# IEEE Copyright Notice

# Evaluating Mixed-Initiative Procedural Level Design Tools using a Triple-Blind Mixed-Method User Study

Sean P. Walton, Alma A. M. Rahat and James Stovold

*Abstract*—Results from a triple-blind mixed-method user study into the effectiveness of mixed-initiative tools for the procedural generation of game levels are presented. A tool which generates levels using interactive evolutionary optimisation was designed for this study which (a) is focused on supporting the designer to explore the design space and (b) only requires the designer to interact with it by designing levels. The tool identifies level design patterns in an initial hand-designed map and uses that information to drive an interactive optimisation algorithm. A rigorous user study was designed which compared the experiences of designers using the mixed-initiative tool to designers who were given a tool which provided completely random level suggestions. The designers using the mixed-initiative tool showed an increased engagement in the level design task, reporting that it was effective in inspiring new ideas and design directions. This provides significant evidence that procedural content generation can be used as a powerful tool to support the human design process.

## I. Introduction

GAME developers are under increasing pressure not only to launch games with hours of unique content, but to continue to add new fresh content post launch [1]–[3]. This provides motivation [4]–[6] to develop tools which can support content generation, which is the aim of procedural content generation (PCG) algorithms [7]. It is also important to look beyond this commercial motivation and ask how PCG algorithms can support designers in their creative process for the sake of creativity [8]. PCG algorithms have been developed to create a wide variety of content [4]. In addition to supporting designers, PCG algorithms can also benefit the players, resulting in an increased diversity of content [4], [5], [9] and creating a source of curiosity and unpredictability [10]. Perhaps the most notable example of this in recent years is Hello Game's title *No Man's Sky*[1], a space exploration game in which almost everything is procedurally generated [11]. There are even examples of using PCG as a game mechanic itself, such as in the game Petalz [12] where players breed and share flowers, becoming part of the PCG algorithm itself.

Despite the clear benefits of PCG algorithms, there are still a number of open challenges in the field. For example, the vast majority of PCG algorithms are highly problem specific, often designed for a single genre of game [9] or limited to

S. Walton and A.A.M. Rahat are with the Department of Computer Science, Swansea University, Wales, UK e-mail: s.p.walton@swansea.ac.uk.

J. Stovold is with the Department of Computer Science, British University Vietnam, Vietnam

[1]https://www.nomanssky.com/

specific geometries [13]. A frequently-cited limitation is the lack of control human designers have when generating content using PCG [5], [13]. PCG algorithms are often non-intuitive, requiring designers to tweak and adjust tuning parameters which are difficult to relate to their goals. This ultimately limits the control designers have over the generation process [7] and builds a knowledge barrier [6].

Despite significant investment into researching new methods for PCG, there is little research on how designers interact with these tools [14]. In an attempt to address this gap, Craveirinha and Roque [14] undertook a participatory design process involving game designers and researchers to design an interface for a PCG algorithm. In doing so they explored the attitudes of game designers toward PCG tools. They found that many PCG algorithms work by optimising certain metrics which the algorithm designers have identified as being important for player experience. These metrics, and target values for them, are determined *a priori*. Designers do not operate in this way, but instead explore the design space to determine metrics which can then be used to optimise player experience. The findings were summarised with two key observations which will inform our work: *(1) the tool needs an understandable metaphor,* and *(2) exploration is needed before optimisation.*

### A. Our Contribution

When investigating the existing work in PCG of levels we found that *most* contributions included **no user study** of the created artefact [4]–[6], [9], [15]–[18]. In contributions which did include a user study, the studies were performed in an ad-hock fashion with **no rigorous qualitative analysis** and **no control group** [2], [19], [20]. *The main contribution of our work is a mixed method triple-blind user study comparing a mixed-initiative PCG tool to a tool which gives the designer random suggestions. Using a reflexive thematic analysis approach [21] we find that our mixed-initiative tool does support the creative process, compared to the control group, which adds strength to the findings of the above mentioned studies.*

The tool designed for this study was rooted in the approach introduced by Baldwin et al. [6], but placed into the context of the two observations of Craveirinha and Roque [14] which we re-framed into two design pillars:

1) The designer must interact with the algorithm by designing content, rather than adjusting parameters.
2) The designer will be supported to explore the design space.

## II. Background

### A. Search-Based Procedural Content Generation

There are numerous approaches to PCG [22]. In our work we adopt a search-based approach to PCG as it aligns well with our second design pillar to support exploration. In search-based PCG an algorithm generates a large volume of content and evaluates each item created using a fitness function. There are two key identifying characteristics of a search-based approach: (a) the fitness function allows the comparison and ranking of content, and (b) this ranking is used to inform the generation of new content [22]. Search-based PCG approaches are often implemented using evolutionary algorithms (EAs); optimisation algorithms which aim to minimise a fitness function over several generations. In the context of PCG, an EA will initialise a population of potential designs, rank these according to the quality defined by the fitness function, then create the next generation through stochastic mutation and interbreeding [6]. Search-based approaches have been used to generate a wide range of content including mazes [17], [18], race tracks [23] and dungeon maps [24]. A common aspect of these contributions is that the authors design and specify a fitness function which they argue will result in a good player experience. Our aim is to allow the designer to directly influence the fitness function through design, rather than relying on the fitness function to dictate what is good.

### B. Mixed-Initiative Approaches to Content Generation

As mentioned in the introduction, one of the key challenges of PCG algorithms is that level designers often do not have knowledge of how to control them. This challenge is directly related to our first design pillar, that designers should interact with our system by designing content. Many researchers [2], [5], [6], [13] have made contributions towards addressing this challenge. The work by Liapis et al. [2] and Baldwin et al. [6] are particularly relevant to our goals and inform our approach.

Liapis et al. [2] introduced the Sentient Sketchbook, a tool for supporting designers creating levels for games. As the designer sketches ideas via the tool's interface, real-time feedback is given to the designer based on a number of game play relevant metrics. The tool suggests alternative map designs based on the sketch the designer creates. This is achieved through a genetic search algorithm which attempts to maximise the map's score based on a number of metrics, or a diversity measure. The results of all these searches are presented to the designer. The general feedback from their user study was positive, with users reporting that the tool started pushing them in design directions they did not initially expect.

Baldwin et al. [6] present a mixed-initiative tool for generating dungeon levels using evolutionary algorithms. Their aim was to allow the designer to control the algorithm using parameters with which they are familiar with, based on what they term *game design patterns*, such as mean corridor length or number of enemies. We suggest a slight change in terminology by referring to these as *level design patterns* hereafter. Game design often refers to the design of mechanics in a game rather than the level geometry, so we feel it is clearer to use the term level design patterns when describing these
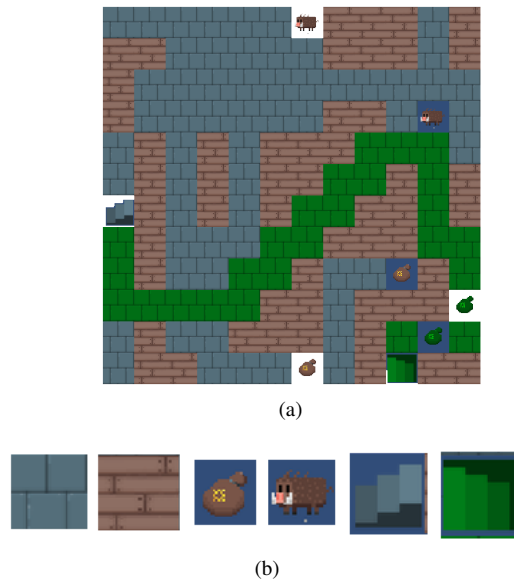


(a)



(b)

Fig. 1. Artwork used to represent map layout and tiles. Assets are distributed by LazerGunStudios without license at https://lazergunstudios.itch.io/roguelike-asset-pack. In (a), a complete map with a path from entrance to exit is shown, while in (b) we show different adjustable components of the map (Floor, Wall, Treasure, Enemy, Entrance and Exit).

metrics. In fact, Baldwin et al. [20] found that the participants in their study found the term game design confusing in this context. Essentially, the level designer specifies targets for the various level design pattern metrics and an evolutionary algorithm attempts to optimise a fitness function based on this. Their results show an impressive ability of control based on these patterns. The tool has been developed over several years since the original paper. More design pattern detection has been implemented [20] and metrics based on visual aesthetics have been added to the tool [15]. Presently, the tool is quite sophisticated, supporting designers to design an entire dungeon rather than single rooms. Efforts were recently made to model designer preference by training a neural network while a designer interacts with the system [19]. Although this is an interesting approach which shows promising results, training a neural network to model the preference of a human is a challenging task which requires lots of data. In our work we have decided to take the original approach presented by Baldwin et al. [6] and their findings on the importance of visual aesthetics [20] and focus on building a tool which only requires the designer to design levels - without tweaking parameters.

