

Measuring Quality of Grammars for Procedural Level Generation

Riemer van Rozen
Amsterdam University of Applied Sciences
Amsterdam, The Netherlands
r.a.van.rozen@hva.nl

Quinten Heijn
University of Amsterdam
Amsterdam, The Netherlands
samuel.heijn@gmail.com

ABSTRACT

Grammar-based procedural level generation raises the productivity of level designers for games such as dungeon crawl and platform games. However, the improved productivity comes at cost of level quality assurance. Authoring, improving and maintaining grammars is difficult because it is hard to predict how each grammar rule impacts the overall level quality, and tool support is lacking. We propose a novel metric called Metric of Added Detail (MAD) that indicates if a rule adds or removes detail with respect to its phase in the transformation pipeline, and Specification Analysis Reporting (SAnR) for expressing level properties and analyzing how qualities evolve in level generation histories. We demonstrate MAD and SAnR using a prototype of a level generator called Ludoscope Lite. Our preliminary results show that problematic rules tend to break SAnR properties and that MAD intuitively raises flags. MAD and SAnR augment existing approaches, and can ultimately help designers make better levels and level generators.

CCS CONCEPTS

•**Software and its engineering** → *Integrated and visual development environments; Design languages; Domain specific languages;*
•**Computing methodologies** → *Generative and developmental approaches; Applied computing* → *Computer games;*

KEYWORDS

game development, level design, PCG, automated game design, grammars, quality, metrics, domain-specific languages

ACM Reference format:

Riemer van Rozen and Quinten Heijn. 2018. Measuring Quality of Grammars for Procedural Level Generation. In *Proceedings of Foundations of Digital Games 2018, Malmö, Sweden, August 7–10, 2018 (FDG '18)*, 8 pages. DOI: 10.1145/3235765.3235821

1 INTRODUCTION

Grammar-based level generation is a form of Procedural Content Generation (PCG) that raises the productivity of game level designers. Instead of hand-crafting levels, designers create a level transformation pipeline that generates levels for them by authoring modules, grammars and rewrite rules. The grammar rules work on data

This research is carried out at the SWAT group of CWI, Amsterdam, the Netherlands. Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.
FDG '18, Malmö, Sweden

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM. 978-1-4503-6571-0/18/08...\$15.00
DOI: 10.1145/3235765.3235821

structures such as strings, tile maps and graphs, which can be used for generating names, level layouts and missions. These artifacts are step-by-step transformed and combined until a final detailed and fully populated level is generated, with missions, power-ups, challenges, enemies, hidden treasures, secret pathways, encounters, etc. Ideally, each generated level has the intended qualities.

Unfortunately, improving the productivity of level designers comes at the cost of quality assurance. In practice, many small problems arise, such as levers in walls, blocked pathways, missing encounters and lava adjacent to water. A lack of direct manipulation compromises the ability of designers to isolate and improve level qualities, e.g., when authoring bridges, forests or paths. As a result, some generated levels may lack intended goals, challenges and missions.

The qualities of generated levels depend on the composition of grammar rules and how they are combined in sequence. Therefore, potential bugs often remain unknown until they are observed during playtesting. Additionally, the combinatorial explosion resulting from recursive rule expansions complicates forming mental models required for reasoning about intended qualities, and how they are represented in the grammar or intermediate data. Moreover, it is hard to predict how individual rules affect the overall level quality.

Grammars are brittle, i.e. code that is liable to break easily. Designers require special measures to ensure that qualities once introduced, remain intact, preventing successive rewrites from breaking levels. Fixing one level with a rule that prevents an occurrence may introduce new problems in others. In general, there is a lack of tools and techniques for authoring, debugging, testing and improving rules that introduce and preserve design intent. As a result, the full potential of these techniques has not yet been realized.

We aim to improve the quality of grammar-based procedural level generation in general, and focus on grammars that work on tile maps in particular. We motivate our research by studying and improving Ludoscope, a state-of-the-art development environment for generating very diverse game levels. Since its inception, Ludoscope was developed by Ludomotion for indie game development, and successfully applied to a rogue-like dungeon crawler called Unexplored. We address the need of developers for better tools. This paper proposes and contributes two enabling techniques:

- (1) Metric of Added Detail (MAD), a novel metric that indicates if a grammar rule adds or removes detail to a tile map. We hypothesize that grammars gradually add detail. MAD leverages a detail hierarchy, a binary relation on alphabet symbols indicating which symbol is more detailed, which can easily be derived from transformation pipelines.
- (2) Specification Analysis Reporting (SAnR), a technique that offers a level property language for expressing level qualities. SAnR analyzes and reports how these properties evolve over time in level generation histories.

