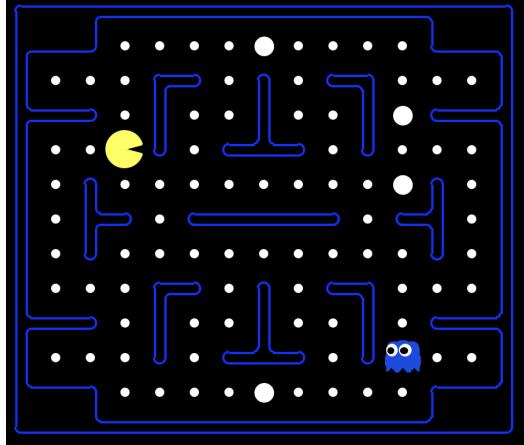
660
661
662
663
664
665
666
667
668
669
6.3. Results: PCG with GA

660
661
662
663
664
665
666
667
668
669
In an attempt to evaluate as well visualize the results of the PCG, we used the Berkeley *Pac-Man* project with the generated layouts. It was important to test different sizes of mazes, in order to ensure that our restrictions are met throughout, as well as manage to create different randomly generated mazes. These different sizes can provide further replayability as well as varying degrees of challenge to improve the flow of the game.

670
671
Some of examples of the retrieved mazes are shown in Figures 9, 10, 11.
672



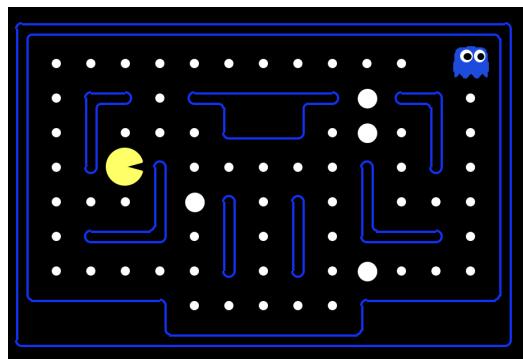
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
(a) An example of 13×15 generated maze



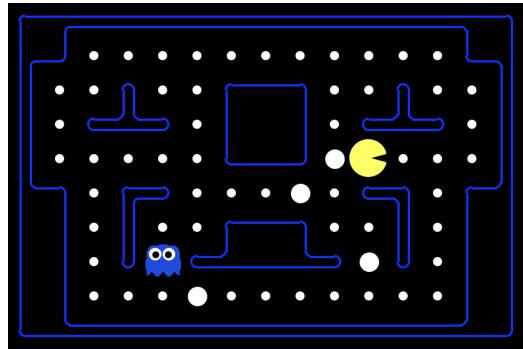
703
704
(b) An example of 13×15 generated maze

705
706
707
708
709
710
711
712
713
714
Figure 9. Some of the generated mazes are shown to present the validity and varying sizes capabilities of the algorithm

714
Closer inspection of these mazes leads to the conclusions that the generation of the mazes with the GA approaches is successful and satisfies the constraints that were listed earlier. However, some minor problems can arise when mirroring the maze. There is a possibility of creating some shapes that are not within the valid ones. Restrictions re-

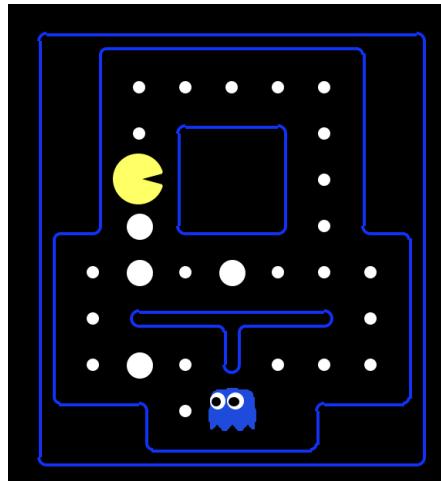


673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
(a) An example of 10×15 generated maze



703
704
705
706
707
708
709
710
711
712
713
714
(b) An example of 10×15 generated maze

705
706
707
708
709
710
711
712
713
714
Figure 10. Some of the generated mazes of size 10×15 are shown here



705
706
707
708
709
710
711
712
713
714
Figure 11. An example of 10×9 generated maze.

715 garding the path and dead ends are fully met, which would
 716 render the possible existence of a shape not in the original
 717 ones not of huge importance since it does not affect the
 718 playability of the game. It should be noted however as a
 719 violation and further investigation for potential solutions can
 720 be performed. Initial attempts indicate that this is not such
 721 a trivial task since we need to also retain the satisfaction
 722 of other constraints and altering the shape will lead to a
 723 violation of another rule, such as the single-width path.
 724

725 **6.4. Results: RL Pac-Man agent**

726 When starting the training process it takes a while to train
 727 an agent to avoid violation of the game rules and learn legal
 728 moves. After that, the agent starts to not only do legal moves
 729 but also search for food. As already mentioned in Section
 730 4.4, another problem we faced was getting stuck in local optima,
 731 as the agent learned to collect food at a limited part of the
 732 maze, and afterwards started to wander around aimlessly
 733 without exploring the rest parts of the maze.
 734