## III. Methodology

### A. Specification and Design of the System/Tool

A system was designed in the context of two design pillars – described and justified in Section I-A – to support a level designer in creating a series of 2D maps/levels for a simple dungeon game. An example of a map is shown in Figure 1a. In this study the dungeon maps are made up of 12 by 12 tiles. Each tile has one of six possible values. *Wall*: this is impassable by the player. *Floor*: this is passable by the player. *Treasure*: this is an item which is desirable for the player to

reach. *Enemy*: this is a non-player character which can damage the player, something the player wishes to avoid. A single *Entrance* tile where the player enters the level and a single *Exit* representing the players goal. There must be a passable path between the entrance and exit for a level to be valid. The graphical representation of these tiles is shown in Figure 1b.

When surveying the search-based PCG literature we observed two key points:

1) Search-based PCG is inherently a multi-objective problem
2) The majority of researchers tackle this multi-objective problem by combining the results from multiple fitness functions into one scalar value through a weighted sum.

An exception to this is the work by Loiacono et al. [23] who used a multi-objective optimisation algorithm without scalarisation. They found an interesting diversity of solutions along the Pareto fronts, which has the potential to support our second design pillar. Although there are many advanced techniques for multi-objective optimisation and finding the Pareto front [25] we opt for a simple approach which is detailed in Section III-D.

Since we wish our designers to interact with our system through designing levels, we turn to the approach by Liapis et al. [2] as a starting point. In their approach suggestions are presented to the designer by optimising predetermined fitness functions with the designer's initial design as a starting point. In our approach the level designer will design the first level, the system will then calculate some metrics which describe that level and record those as targets. An evolutionary optimisation algorithm will then randomly initialise a population and try to match the metrics from the user-designed level. Preuss et al. [16] found that restarting their evolutionary algorithm performs as well as advanced approaches to increasing novelty and diversity. Therefore we will restart our algorithm at regular intervals and use this opportunity to allow the level designer to influence the target metrics at run time. This will be achieved by allowing the level designer to edit and select maps produced by the system which are desirable. The system will store the metrics of these *liked* maps and use them in fitness function evaluations.

### B. System Overview

Algorithm 1 gives an overview of the final system. The user is initially presented with a blank canvas to design an initial level, once finished the user clicks the submit button. After the optimisation algorithm has finished running the user is presented with the view shown in Figure 2. The eight maps displayed are a selection from the feasible population of the final generation produced by the optimisation algorithm. Each of the eight maps can be edited by the user; clicking a tile in a map cycles it between all the values possible in turn. The user can additionally tag any number of these maps as *like* or *keep* using the checkboxes below the maps. The *like* tag is used as part of the optimisation process, and the *keep* tag indicates that this map should be included in the final set of levels designed. At the top of the view the levels tagged keep are shown. Once happy with their edits and tags the user

---

**Algorithm 1** System Overview

1: user designs first level $\mathbf{x}_1$
2: store $\mathbf{x}_1$ in the list of liked maps and the list of levels
3: **repeat**
4:     run optimisation algorithm
5:     display a subset of maps from the final generation of the GA
6:     user may edit maps and tag them as like and/or keep
7:     **for** each map $\mathbf{x}_i$ **do**
8:         **if** $\mathbf{x}_i$ is tagged like or keep **then**
9:             store $\mathbf{x}_i$ in the list of liked maps
10:             **if** user has tagged $\mathbf{x}_i$ to keep **then**
11:                 add $\mathbf{x}_i$ to the list of game levels
12:             **end if**
13:         **end if**
14:     **end for**
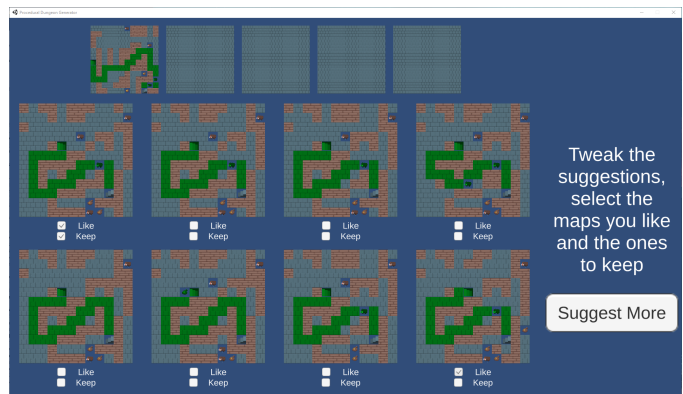15: **until** the list of game levels is full

---



Fig. 2. Feedback view of the system. At the top row (five smaller windows), the user can see the designs that they have already chosen or created. In middle and bottom rows, we show the generated levels, and provide options for keeping or liking designs. These generated levels can be edited by the user by clicking on their tiles. On the right, the user has the option to request further suggestions.

clicks "Suggest More" and the tool generates eight more maps replacing those displayed previously. The tool is open source and can be downloaded via GitHub[2].

### C. Metrics used to Define the Fitness Functions

The fitness functions based on level design patterns designed by Baldwin et al. [6] show an impressive ability to control the types of maps generated by their search algorithm. Therefore, we have opted to use these functions along with visual impression metrics which Preuss et al. [16] found to be highly effective. In total there were 31 metrics used to characterise a map design. The metrics are split into two broad categories: level design patterns (III-C1 to III-C7) and visual impression metrics (III-C8 to III-C9). We use the notation that $M_i(\mathbf{x}_k)$ is metric $i$ calculated for the map $\mathbf{x}_k$. It is worth emphasising that we are not optimising these metrics directly, but using them to

---

[2]https://github.com/seanwalton/mixed-initiative-procedural-dungeon-designer

construct the fitness function which will drive the optimisation process.

*1) Path Length:* $M_1$ is the path length, $P(Entrance, Exit)$, measured in number of tiles, divided by the total number of tiles in the map, $N_{total}$.

*2) Global Wall to Passable Tile Ratio:* $M_2$ is the ratio of walls to non–wall tiles in the map.

*3) Corridor Metrics:* Corridors are defined as horizontal or vertical series of passable tiles enclosed by impassable tiles on either side [6]. In our implementation corridors of length one are counted. The metrics $M_3$ to $M_6$ are the number of corridors followed by the maximum, minimum and mean corridor lengths.

*4) Chamber Metrics:* A chamber is defined as a continuous block of passable tiles which are wider than a corridor. A less rigid definition is followed than the one outlined by Baldwin et al. [6]. In their work these metrics are used to generate dungeons using user inputs such as chamber size, therefore they have to consider what a user might expect a chamber to look like. In our work these metrics are only used to compare the structure of two maps, we do not want to assume a minimum chamber size. Chambers are identified following corridor identification. Once chambers are identified two qualities for each chamber is calculated, the area $k_i^A$ and squareness $k_i^S$ given by:

$$k_i^A = k_i^h k_i^w \tag{1}$$

$$k_i^S = \frac{k_i^A}{\min(k_i^h, k_i^w)^2} \tag{2}$$

Where $k_i^h$ and $k_i^w$ are the height and width of chamber $i$. This then leads to 7 metrics ($M_7$–$M_{13}$) for chambers. The total number of chambers, the maximum, minimum and mean chamber areas and maximum, minimum and mean chamber squareness.

*5) Dead Floor Tiles:* A dead tile is defined as a passable tile which has not been identified as a chamber or corridor. These often appear as tiles which connect multiple corridors or chambers. The metric, $M_{14}$ is simply the number of these tiles divided by the total number of tiles.

*6) Entrance Metrics:* Two metrics are defined for the number of treasure and enemy tiles around the entrance [6]. $M_{15}$ is the minimum area around the entrance tile which does not contain an enemy tile, and $M_{16}$ is the minimum area around the entrance tile which does not contain a treasure tile.

*7) Enemy and Treasure Metrics:* $M_{17}$ and $M_{18}$ are simply the fraction of enemy and treasure tiles respectively. In addition a safety measure, defined in [6], is calculated for each treasure and $M_{19}$ and $M_{20}$ are the mean and standard deviation of this.

*8) Visual Symmetry of Wall Tiles:* Preuss et al. [16] introduced a number of visual symmetry metrics which we have adapted for use here. Two lines of symmetry are defined along the centre of the map horizontally and vertically. The number of a specific type of tile is counted either side of these lines then used to calculate ratios. For example,

$$N_{wall}^{top}$$

is the number of wall tiles in the top half of the map and

$$N_{wall}^{left}$$

is the number of wall tiles in the left half of the map. A total of 8 metrics are defined based on these ratios, for example the left to right wall tile ratio is:

$$M_{21}(\mathbf{x}_k) = \frac{\left| N_{wall}^{left} - N_{wall}^{right} \right|}{N_{wall}} \tag{3}$$

There is also a top to bottom wall ratio $M_{22}$, left to right and top to bottom enemy and treasure ratios ($M_{23}$–$M_{26}$), and treasure to enemy ratios defined as:

$$M_{27}(\mathbf{x}_k) = \frac{\left| N_{treasure}^{left} - N_{enemy}^{right} \right|}{N_{treasure} + N_{enemy}} \tag{4}$$

$$M_{28}(\mathbf{x}_k) = \frac{\left| N_{treasure}^{top} - N_{enemy}^{bottom} \right|}{N_{treasure} + N_{enemy}} \tag{5}$$

For equations 3 to 5 if the denominator would be zero the metric is given a value of zero.