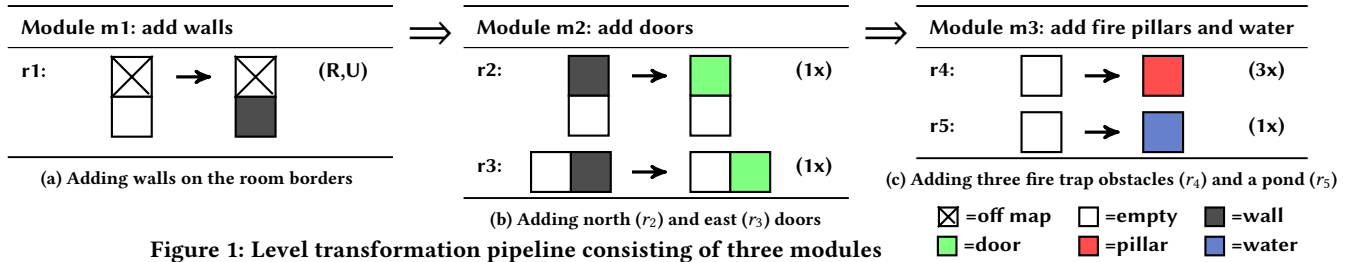


Figure 1: Level transformation pipeline consisting of three modules

We demonstrate the feasibility of MAD and SAnR by implementing Ludoscope Lite (LL), a light-weight version of Ludoscope intended to study level quality. LL is implemented using RASCAL, a meta-programming language and language workbench. Our preliminary evaluation shows that SAnR can express and analyze simple level properties, and that MAD raises flags for rules that remove detail. MAD and SAnR augment existing approaches by supporting gradually adding detail and analyzing level generation histories, which ultimately helps designers make better levels and level generators.

2 RELATED WORK

Evaluating content generators and their output is a key open research problem [12, 15, 19]. Generators can be analyzed in terms of generated content, e.g., Summerville *et al.* evaluate metrics for difficulty, visual aesthetics and enjoyment of platform games [15].

We take an authoring perspective on level grammars. Our approach stands apart by also taking into account how generated levels are generated. This enables level designers to relate qualities of generated levels back to the source code of the generator (grammar rules) and make targeted improvements.

Level grammars are under-specified, since they also generate many levels that are bad with respect to design constraints. The challenge is authoring a set of rules that efficiently generates varied and well-structured results capturing design intent while limiting the recursion. Smith and Mateas propose explicitly describing design spaces as an answer set programs, and show generators can be sculpted for a variety of content domains [13]. Van der Linden *et al.* focus on improving authoring and controlling level generators by expressing gameplay-related design constraints. They use graph grammars to encode these constraints, and generate action graphs that associate player actions and content for generating complete layouts of game levels [16]. We refer to a survey of van der Linden *et al.* for a wider discussion on techniques for procedural dungeon generation [17].

We relate our work to other content generators that use grammars. Tracery is a grammar-based tool for authoring stories and art as structured strings that has been used for generating names, descriptions, stories in poetry, Twitter bots and games [1, 2]. PuzzleScript is a language and authoring environment which uses rewrite rules to express puzzle mechanics [10]. Ludoscope is a visual environment for authoring level transformation pipelines as grammars that builds upon the mission and spaces framework [3, 4]. Pipelines consist of modules that contain grammars, alphabets and recipes that transform level artifacts such as strings, tile maps, graphs and Voronoi diagrams. In particular, recipes are crucial to control the generation and focus the application of rules for obtaining