735 Despite the fact that we managed to train the agent to
 736 make legal moves and explore a part of the maze, it was not
 737 possible yet to solve the game. Nevertheless, it was obvious
 738 from our experimentation that with further hyperparameter
 739 optimization and training the agent will be able to reach
 740 convergence and solve the maze using the implemented
 741 policy-based RL approaches. Due to limited time and lack
 742 of computational resources though, the next steps needed
 743 to train a sufficient agent were not made yet, but they are
 744 proposed as future work in Section 7.2.
 745

746 **7. Conclusions and Future Work**

747 **7.1. Conclusions**

748 To conclude, the additional player experience models pro-
 749 vided replayability as well as more natural ways to interact
 750 with the game. Our small test group unanimously said that
 751 the addition of new ways to interact with this classic game
 752 has been a fun experience and that the learning curve for
 753 both vision and audio assisted controls has been smooth
 754 and minimal. Therefore, these have not hindered their ex-
 755 perience, and on the contrary it allowed the game to retain
 756 replayability. Furthermore, even though we were unable
 757 to test the implementations with a group of people with
 758 disabilities, we strongly believe that these two modes of
 759 interactions will allow the target demographic to expand
 760 and provide means of interactions to those who previously
 761 lacked it.

762 Based on our search space, the suggested Speech-to-Text
 763 service for the voice game interaction is the Live transcription
 764 service from AssemblyAI, whereas for the vision assisted
 765 implementation the retrieved optimal settings, when using
 766 the mediapipe and OpenCV modules, are the following:
 767

768 **time interval between recorded moves** : 500ms

769 **move threshold** : $scaling_factor \times 4$ pixels, where
 770 $scaling_factor = \frac{z_{mean} - 0.0001}{0.4}$

771 **central area** : square window defined by a distance of 10%
 772 of resolution's width or height from the center point

773 In regards to the PCG with GA, we can conclude that
 774 overall the performance is good since generated mazes seem
 775 to be valid. However there are some cases where due to the
 776 mirroring some undesired shapes may occur. Our current
 777 efforts to fix these cases have not been successful since it
 778 is not such a trivial task to make changes whilst retaining
 779 the previous constraints satisfaction. However, aside from
 780 the issue that arose through mirroring, which might not
 781 necessarily manifest itself with most of the generated mazes,
 782 all of the constraints are fully met and the believability of the
 783 authenticity of the generated mazes is more than satisfactory.

784 The identified optimal settings for the PCG are as follows:

785 **Next Generation selection** +-selection.

786 **Population Size and Offspring Size** 8 for both.

787 **Mating Selection** Both tournament and proportional selec-
 788 tion were deemed equally favorable. Factor of 1.0 and
 789 $k = [2, 3]$.

790 **Crossover** Uniform crossover with $p_c = 0.7$.

791 **Mutation** $p_m = 0.5$ and $s_{p_m} = 0.05$.

792 Current efforts to tune the actor-critic model for the *Pac-*
Man environment were unsatisfactory but displayed a learning
 793 process even though convergence to a winning state was
 794 not achieved. We believe that the current model would be
 795 satisfactory for solving this environment, however further
 796 tuning will be necessary. In an attempt to perform a more
 797 extensive search for optimal hyperparameter settings, we
 798 have developed a framework to be used for experimenting
 799 with large search spaces. The respective hyperparameter
 800 settings and retrieved scores can be accessed and evaluated.
 801 The model which lead to the optimal settings will also be
 802 accessible so retraining would not be necessary. Due to lack
 803 of resources and time, this large-scale training and hyperpar-
 804 ameter tuning was not feasible but it should be useful for
 805 any future efforts to tune this model and train a successful
 806 agent.

807 **7.2. Future Work**

808 The player experience models developed seem to be
 809 achieving our initial goals and expectations, thus it would
 810 be interesting to explore further this area. Typically *Pac-*
Man is a single-player game, but it would be interesting to

770 take advantage of these alternative and natural methods of
771 input in order to allow the game to be played by two players
772 simultaneously. One player could be controlling the ghost
773 agent whilst the other one the eponymous character. This
774 can be achieved by taking advantage of both vision and
775 audio controls at the same time. This idea could be interesting
776 and it can provide more game related content beyond the
777 classical game.

778 For the PCG with GA method to generate *Pac-Man* mazes,
779 it would be advised to investigate potential solutions to the
780 problem that could arise through mirroring. Additional tun-
781 ing of the current hyperparameters can also further optimize
782 the implementation. Currently multi-threading has also been
783 implemented, but just for the creation of mazes in parallel.
784 It would be a good idea to parallelize the maze creation by
785 running each individual's generational procedure in par-
786 allel. Furthermore, different Evolutionary Strategies can also
787 be experimented with for the maze creation, or additional
788 mating selection, crossover and mutation strategies can be
789 added.