*9) Exact Symmetry Metrics:* As well as the visual symmetries we also introduce and define 3 metrics which give a measure of exact reflection over the symmetry lines used for ($M_{21}$–$M_{28}$). In addition a measure of rotational symmetry is considered by comparing the map against its transpose. These metrics ($M_{29}$–$M_{31}$) are calculated by simply counting the number of tiles that exactly match their reflected counterpart across the various symmetry lines (or to the tile at its transposed location) and express them as a fraction of the total number of tiles.

### D. Genetic Algorithm

Following the approach of Baldwin et al. [6] we use a feasible–infeasible two-population (FI-2Pop) genetic algorithm (GA) [26] as our evolutionary optimisation algorithm. FI-2Pop is the same as a standard GA but splits the population into feasible and infeasible sub-populations. Maps are considered feasible if a valid path from the entrance to exit exists, and are automatically entered into the correct sub-population once created. In our system the only difference between evaluating fitness in these two populations is that $M_1$, the path length, is not considered for the infeasible maps. A tournament selection approach is used to select individuals (from the same sub-population) to reproduce and a number of elite individuals survive from one generation to the next.

*1) Fitness Functions and Ranking Procedure:* The fitnesses of an individual map are defined as

$$f_i(\mathbf{x}_k) = \min_{t=1}^{T} |M_i(\mathbf{x}_k) - M_i(\mathbf{x}_t)| \tag{6}$$

where $i \in [1, 31]$, $\mathbf{x}_t$ is one of the $T$ levels a user has liked or stored. This is a form of goal programming approach, where our aim is to generate maps that have metrics similar to those provided by the designer. Thus, lower fitness values are desirable.

Map $\mathbf{x_j}$ is said to dominate map $\mathbf{x_k}$, denoted as $\mathbf{x_j} \prec \mathbf{x_k}$, iff $f_i(\mathbf{x_j}) \le \mathbf{f_i}(\mathbf{x_k})$ for all $i \in [1, 31]$, and there is at least one

fitness function $f_l$ for which $f_l(\mathbf{x_j}) < \mathbf{f_l(x_k)}$ [27]. The Pareto set of mutually non-dominated solutions is defined as:

$$\mathcal{P} = \{\mathbf{x} \mid \mathbf{x}' \nprec \mathbf{x}, \forall \mathbf{x}, \mathbf{x}' \in \mathcal{X} \wedge \mathbf{x} \neq \mathbf{x}'\}, \qquad (7)$$

where $\mathcal{X}$ is the feasible decision space. Typically, it is impossible to exactly locate the Pareto set, so a representative approximation, $\mathcal{P}^*$, is often sufficient.

Optimising more than three objectives simultaneously is referred to as a many-objective optimisation problem [28]. As the number of objectives increase so does the probability of finding solutions which improve at least one of them, compared to existing solutions, leading to almost all solutions in the search space becoming Pareto optimal [29]. Thus locating a representative approximation of the Pareto set is extremely challenging [30]. An *ad hoc* approach to identify a smaller subset of $\mathcal{P}^*$ is often required. Osyczka *et al.* proposed to prune $\mathcal{P}^*$ such that only one solution within a predefined interval is retained and others in close proximity are deleted [31]. As identifying a sensible interval is not straightforward in a high-dimensional objective space, such an approach may still result in a large number of solutions in the $\mathcal{P}^*$.

In this paper, we, therefore, take a more aggressive approach in order to keep candidates of interest in our population. We aim to only retain maps that are **more similar** to the ones selected by the user: this is a *majority voting* approach, popular in decision making with ensemble of systems [32]. We deliberately ignore the magnitude of the differences in various objectives for two reasons. Firstly, the standard euclidean distance in high-dimensions loses its efficacy in objectively identifying how different two maps are [33]. Secondly, it is not obvious if there is a natural preference between objectives. As such our *ad hoc* measure for ranking and retaining solutions in the population is: given $\mathbf{x_j}, \mathbf{x_k} \in \mathcal{X}$, if more $f_i(\mathbf{x_j}) \leq f_i(\mathbf{x_k})$ than $f_i(\mathbf{x_k}) \leq f_i(\mathbf{x_j})$ for $i \in [1, 31]$, then we prefer $\mathbf{x_j}$ over $\mathbf{x_k}$.

*a) Majority Voting and Pareto Optimality:* For majority voting to be effective in ranking and retaining solutions of interest, it must have the following characteristic:

$C$. The dominance relationship must be preserved, that is if for two arbitrary solutions $\mathbf{x}$ and $\mathbf{x}'$, $\mathbf{x} \prec \mathbf{x}'$, then majority voting must rank $\mathbf{x}$ better than $\mathbf{x}'$.

*Proof.* Comparing two solutions $\mathbf{x}$ and $\mathbf{x}'$ with $K$ objectives, we can evaluate sets of the following relationships: $\mathcal{L} = \{i \in [1, K] \subseteq \mathbb{N}_1 | f_i(x) < f_i(x')\}$, $\mathcal{E} = \{j \in [1, K] \subseteq \mathbb{N}_1 | f_j(x) = f_j(x')\}$ and $\mathcal{G} = \{k \in [1, K] \subseteq \mathbb{N}_1 | f_k(x) > f_k(x')\}$, where $i \neq j \neq k$ and $\mathbb{N}_1$ is the set of positive integers.

According to majority voting, we can conclude that $\mathbf{x}$ is better than $\mathbf{x}'$ iff:

$$|\mathcal{L}| + |\mathcal{E}| > |\mathcal{G}|. \qquad (8)$$

By definition of Pareto optimality (see (7)), $|\mathcal{G}| = 0$ when $\mathbf{x} \prec \mathbf{x}'$. Hence, for any number of objectives greater than 0, (8) must be true, and thus $\mathbf{x}$ will be ranked higher than $\mathbf{x}'$. Thus, $C$ holds. $\qquad \square$

*b) Implications on the Search Population:* We start our search with an arbitrary population of fixed size. Each individual in the population is ranked using the majority voting approach. Since $C$ holds, the best ranked individual is the best approximation of a Pareto optimal solution in the population, and there may be mutually non-dominated solutions in the population based on the balance reached between the sizes of the sets $\mathcal{L}$, $\mathcal{E}$ and $\mathcal{G}$.

When we update the population, we only allow a child $\mathbf{x}$ to replace a solution $\mathbf{x}'$ in the population if and only if $\mathbf{x}$ is better ranked. Clearly, this means one of following two possibilities is met for replacing an individual in the population:

(i) $\mathbf{x} \prec \mathbf{x}'$, i.e. individual is dominated by the child.
(ii) $\mathbf{x} \nprec \mathbf{x}'$, i.e. individual and child are mutually non-dominated, but child satisfies (8).

This way we are creating evolutionary pressure to eradicate dominated individuals from existing population. Thus we are highly likely to retain only mutually non-dominated solutions in the final population, and generate a fixed sized approximation of $\mathcal{P}$.

*2) Crossover and Mutation:*

*a) Crossover:* Two parent maps $\mathbf{x}_i$ and $\mathbf{x}_j$ are crossed to create a child map $\mathbf{x}_c$ by first randomly picking an entrance and exit from the parents such that $\mathbf{x}_c$ has exactly one entrance and one exit. The remaining tiles in $\mathbf{x}_c$ are randomly selected from $\mathbf{x}_i$ and $\mathbf{x}_j$ with equal probability.

*b) Mutation:* To mutate $\mathbf{x}_c$, a random tile is selected and then swapped with a random adjacent tile to create the mutated map.

*3) Selecting Tuning Parameters:* The performance of GAs is highly problem dependent [34]. It is therefore crucial to carry out a parameter sensitivity study for each new application to maximise performance. [35]. The map shown in Figure 1a was selected for use in the parameter study since it has a balance of corridors and chambers. For brevity we do not present the detailed results of our studies, but explain our process and present the final parameters used. For each combination of parameters we performed 30 tests with different random seeds for the random number generator. We then compared mean performance to select the final set of parameters. It is challenging to define good performance in the context of our aims which are primarily human focused. As an approximate measure we compared the sum of fitnesses of the best level generated for different combinations of parameters. In all tests the number of objective function evaluations was kept constant at $10,000$. This was a decision made based on the time taken for an optimisation run to complete on the machines used in the user study, with the aim of limiting the participation time of the user study to 30 minutes. The best performing set of parameters and methods were found to be: Mutation Rate: 0.5, Tournament Size: 2, Number of Elite: 1, Population Size: 20 and Number of Generations: 500. Examples showing how the GA performs when driven without human input are provided in the supplementary materials for this paper.