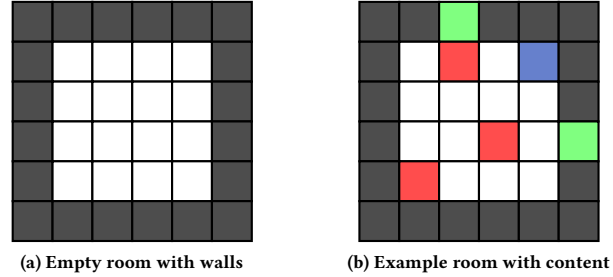


Figure 2: Tile maps that are input and output of the pipeline

aesthetically pleasing levels. Recipes parameterize modules with instructions, that determine the ordering of rules and limit how often rules work. Member values annotate tiles with extra information. Both help reducing the generation space, but neither are well-suited to check qualities off-line and independently. Ludoscope is neither extensively documented, nor currently available as open source software. Karavolos *et al.* report experiences on applying Ludoscope to a platformer and a dungeon crawl game, which require very different transformation pipelines [7]. Our approach closely follows the pipeline structure of Ludoscope, but it improves upon its capabilities for analyzing grammar and level quality.

3 GRAMMARS FOR LEVEL GENERATION

Here, we introduce quality issues in grammar-based level design using a simple example that generates a room for a dungeon crawler, which illustrates some of the challenges that arise during authoring grammars. It isolates problems that have larger more complex forms in practice, e.g., in Unexplored. We relate questions designers might have in Section 3.2 to technical challenges in Section 3.3.

3.1 Introductory Example

In dungeon crawlers, tile maps often represent rooms connected by pathways. Our level generation pipeline, shown in Fig. 1, generates rooms with two doors connecting to a larger dungeon. It consists of three modules of grammar rules that represent sequential level transformation phases. The grammar rules rewrite pieces of the tile map matched by their pattern on their left hand side to the pattern on their right hand side. The pipeline takes an empty tile map as input, e.g., of 6x6 tiles. Each phase randomly selects and applies rules, gradually adding detail. Many levels can result, and as we will see, not all of these are what a designer might deem desirable.

First, module m_1 adds walls on the borders of the tile map (Fig. 1a). It contains one rule called r_1 , whose left hand pattern matches on an empty tile on the north edge of the map. Grammar rule r_1 replaces

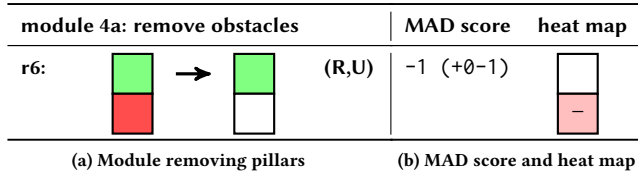


Figure 3: Module for removing pillars that block doors

an empty tile on the north edge of the map with a wall. Rules can have modifier symbols to its right. The (U) symbol to the right indicates that rule r_1 is applied as many times as possible. The (R) symbol indicates that rule r_1 is also applied to the east, south and west borders of the map. The result of module m_1 is always a tile map with walls on its borders, e.g., Fig. 2a is the output at 6x6.

Next, module m_2 adds doors in the north and east walls that connect the room to other parts of the dungeon (Fig. 1b). The rules r_2 and r_3 respectively add a door in the north and east walls. These rules are applied exactly once (1x). Finally, module m_3 introduces challenge (Fig. 1c). Rule r_4 places three fire pillars, traps that set players on fire if they remain close too long. In addition, rule r_5 adds a pool of water the player can use to extinguish the flames.

3.2 Level designer questions

The pipeline of Fig. 1 can also generate problematic levels. For instance, in Fig. 2b, a fire pillar in front of the north door prevents players from passing. One way to fix this is to *patch* the level by removing obstacles, as shown in Fig. 3a results in Fig. 5a. However, fewer pillars than intended may reduce the difficulty. Another way is moving obstacles away from doors, as shown in Fig. 4a, which results in Fig. 5b. Unfortunately, other problematic output still exists, e.g., Fig. 6. Authoring level grammars is hard, even for this tiny example. Questions about quality a designer might have are:

- (1) **Efficiency.** Do the grammar rules efficiently generate levels, or is time wasted on overwritten dead content?
- (2) **Effectiveness.** Do the grammar rules effectively generate levels that contain all the intended objects, composite structures, problems and solutions, or are some parts missing?
- (3) **Root-cause analysis.** Given a level with a problem, by which rules were the affected tiles generated?
- (4) **Bug-fixing.** Does changing a rule improve levels, or does it also introduce new problems?
- (5) **Bug-free.** How can unwanted situations be prevented and removed from the level generation space?

Other relevant questions not further discussed here are, e.g.,

- **Playability.** Are the challenges of all generated levels solvable, or are there ways in which players can get stuck?
- **Challenge.** Are the levels challenging to play?