790 Finally, as aforementioned, for the RL agent, it would
791 be advised to perform extensive tuning and the tool created
792 for exhaustive search of parameters might be useful for the
793 search of optimal settings. Once the model is tuned, leading
794 to a successful agent using Actor-critic, we can expand to
795 other policy-based approaches or implement further meth-
796 ods that would optimize the current approach. Another idea
797 for extending further the DRL implementation would be
798 capturing the last N frames, thus keeping track of the re-
799 cent history of states, instead of keeping only the last one.
800 The benefit of such approach is that the network can track
801 in which directions ghosts are moving and extrapolate the
802 ghosts' trajectory. To implement such network behavior, a
803 convolutional LSTM network can be used. Different policy
804 optimization such as **Proximal Policy Optimization with**
805 **Covariance Matrix Adaptation** (PPO-CMA) can also be ex-
806 plored. An interesting way that could lead to a more generic
807 agent (that could solve multiple mazes) would be to expand
808 the training procedure, by allowing the agent to interact
809 with more than one environmental structures. This can be
810 achieved by utilizing our PCG with GA model and using
811 multiple generated mazes, of varying sizes, for the training
812 of the agent.

813

814 References

815 Assemblyai. Assemblyai speech-to-text api: Auto-
816 automatic speech recognition. URL [https://www.](https://www.assemblyai.com/)
817 [assemblyai.com/](https://www.assemblyai.com/).

818

819 Bambocher, P. Python interface to cmu sphinxbase and pock-
820 etsphinx libraries. URL [https://github.com/](https://github.com/bambocher/pocketsphinx-python)
821 [bambocher/pocketsphinx-python](https://github.com/bambocher/pocketsphinx-python).

822

823

824

Berkeley, U. Berkeley: The pac-man projects.
URL http://ai.berkeley.edu/project_overview.html.

Csikszentmihalyi, M. *Flow: The Psychology of Optimal Experience*. 01 1990.

Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., et al. An image is worth 16x16 words: Transformers for image recognition at scale. *arXiv preprint arXiv:2010.11929*, 2020.

Haarnoja, T., Zhou, A., Abbeel, P., and Levine, S. Soft actor-critic: Off-policy maximum entropy deep reinforcement learning with a stochastic actor. In *International conference on machine learning*, pp. 1861–1870. PMLR, 2018.

O'Donovan, J., Ward, J., Hodgins, S., and Sundstedt, V. Rabbit run: Gaze and voice based game interaction. In *Eurographics Ireland Workshop, December*, 2009.

PAC-MAN. History the official site for pac-man - video games amp; more. URL <https://www.pacman.com/en/history/>.

Park, H. S., Jung, D. J., and Kim, H. J. Vision-based game interface using human gesture. In *Pacific-Rim Symposium on Image and Video Technology*, pp. 662–671. Springer, 2006.

Plaat, A. Deep reinforcement learning, a textbook, 2022.
URL <https://arxiv.org/abs/2201.02135>.

Spacco, J. Jspacco/pac3man: Python3 port of berkeley pacman. URL <https://github.com/jspacco/pac3man>.

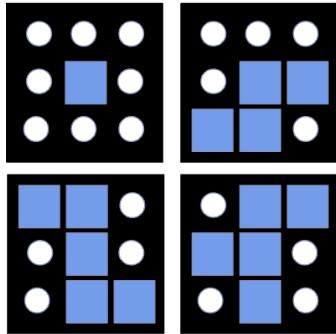
Sutton, R. S. and Barto, A. G. *Reinforcement Learning: An Introduction*. The MIT Press, second edition, 2018. URL <http://incompleteideas.net/book/the-book-2nd.html>.

Uberi. Speech recognition module for python, supporting several engines and apis, online and offline. URL https://github.com/Uberi/speech_recognition.

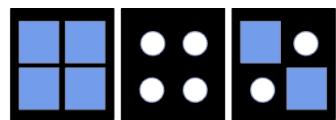
Yannakakis, G. N. and Togelius, J. A panorama of artificial and computational intelligence in games. *IEEE Transactions on Computational Intelligence and AI in Games*, 7(4):317–335, 2015. doi: 10.1109/TCIAIG.2014.2339221.

825 **A. Invalid Blocks**

826 This section includes the set of invalid blocks used for
827 evaluating the individuals.



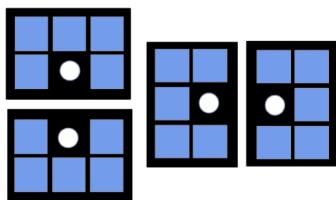
841 *Figure 12.* The 3×3 invalid blocks.
842
843



849 *Figure 13.* The 2×2 invalid blocks.
850
851



858 *Figure 14.* The 4×1 invalid blocks.
859
860
861
862
863
864



865 *Figure 15.* The 2×3 invalid block on the left hand side and 3×2
866 on the right.
867
868
869
870
871
872
873
874
875
876
877
878
879