## IV. User Study

### A. Methodology

Yannakakis et al. [36] introduce an assessment methodology for mixed-initiative systems. They recommend evaluating how often the computational creations are used by the designer, and whether or not those creations changed the thinking process of the designer; our user study was designed to evaluate these aspects. Ethical approval for the study was obtained from the Swansea University College of Science ethics committee[3]. Our original plan was to perform the study in lab conditions, however due to the COVID-19 pandemic we had to change our methodology and carry out the study on-line. Four participants did complete the study in lab conditions prior to the UK lockdown, every effort was made to ensure parity between the lab and on-line experiments. Participants were recruited through social media and the research team's professional networks. The only requirements for taking part in the study were that you had to be aged 18 or over and have access to an internet-connected computer running Windows or Linux. Participants were each given an information sheet which explained that we were investigating approaches people take when designing levels for video games, with the aim to better understand this process to enable us to make level design tools. They were then asked to "create 5 levels for a simple dungeon game using a computer assisted tool." A set of instructions for using the tool and what constituted a valid level were provided along with the tool itself.

Some slight modifications to the tool were made for the user study. Alongside the suggestions from the system, participants were given a blank canvas where they could design a new level from scratch. Before starting the process the tool asked the participant to enter a unique ID. Based on this ID there was a 50% chance the tool used the GA designed in this paper to generate suggested maps. In all other cases maps were randomly generated with no optimisation at all. This was done using a triple-blind approach, neither the participant or researchers knew which algorithm had been selected until after the data was analysed. The result is that we have two groups of participants to compare, the GA group and the control group (who were given random suggestions). Once the participant completed the game design task they were asked to upload log files which contained quantitative results and answer a series of free response questions.

*a) Quantitative Measures:* Each participant submitted a log file which contained the following quantitative measures:

- Which participant group they belong to (GA or control)
- The number of maps the participant marked as like or keep at each iteration.
- The number of times the participant created a map from scratch using the blank editor.
- How much a participant tweaked a suggested design if they decided to keep or like it.

Participants were also required to submit a screenshot of the final screen of the tool, which includes the 5 levels they created, these are included in the supplementary materials for this paper.

*b) Qualitative Questions:* Each participant was then asked 4 questions with a free text response. The questions were:

1) Describe the process you took to design a new level.
2) Was designing 5 levels challenging, or could you have easily designed many more? Explain your answer.
3) Did the tool affect the way you designed your levels? Explain your answer.
4) How would you describe the tool to someone else?

To analyse the responses an inductive coding approach was adopted. Codes were created by reading through all responses, to all questions, independently by each member of the research team. These codes were combined into a final set of codes for each question, which were used for the final coding which was performed by SW. This analysis was all carried out before participant responses were linked to their group, making our study triple-blind.

### B. Materials

A total of 24 participants took part in the study. Of those 17 (71%) were male, 6 (25%) female and one (4%) did not disclose their gender. The mean age of participants was 25.2 years (SD = 7.81, range = 18 to 48). Participants were asked two questions relating to the frequency with which they play video games and their experience with designing levels. The majority (83%) of participants reported that they play games frequently, more than once a month. Participants rated their level design experience on a Likert scale from 1 *No Experience* to 5 *Level Design is my primary profession*. The mean self reported experience of the participants was $2.2 \pm 0.2$, with range 1 to 5. When reporting experience values the standard error in the mean is presented. A total of 14 (58.4%) participants were given suggestions from the GA and 10 (41.6%) were in the control group. The self reported experience of the two groups was comparable, $2.1 \pm 0.2$ for the GA and $2.2 \pm 0.4$ for the control. 5 of the 24 participants failed to correctly upload log files following the user study resulting in a total of 19 quantitative data points, of which 11 (58%) were given level suggestions by the GA and 8 (42%) were in the control group. When analysing the quantitative data Welch's t-test was used to determine statistically significant deviations between the means of the two groups, p-values of less than 0.05 were considered statistically significant.

### C. Results

Tables I to IV show the codes and frequencies for all the qualitative questions, along with the mean self reported experience (as described in IV-B) of participants who responded with each code. When answering Q1 participants predominately (N=20, 83.3%) reported considering the player experience when designing levels. For example, *"The first level ensured an easy layout where everything is encountered and choice is allowed..."*. There were some notable differences between the ways the two groups answered this question. More participants in the GA group were interested in creating levels which rewarded and encouraged exploration (57% compared to 30%),

### TABLE I
### Q1: Describe the Process you Took to Design a New Level

| Code | Control | GA | Experience |
|---|---|---|---|
| *Thoughts relating to level design approach* | | | |
| Considered player experience/game mechanics | 8 (80%) | 12 (86%) | $2.3 \pm 0.3$ |
| Creating Risk-Reward Trade-off/Balance | 3 (30%) | 8 (57%) | $2.4 \pm 0.3$ |
| Encourage/reward exploration | 3 (30%) | 8 (57%) | $2.3 \pm 0.3$ |
| Focused on the path from entrance to exit | 4 (40%) | 6 (43%) | $1.8 \pm 0.1$ |
| Creating interesting decisions for the player | 4 (40%) | 6 (43%) | $2.2 \pm 0.3$ |
| Incremental complexity/difficulty | 2 (20%) | 4 (29%) | $2.0 \pm 0.5$ |
| Considered visual aesthetics | 4 (40%) | 1 (7%) | $2.6 \pm 0.5$ |
| Aimed to create diversity | 0 (0%) | 3 (21%) | $2.7 \pm 0.5$ |
| Used prior experience | 1 (10%) | 1 (7%) | $1.5 \pm 0.4$ |
| Unstructured approach | 1 (10%) | 1 (7%) | $1.5 \pm 0.4$ |
| *Thoughts relating to the system/tool* | | | |
| Tweaked/edited suggestions from the system | 1 (10%) | 3 (21%) | $2.5 \pm 0.4$ |
| Not satisfied by the suggested levels | 0 (0%) | 1 (7%) | $2.0 \pm 0.0$ |
| Used suggestions from the system | 1 (10%) | 0 (0%) | $2.0 \pm 0.0$ |

### TABLE II
### Q2: Was Designing 5 Levels Challenging?

| Code | Control | GA | Experience |
|---|---|---|---|
| *Comments related to challenge* | | | |
| It was challenging to design multiple levels | 4 (40%) | 6 (43%) | $1.9 \pm 0.3$ |
| It was easy to produce lots of maps | 3 (30%) | 7 (50%) | $1.9 \pm 0.3$ |
| The designs I created ended up similar | 1 (10%) | 2 (14%) | $2.0 \pm 0.5$ |
| *Comments related to tool/system* | | | |
| The tool was useful/helped | 3 (30%) | 2 (14%) | $2.2 \pm 0.4$ |
| The levels generated by the system changed my approach | 1 (10%) | 1 (7%) | $3.0 \pm 0.7$ |
| The tool made it difficult | 0 (0%) | 1 (7%) | $4.0 \pm 0.0$ |
| *Comments related to the task* | | | |
| The limited design space/options made it challenging | 2 (20%) | 3 (21%) | $2.8 \pm 0.7$ |
| It was enjoyable/fun/interesting | 1 (10%) | 3 (21%) | $3.0 \pm 0.5$ |
| The rules of the game were not well defined, so it was difficult | 2 (20%) | 1 (7%) | $2.3 \pm 0.7$ |
| Took longer than expected | 0 (0%) | 2 (14%) | $2.0 \pm 0.0$ |

### TABLE III
### Q3: Did the Tool Effect the way you Designed your Levels?

| Code | Control | GA | Experience |
|---|---|---|---|
| *Description of the effectiveness* | | | |
| It did effect my approach | 1 (10%) | 4 (29%) | $1.8 \pm 0.2$ |
| It moderately effected my approach | 1 (10%) | 3 (21%) | $3.0 \pm 0.5$ |
| It did not effect my approach | 2 (20%) | 1 (7%) | $1.3 \pm 0.3$ |
| *Discussion of the suggestions presented by tool/system* | | | |
| I tweaked suggestions from the system | 2 (20%) | 4 (29%) | $2.0 \pm 0.2$ |
| The suggestions changed my approach | 2 (20%) | 4 (29%) | $3.0 \pm 0.4$ |
| It is good for generating starting points | 2 (20%) | 4 (29%) | $2.0 \pm 0.2$ |
| The suggestions seemed random | 1 (10%) | 2 (14%) | $2.3 \pm 0.7$ |
| I kept generating maps until something good appeared | 2 (20%) | 1 (7%) | $1.7 \pm 0.3$ |
| Suggestions not varied enough | 0 (0%) | 2 (14%) | $2.0 \pm 0.0$ |
| No suggestions were useful/helpful | 1 (10%) | 1 (7%) | $1.5 \pm 0.4$ |
| I had to significantly modify the suggestions | 0 (0%) | 2 (14%) | $3.0 \pm 0.7$ |
| Suggestions rarely got the treasure/enemy layout right | 0 (0%) | 2 (14%) | $2.5 \pm 0.4$ |
| I tried to influence the suggestions | 0 (0%) | 1 (7%) | $2.0 \pm 0.0$ |
| Some of the generated maps were unsuitable | 0 (0%) | 1 (7%) | $2.0 \pm 0.0$ |