3.3 Challenges

Here, we identify technical challenges that need to be addressed for answering questions of level designers described in Section 3.2.

- (1) **Static analysis and metrics.** Profiling the applications of rules helps to assess efficiency measuring (relative) times and amounts. However, static analysis may also help predict rule efficiency. Upper bounds on rule applications

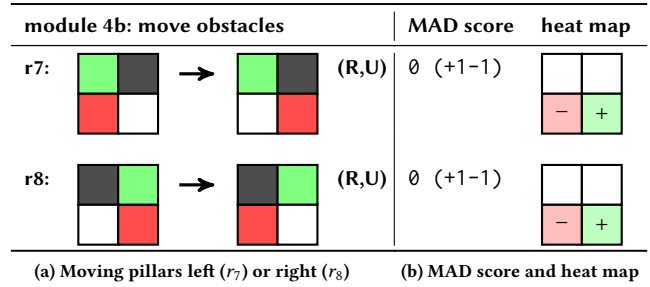


Figure 4: Module for moving pillars that block doors

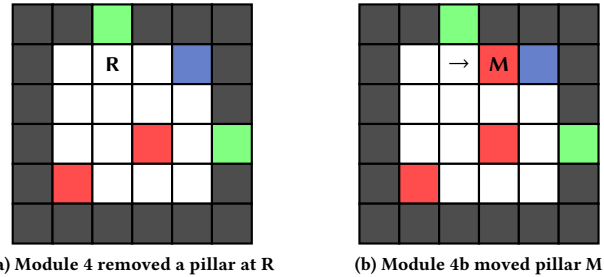


Figure 5: Repairing the example level of Fig. 2b in two ways

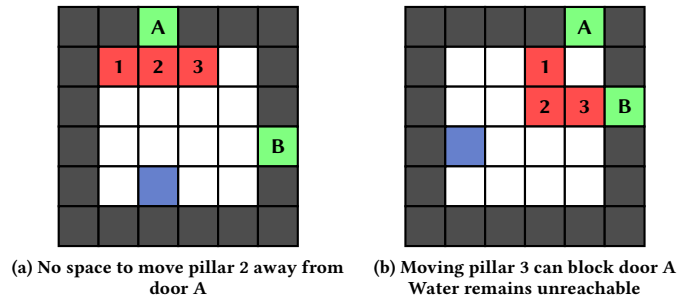


Figure 6: Levels that cannot be repaired by Module 4b

enable reasoning about worst-case scenarios. Left hand patterns that can never match indicate dead code. In addition, metrics can help assess to which extent rules contribute to generating an intended result, to find bad rules.

- (2) **Analyzing the level generation space.** Viewed as a state-space exploration problem, rules might rewrite levels to prior states. For a given level, the shorter its trace of rewrites, the more efficient its generation.
- (3) **Expressing and analyzing level qualities.** Grammar rules lack ways to specify properties at specific points in the pipeline, e.g., if objects are (not) adjacent, contained, intact or missing. Designers need an additional formalism for effectively specifying properties that intuitively capture design intent. To see how qualities evolve, levels can be checked against these properties after each transformation.
- (4) **History analysis.** Generators produce tile maps by applying grammar rules in sequence, e.g. Fig. 7. However, these generation histories are usually not stored. For identifying rules that impact tiles, or groups of tiles, designers require an analysis of the level transformation history.

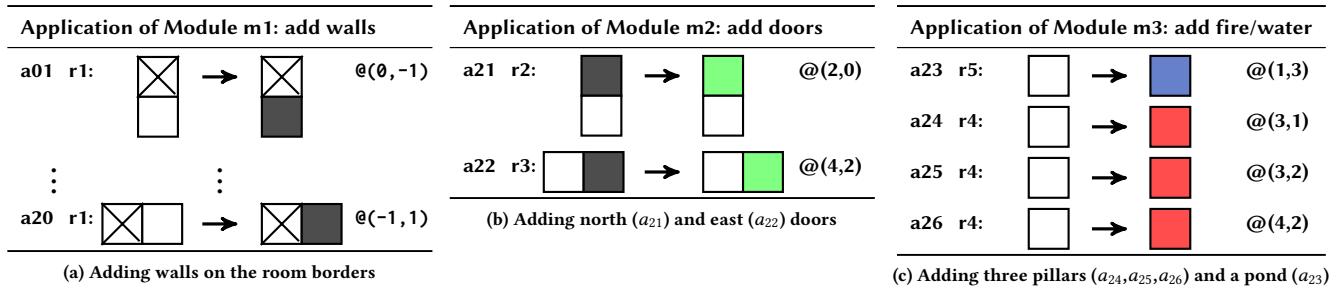


Figure 7: Level generation history showing how rules generated the example level shown in Fig. 2b