### TABLE IV
### Q4: How Would you Describe the Tool to Someone Else?

| Code | Control | GA | Experience |
|---|---|---|---|
| *It is a tool which works with the designer* | | | |
| It learns from your seed designs | 2 (20%) | 5 (36%) | $1.9 \pm 0.2$ |
| It generates starting points - you'll need to edit them | 1 (10%) | 4 (29%) | $1.6 \pm 0.2$ |
| It suggests different levels to you | 1 (10%) | 4 (29%) | $2.4 \pm 0.6$ |
| It helps inspire new ideas | 1 (10%) | 1 (7%) | $1.5 \pm 0.4$ |
| It is a rapid prototyping tool | 1 (10%) | 1 (7%) | $3.0 \pm 0.7$ |
| It is an interactive tool for PCG | 0 (0%) | 1 (7%) | $2.0 \pm 0.0$ |
| *It is a tool which works independently from the designer* | | | |
| It randomly generates levels | 2 (20%) | 2 (14%) | $1.7 \pm 0.3$ |
| No inclusion of human approach to games | 0 (0%) | 1 (7%) | $3.0 \pm 0.0$ |
| *Description of UI/UX* | | | |
| Functional description of UI | 3 (30%) | 2 (14%) | $1.4 \pm 0.2$ |
| It is fun/enjoyable | 1 (10%) | 1 (7%) | $1.5 \pm 0.4$ |
| The tool can be tedious | 0 (0%) | 1 (7%) | $2.0 \pm 0.0$ |

and creating a diverse set of levels (21% compared to 0%). The participants focusing on diversity tended to be those with more design experience. Furthermore, more participants in the GA group reported tweaking suggestions from the system (21% compared to 10%).

When answering Q2 50% of the GA group described the task as easy, compared to 30% of the control group. 14% of the GA group reported that the tool was helpful, compared to 30% of the control group. However, when answering Q3, 50% of the GA group reported that the tool had an effect on their design approach compared to 20% of the control group.

The only statistically significant difference between the two groups in the quantitative data was the number of edits of suggestions from the system. The mean number of edits of liked maps was 13.00 (SD = 14.59) and 2.82 (SD = 5.34) for the GA and control groups respectively. For kept maps the mean number of edits was 14.81 (SD = 14.89) for GA, and 4.10 (SD = 6.02) for the control group. In both cases the p-value was less than 0.01, the full distributions are shown in Figure 3. When put in the context of our qualitative findings, which suggest that the participants in the GA group were more engaged with the task, we can interpret this increased number of edits as increased engagement. When answering Q3, partic-
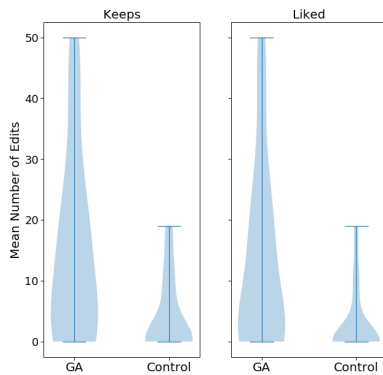
Fig. 3. Comparing the number of edits of liked and kept levels between the GA and control group. The differences in the means of these distributions is statistically significant (p-value < 0.01).


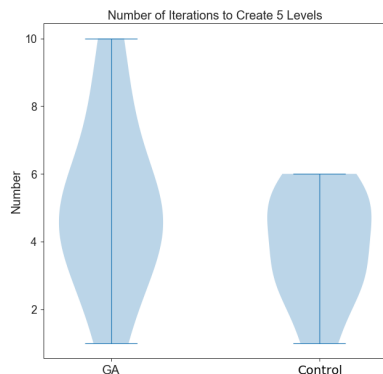
Fig. 4. Comparing the number of iterations taken to design 5 levels, when users were given level suggestions by the GA to random suggestions (p-value = 0.22)

ipants in the GA group were more likely to discuss tweaking suggestions, how the suggestions changed their approach and how the suggestions were good starting points (29% compared to 20% for all three responses). Participants who reported that the suggestions changed their approach had more design experience. When discussing the level of challenge of the task, 21% of the GA group described the task as enjoyable compared to 10% of the control group, these participants tended to have more design experience. When answering Q2, 14% of the GA group reported that the task took longer than expected, compared to 0% of the control group. This could be evidence of increased engagement, but it was unclear if the participants saw this positively or negatively. There was a trend that participants in the GA group took more iterations to design 5 levels, but there is not enough data to show statistical significance. The mean number of iterations taken to design 5 levels by the GA group was 5.08 (SD = 2.39) and the control group 3.88 (SD = 1.53), the distribution is plotted in Figure 4.

The mean number of likes per iteration for GA group was 1.31 (SD = 1.24) and control group 1.84 (SD = 1.15). The p-value for this comparison was 0.71 meaning that there was no statistical significance. In the group of participants who were presented suggestions by the GA there were 2 cases where a user used the blank canvas to create a new design from scratch, and 1 case in the control group. We can not conclude

a statistically significant difference from this data. This data is plotted in Figure 5.

In the final question participants were asked to describe the tool, Table IV shows the results from coding the answers to this question. There were two identifiable groups of description based on how a level designer interacts with the tool. Participants either described it as a tool which works with or independently to the designer. More participants from the GA group described the tool as something which works with the designer. 36% of the GA group described the tool as learning from your designs compared to 20% from the control group. Interestingly, 29% of the GA group stated that the tool generates starting points which you are required to edit, compared to 10% of the control group. This puts the responses for Q3 into context, which suggest that the GA group were more engaged with modifying the suggestions from the system. 14% of the GA group described the tool as something which randomly generates levels, compared to 20% from the control group. Participants from the GA group were less likely to simply describe the UI features of the tool than those from the control group (14% compared to 30%).

## V. DISCUSSION

A common thread throughout the qualitative data was that those participants who were given suggestions by the GA talked a lot more about the suggestions the system gave them. They described the tool as learning from the designs they created and as a tool which works with the designer to support prototyping. The participants from the control group focused more on the functional description of the UI and generally provided less detailed responses to questions. This general lack of engagement from participants in the control group is further supported by the quantitative data which showed that the GA group edited suggestions by the system more than the control group. Initially we thought this was an indication that the GA was doing a bad job, but when taken in context of the qualitative data we found that these participants were considering their designs much more—as one participant stated, the suggestions sparked new ideas. This highlights that a mixed-methods approach is essential when evaluating mixed-initiative systems, quantitative data only tells half the story. The engagement narrative is further supported by the answers to Q1 where participants from the GA group were more likely to consider higher level design concepts, such as rewarding exploration, in their design process. There was also a general trend that the GA group took more iterations to complete the task. We were surprised to find that as many participants in the control group described the suggestions as being useful as in the GA group, suggesting that any suggestions are helpful to the creative process. When comparing the self reported experience levels of participants we found that more experienced participants (a) were more concerned about diversity, (b) reported that the tool changed their approach and (c) enjoyed the task. Overall our data shows that the system we designed does support designers through the design process and is more effective than random suggestions.
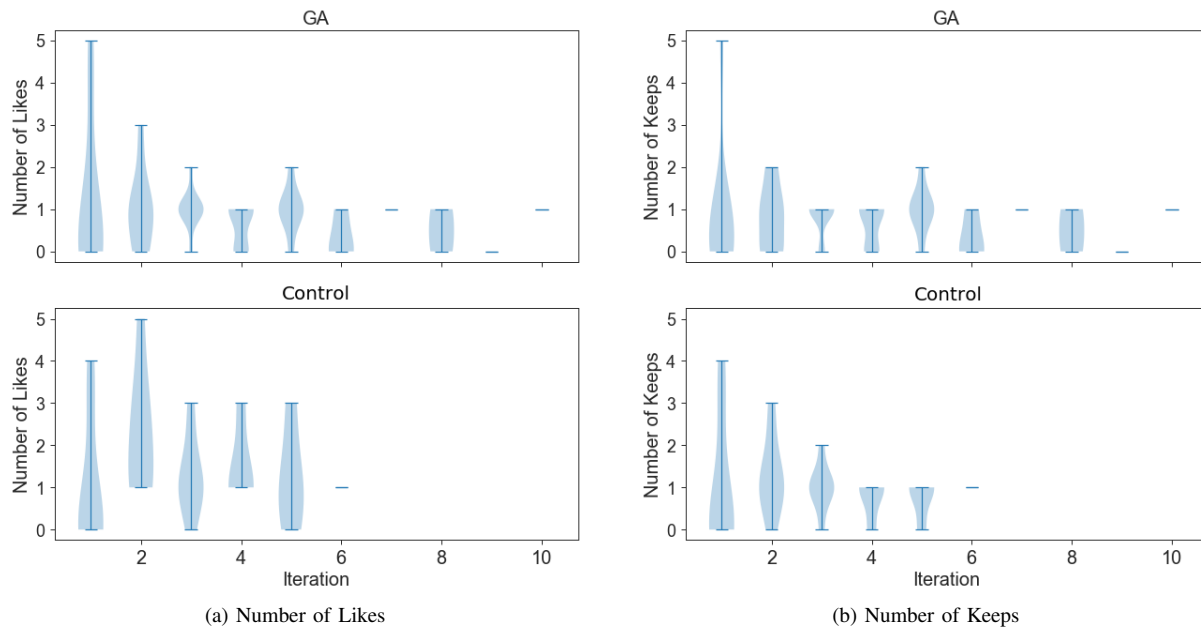
Fig. 5. Comparing the distribution of the number of Likes and Keeps at each iteration between the GA and Control groups