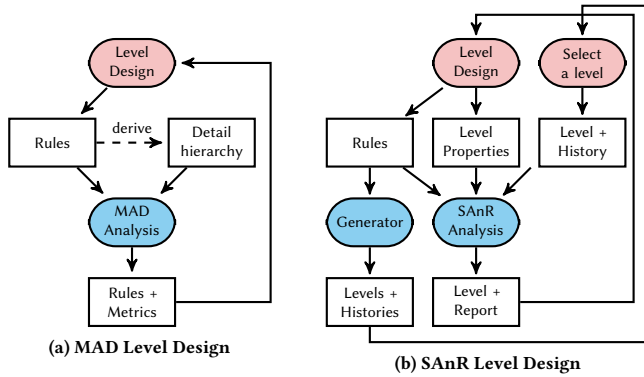


Figure 8: Producing MAD and SANr level design artifacts

- (5) **Impact analysis.** Assessing the impact of rules on many generated results requires isolating rule effects. The position in the pipeline scopes the locality of impact, and a dependency analysis can exclude side-effects, but an exhaustive impact analysis requires generating examples.
- (6) **Test Automation.** Testing the impact of changes on all possible levels is not feasible. As a result, levels may exist that contain bugs. The challenge is devising a test harness that generates representative levels for finding bugs.
- (7) **Debugging.** Identifying and fixing bugs requires appropriate views and tools for setting break points and making modifications, e.g., selecting one or more adjacent tiles to filter and analyze selected properties.

4 GRAMMAR ANALYSIS AND DEBUGGING

We approach the challenges of Section 3.3 from a software evolution perspective. We propose two solutions, Metric of Added Detail (MAD) and Specification Analysis Reporting (SANr). Fig. 8 schematically shows how designer activities and algorithmic processes (respectively shown as pink and blue rounded rectangles) produce (outgoing arrows) and consume (incoming arrows) artifacts (rectangles). The field of software evolution studies how software evolves over time [11]. As software ages, it conforms less and less to the changing expectations of its users. In addition, for developers it also becomes harder over time to adjust software and maintain its quality. Research includes methods and techniques for analyzing source code and for making changes to improve the software

```

1  module util::mad::Metric
2  alias Detail = rel[str greaterSymbol, str lesserSymbol];
3  alias Rule = lrel[str lhs, str rhs];
4  alias RuleScore = lrel[str lhs, str rhs, int score];
5
6  RuleScore getRuleScore(Rule r, Detail d)
7    = [<lhs, rhs, getTileScore(lhs, rhs, d)> | <lhs, rhs> ←r];
8
9  int getTileScore(str lhs, str rhs, Detail d) { //rewriting a tile
10   if(<lhs,rhs> in d) return -1; //removes detail
11   else if(<rhs,lhs> in d) return 1; //adds detail
12   else return 0; //retains detail
13 }

```

Figure 9: Metric of Added Detail as a RASCAL program

quality. Since game requirements are mainly non-functional and evolve rapidly, these techniques are also vital for game quality.

4.1 Metric of Added Detail

Metrics have been proposed to analyze how changes to source code impact software quality. Volume (or size) can be measured by counting Lines Of Code (LOC), and branch points in the control flow of methods can be measured using Cyclometric Complexity (CC). At any moment, metrics are just abstract values, but when studied over time they can provide insight into phenomena and quality, in particular when developers have questions regarding the effect of maintenance and new requirements that require programming. Heitlager *et al.* describe a software maintainability model [6], which requires that measures are 1) technology independent; 2) simply defined; 3) easy to understand and explain; and 4) enablers of root cause-analysis, relating source code properties to system qualities.

Here we introduce the Metric of Added Detail (MAD), a simple metric for grammars operating on tile maps, which is easy to explain and understand. MAD does not directly predict level quality, but instead measures the effect on detail of individual rules by leveraging the assumption that details are gradually added (Fig. 8a).