### A. Evaluation of Our Scientific Approach

In hindsight it would have been appropriate to have included some Likert scale questions as part of our user survey. In particular, with Q3 we found that not all participants clearly stated if the tool affected their approach, which would have been captured by a scale response. Performing the study on-line introduces problems such as not all participants correctly submitting log files and possible minor differences in experience based on the hardware they are running. A larger group of participants would have resulted in the generation of more or different codes during the thematic analysis, but that does not mean those codes would have been better. The objective of thematic analysis is not to determine all possible themes, but to generate themes based on the data collected, and to use those to identify patterns of meaning to answer research questions [21], [37]. The patterns we identified align well with other research in this field and add further evidence that mixed-initiative tools can support the creative process.

### B. Future Work

One limitation of the tool was that a number of participants in the GA group noted that the suggestions which were given were all too similar to each other. In the future it would be interesting to add more sophisticated mechanisms [4], [16], [38], [39] to ensure diversity in the suggestions. A further line of enquiry would be to take a model based approach as explored by Alvarez et al. [19] and build a model to predict which maps the designer has a preference for. We suggest that Bayesian Optimisation would be a good avenue to explore in addition to machine learning, this is a model based approach that builds a surrogate model of the function it is optimising. In this application that surrogate model could be of the designers design preference.

### REFERENCES

[1] J. Roberts and K. Chen, "Learning-Based procedural content generation," *IEEE Trans. Comput. Intell. AI Games*, vol. 7, no. 1, pp. 88–101, Mar. 2015.
[2] A. Liapis, G. N. Yannakakis, and J. Togelius, "Sentient sketchbook: Computer-aided game level authoring," in *FDG*, 2013, pp. 213–220.
[3] M. Cook and S. Colton, "Multi-faceted evolution of simple arcade games," in *2011 IEEE Conference on Computational Intelligence and Games (CIG'11)*, Aug. 2011, pp. 289–296.
[4] A. S. Melotti and C. H. V. de Moraes, "Evolving roguelike dungeons with deluged novelty search local competition," *IEEE Trans. Comput. Intell. AI Games*, vol. 11, no. 2, pp. 173–182, Jun. 2019.
[5] A. S. Ruela and K. Valdivia Delgado, "Scale-Free evolutionary level generation," in *2018 IEEE Conference on Computational Intelligence and Games (CIG)*, Aug. 2018, pp. 1–8.
[6] A. Baldwin, S. Dahlskog, J. M. Font, and J. Holmberg, "Mixed-initiative procedural generation of dungeons using game design patterns," in *2017 IEEE Conference on Computational Intelligence and Games (CIG)*, Aug. 2017, pp. 25–32.
[7] R. van der Linden, R. Lopes, and R. Bidarra, "Procedural generation of dungeons," *IEEE Trans. Comput. Intell. AI Games*, vol. 6, no. 1, pp. 78–89, Mar. 2014.
[8] M. Cook, "A vision for continuous automated game design," in *Proceedings of the 13th Experimental AI and Games Workshop*. AIIDE, Jul. 2017.
[9] S. Snodgrass and S. Ontañón, "Learning to generate video game maps using markov models," *IEEE Trans. Comput. Intell. AI Games*, vol. 9, no. 4, pp. 410–422, Dec. 2017.
[10] M. Cook, "Make something that makes something: A report on the first procedural generation jam," in *Proceedings of the Sixth International Conference on Computational Creativity (ICCC 2015)*, H. Toivonen, S. Colton, M. Cook, and D. Ventura, Eds. Park City, Utah: Brigham Young University, Jun. 2015, pp. 197–203.
[11] H. Alexandra, "A look at how no man's sky's procedural generation works," *Kotaku*, Oct. 2016, accessed: 2020-1-22. [Online]. Available: https://kotaku.com/a-look-at-how-no-mans-skys-procedural-generation-works-1787928446

[12] S. Risi, J. Lehman, D. B. D'Ambrosio, R. Hall, and K. O. Stanley, "Petalz: Search-Based procedural content generation for the casual gamer," *IEEE Trans. Comput. Intell. AI Games*, vol. 8, no. 3, pp. 244–255, Sep. 2016.

[13] B. von Rymon Lipinski, S. Seibt, J. Roth, and D. Abé, "Level graph – incremental procedural generation of indoor levels using minimum spanning trees," in *2019 IEEE Conference on Games (CoG)*, Aug. 2019, pp. 1–7.

[14] R. Craveirinha and L. Roque, "Studying an Author-Oriented approach to procedural content generation through participatory design," in *Entertainment Computing - ICEC 2015*. Springer International Publishing, 2015, pp. 383–390.

[15] A. Alvarez, S. Dahlskog, J. Font, J. Holmberg, and S. Johansson, "Assessing aesthetic criteria in the evolutionary dungeon designer," in *Proceedings of the 13th International Conference on the Foundations of Digital Games*, ser. FDG '18, no. Article 44. New York, NY, USA: Association for Computing Machinery, Aug. 2018, pp. 1–4.

[16] M. Preuss, A. Liapis, and J. Togelius, "Searching for good and diverse game levels," in *2014 IEEE Conference on Computational Intelligence and Games*, Aug. 2014, pp. 1–8.

[17] D. Ashlock, C. Lee, and C. McGuinness, "Search-Based procedural generation of Maze-Like levels," *IEEE Trans. Comput. Intell. AI Games*, vol. 3, no. 3, pp. 260–273, Sep. 2011.

[18] C. McGuinness and D. Ashlock, "Decomposing the level generation problem with tiles," in *2011 IEEE Congress of Evolutionary Computation (CEC)*, Jun. 2011, pp. 849–856.

[19] A. Alvarez and J. Font, "Learning the designer's preferences to drive evolution," in *International Conference on the Applications of Evolutionary Computation (Part of EvoStar)*. Springer, 2020, pp. 431–445.

[20] A. Baldwin, S. Dahlskog, J. M. Font, and J. Holmberg, "Towards pattern-based mixed-initiative dungeon generation," in *Proceedings of the 12th International Conference on the Foundations of Digital Games*, ser. FDG '17, no. Article 74. New York, NY, USA: Association for Computing Machinery, Aug. 2017, pp. 1–10.

[21] V. Braun and V. Clarke, "Reflecting on reflexive thematic analysis," *Qualitative Research in Sport, Exercise and Health*, vol. 11, no. 4, pp. 589–597, Aug. 2019.

[22] J. Togelius, G. N. Yannakakis, K. O. Stanley, and C. Browne, "Search-Based procedural content generation: A taxonomy and survey," *IEEE Trans. Comput. Intell. AI Games*, vol. 3, no. 3, pp. 172–186, Sep. 2011.

[23] D. Loiacono, L. Cardamone, and P. L. Lanzi, "Automatic track generation for High-End racing games using evolutionary computation," *IEEE Trans. Comput. Intell. AI Games*, vol. 3, no. 3, pp. 245–259, Sep. 2011.

[24] V. Valtchanov and J. A. Brown, "Evolving dungeon crawler levels with relative placement," in *Proceedings of the Fifth International C* Conference on Computer Science and Software Engineering*, ser. C3S2E '12. New York, NY, USA: ACM, 2012, pp. 27–35.

[25] A. Zhou, B.-Y. Qu, H. Li, S.-Z. Zhao, P. N. Suganthan, and Q. Zhang, "Multiobjective evolutionary algorithms: A survey of the state of the art," *Swarm and Evolutionary Computation*, vol. 1, no. 1, pp. 32–49, Mar. 2011.

[26] S. O. Kimbrough, G. J. Koehler, M. Lu, and D. H. Wood, "On a Feasible–Infeasible Two-Population (FI-2Pop) genetic algorithm for constrained optimization: Distance tracing and no free lunch," *Eur. J. Oper. Res.*, vol. 190, no. 2, pp. 310–327, Oct. 2008.

[27] K. Deb, *Multi-objective optimization using evolutionary algorithms*. John Wiley & Sons, 2001, vol. 16.

[28] K. Li, R. Wang, T. Zhang, and H. Ishibuchi, "Evolutionary many-objective optimization: A comparative study of the state-of-the-art," *IEEE Access*, vol. 6, pp. 26 194–26 214, 2018.

[29] C. A. C. Coello, G. B. Lamont, D. A. Van Veldhuizen *et al.*, *Evolutionary algorithms for solving multi-objective problems*. Springer, 2007, vol. 5.

[30] D. J. Walker, R. Everson, and J. E. Fieldsend, "Visualizing mutually nondominating solution sets in many-objective optimization," *IEEE Transactions on Evolutionary Computation*, vol. 17, no. 2, pp. 165–184, 2013.

[31] A. Osyczka and S. Krenich, "Evolutionary algorithms for multicriteria optimization with selecting a representative subset of pareto optimal solutions," in *International Conference on Evolutionary Multi-Criterion Optimization*. Springer, 2001, pp. 141–153.

[32] R. Polikar, "Ensemble based systems in decision making," *IEEE Circuits and systems magazine*, vol. 6, no. 3, pp. 21–45, 2006.

[33] C. C. Aggarwal, A. Hinneburg, and D. A. Keim, "On the surprising behavior of distance metrics in high dimensional space," in *International conference on database theory*. Springer, 2001, pp. 420–434.

[34] D. H. Wolpert and W. G. Macready, "No free lunch theorems for optimization," *IEEE Trans. Evol. Comput.*, vol. 1, no. 1, pp. 67–82, Apr. 1997.

[35] S. P. Walton and M. R. Brown, "Predicting effective control parameters for differential evolution using cluster analysis of objective function features," *Journal of Heuristics*, vol. 25, no. 6, pp. 1015–1031, Dec. 2019.

[36] G. N. Yannakakis, A. Liapis, and C. Alexopoulos, "Mixed-initiative co-creativity," in *9th International Conference on the Foundations of Digital Games*. Foundations of Digital Games, 2014.

[37] V. Braun and V. Clarke, "To saturate or not to saturate? questioning data saturation as a useful concept for thematic analysis and sample-size rationales," *Qualitative Research in Sport, Exercise and Health*, pp. 1–16, Dec. 2019.

[38] P. Sampaio, A. Baffa, B. Feijó, and M. Lana, "A fast approach for automatic generation of populated maps with seed and difficulty control," in *2017 16th Brazilian Symposium on Computer Games and Digital Entertainment (SBGames)*, Nov. 2017, pp. 10–18.

[39] A. Alvarez, S. Dahlskog, J. Font, and J. Togelius, "Empowering quality diversity in dungeon design with interactive constrained MAP-Elites," in *2019 IEEE Conference on Games (CoG)*. ieeexplore.ieee.org, Aug. 2019, pp. 1–8.

**Sean P. Walton** received his MPhys degree in Physics from Aberystwyth in 2002. He worked as a school Physics teacher for several years before completing his PhD in numerical methods at Swansea Universities' Zienkiewicz Centre for Computational Engineering in 2013. Currently he works as a senior lecturer in Computer Science at Swansea University. His academic research focus is on using evolutionary optimisation algorithms to support the design process in a number of fields, and investigating game design approaches that are effective for educational games. Outside of academia he is a BAFTA Cymru nominated game designer and founding director of Pill Bug Interactive.

**Alma A. M. Rahat** is a Lecturer in Big Data/Data Science at Swansea University, UK. He has a BEng (Hons) in Electronic Engineering from the University of Southampton, UK, and a PhD in Computer Science from the University of Exeter, UK. He worked as a product development engineer after his bachelor's degree, and held post-doctoral research positions at the University of Exeter. Before moving to Swansea, he was a Lecturer in Computer Science at the University of Plymouth. His current research focus is in the broad areas of fast hybrid optimisation methods, real-world problems and machine learning. In particular, he is developing efficient methods inspired from surrogate-assisted (Bayesian) optimisation for optimising computationally or financially expensive problems (for example, computational fluid dynamics aided design problems).

**James Stovold** received his MEng degree in Computer Science from York in 2012, and completed his PhD at the York Cross-Disciplinary Centre for Systems Analysis (YCCSA) in 2016. After a brief stint in industry, he returned to teaching in 2018. Currently a lecturer at the British University Vietnam, his interests are in distributed cognition, bio-inspired algorithms, swarm intelligence, and robotics.
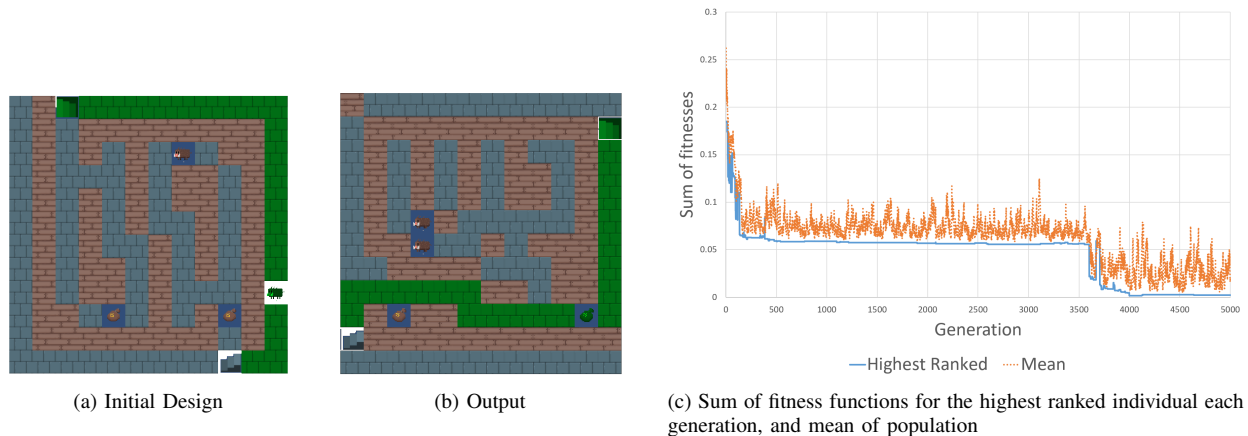
(a) Initial Design



(b) Output



(c) Sum of fitness functions for the highest ranked individual each generation, and mean of population

Fig. 6. Algorithm-Driven Test A.
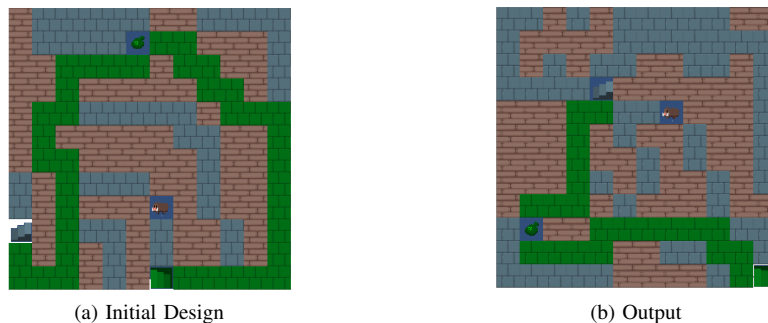


(a) Initial Design



(b) Output

Fig. 7. Algorithm-Driven Test B

## SUPPLEMENTARY MATERIALS

### C. Algorithm-Driven Benchmarks

*1) Methodology:* To test the effectiveness of the GA itself a series of studies were performed using an entirely algorithm driven-approach. For benchmark tests we use six maps presented by Baldwin et al. [6] to show the different styles of maps which could be created by varying the configurations of their approach. Conveniently they represent a range of styles, from maps with no corridors to maps with no chambers. We use these maps as benchmarks to avoid unconscious bias which could result from us designing our own. For each test the target map was entered as the initial user designed level. The GA was then run and the highest ranked map in the population at the end of the optimisation is presented.

*2) Results:* The first-algorithm driven test is one which is made up of corridors with zero chambers. The target map is shown in Figure 6a, and the map created by the GA is shown in Figure 6b. The created map contains only one chamber, is predominately made up of corridors and has the same number of treasure and enemy tiles as the targets. For this first test we have included an optimisation history graph, Figure 6c, constructed by taking the sum of fitnesses for the best individual each generation. It is typical of the behaviour observed in all tests. The results from test B are shown in Figure 7. In this test the target map has a single chamber and many corridors. The resulting output is dominated by corridors, although some of them are unreachable by the player. The output design has a similar ratio of passable to impassible tiles. Both maps have a single treasure and enemy tile. Test C is a map with a comparable number of corridors to chambers. The results of this study are shown in Figure 8. The output has a similar balance of corridors and chambers, and a similar distribution of treasures and enemies. The results for test D are shown in Figure 9. This target design is largely made up of chambers with a few corridors. The GA is capable of matching this distribution. In test E the target map is made up of chambers connected by single tile corridors. The results of this test are shown in Figure 10. Much like the target the output is made up of chambers and single tile corridors with the same number of treasure and enemy tiles. The final test, F, is simply a map with zero wall tiles. Figure 11 shows that the GA handles this edge case. Also notice that the path length is almost the same in both.

### D. Final Designs of Participants

The following images show the final screen after participants had completed their task. As a reminder, the top 5 maps are the participants final 5 levels. The larger maps are the suggestions from the system which may have been modified by the user
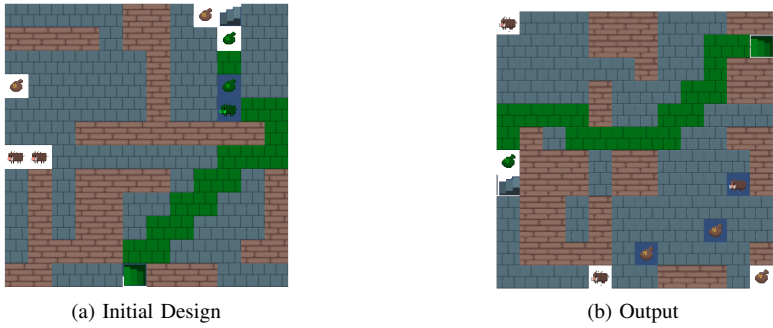
(a) Initial Design

(b) Output

Fig. 8. Algorithm-Driven Test C



(a) Initial Design

(b) Output

Fig. 9. Algorithm-Driven Test D



(a) Initial Design

(b) Output

Fig. 10. Algorithm-Driven Test E



(a) Initial Design
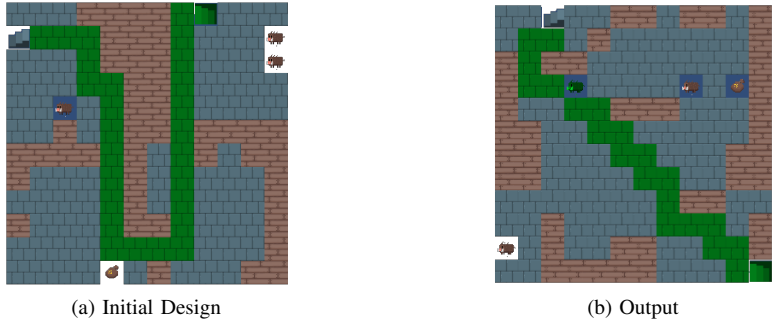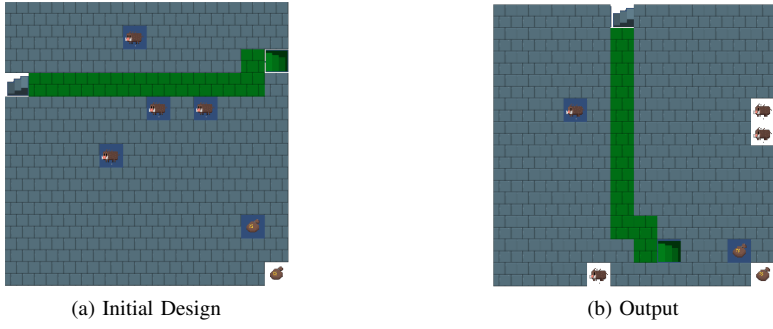
(b) Output

Fig. 11. Algorithm-Driven Test F

before this image was taken.

*1) Control Group:* Figures 12 to 18 show the final screens for participants in the control group.

*2) GA Group:* Figures 19 to 29 show the final screens for participants in the genetic algorithm group.

Fig. 12.  Control Participant A



Fig. 13.  Control Participant B



Fig. 14.  Control Participant C



Fig. 15.  Control Participant D

Fig. 16.  Control Participant E



Fig. 17.  Control Participant F



Fig. 18.  Control Participant G
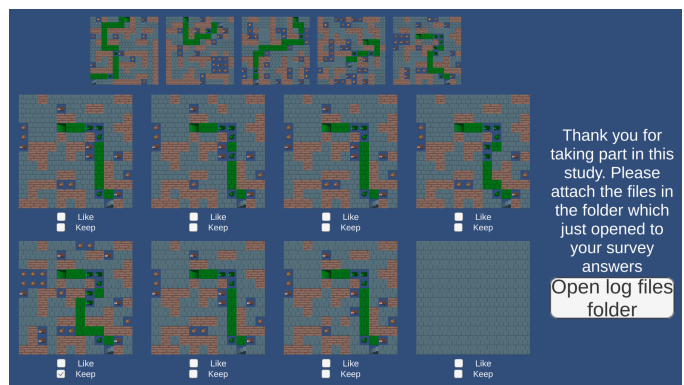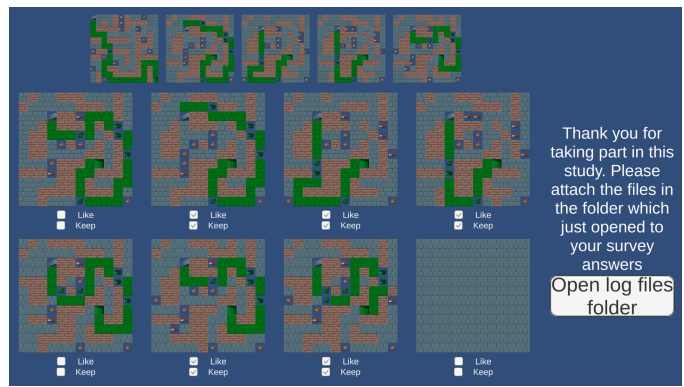


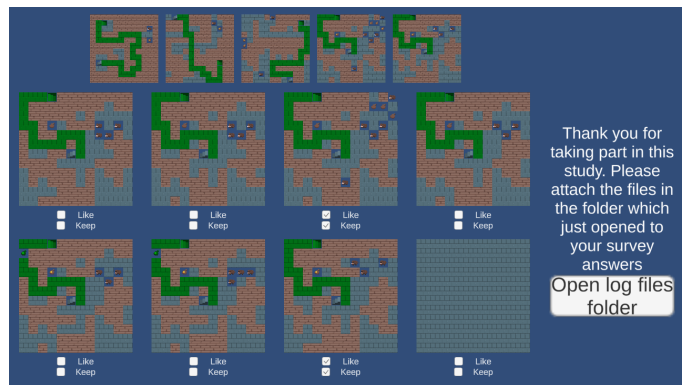Fig. 19.  Genetic Algorithm Participant A

Fig. 20. Genetic Algorithm Participant B



Fig. 21. Genetic Algorithm Participant C



Fig. 22. Genetic Algorithm Participant D



Fig. 23. Genetic Algorithm Participant E

Fig. 24.  Genetic Algorithm Participant F



Fig. 25.  Genetic Algorithm Participant G


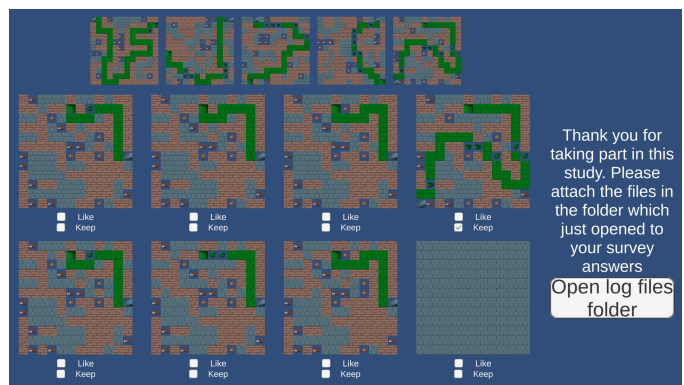
Fig. 26.  Genetic Algorithm Participant H



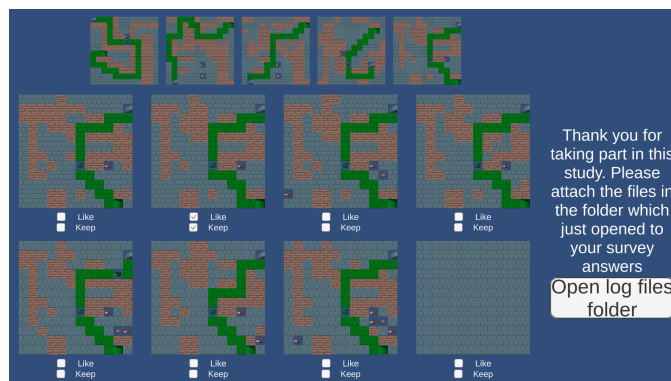Fig. 27.  Genetic Algorithm Participant I

Fig. 28.  Genetic Algorithm Participant J



Fig. 29.  Genetic Algorithm Participant K