We define MAD in Fig. 9, using the concise functional notation of RASCAL. MAD requires a *detail hierarchy*, represented as a binary relation on grammar symbols (line 2). Rules are represented as lists of tuples of source and target symbols that abstract from tile map dimensions (line 3). The result of the metric adds a *score* element to each tuple that records if detail is added (score +1), removed

	a_{01}	...	a_{21}	a_{22}	a_{23}	a_{24}	a_{25}	a_{26}	a_{27} (alternative)	a_{27} (alternative)
2x door in walls //m1: add walls				✓						
1x water				✓						
3x pillar							✓			
no pillar adjacent to door	✓				✗				✗	
no water adjacent to pillar	✓								✓	✗

(a) Level Property Language specification encoding level qualities (b) Level generation report showing how a level (Fig. 2b) evolved over time (Fig. 7) (c) Certain result of Module 4a (d) Possible result of Module 4b

Figure 10: Level properties and a level generation report

(score -1) or persisted (score 0) (line 4). The function *getRuleScore* specifies the rule metric as a list comprehension (lines 6–7). Given a rule and a detail hierarchy, it calculates for each symbol on the left hand side if the right hand side adds or removes detail using the function *getTileScore* (lines 9–13). Displayed as a heat map, the result is aggregated as a sum of tile detail scores.

4.2 Deriving Detail Hierarchies

MAD is tool independent and rule parametric, but it requires a detail hierarchy, which needs to be derived. Modules imply a natural hierarchy for tools that use level transformation pipelines, each phase introducing symbols that are more detailed than the last. Using this approach, we derive the following detail hierarchy for the example of Section 3.1 Fig. 1 {water, pillar} > door > wall > empty, or visually {■, ■} > ■ > ■ > ■.

Competing non-deterministic rules do not sequentially add detail, e.g., r_4 or r_5 adds ■ or ■ first. Therefore, deriving a symbol hierarchy for exposing data generated and overwritten within a module is less straightforward. We see the following alternatives:

- (1) Allow an explicit user-defined detail hierarchy, or derive it from an explicit rule ordering such as a Ludoscope recipe.
- (2) Assume detail is sequential to the rules in the module.
- (3) Add the inverse to the relation for symbols with the same rank in the hierarchy, e.g., ■ > ■ and ■ > ■. However, this is not very intuitive.

4.3 Analyzing Rules with MAD

Using the detail hierarchy derived in Section 4.2 we calculate MAD scores for rules of modules $m4a$ and $m4b$ intended to fix broken levels, shown in Fig. 3a and Fig. 4a. Rule r_6 , which removes fire pillars, has a negative effect on detail, as shown in Fig. 3b. The effect of rules r_7 and r_8 that instead move them, shown in Fig. 4b, is neutral. MAD helps designers assess if rules contribute to generating intended results, and augments intuitions with facts. Rules that remove details may be fixes, but may also cause dead content or regressions in the level generation space that waste time.

4.4 Expressing and Analyzing Level Properties

Here we address the challenges of expressing and analyzing level qualities from a Software Language Engineering perspective [9]. We propose Specification Analysis Reporting (SAnR), a technique for analyzing level grammars against level properties. In the mixed-initiative design process shown in Fig. 8b, designers author a grammar (rules and modules) and SAnR level properties, a generator

```

1  start syntax LevelSpec = spec: Property*;
2  syntax Property = property: Condition TileSet;
3  syntax Condition //required condition
4    = none: "no" //tile set is empty
5    | count: INT size "x"; //tile set size is
6  syntax TileSet //defines a set of tiles (now visible)
7    = tileSet: ID fileName FilterNow FilterWhere;
8  syntax FilterNow //filters the tile set (now visible)
9    = nowAny: //empty alternative, no filter
10   | nowAdjacent: "adjacent to" ID fileName;
11 syntax FilterWhere //filters a tile set (historically)
12   = everAny: //empty alternative, no filter
13   | everRule: "in" ID ruleName; //topological locations

```

Figure 11: Syntax of Level Property Language in RASCAL

generates levels, and the designers selects one level to analyze, for which SAnR generates a report.

SAnR provides a property notation. This is a so-called Domain-Specific Language (DSL), a language that offers appropriate notations and abstractions with expressive power and affordances over a particular problem domain [18], in this case specifying properties of tile maps as correct outcomes of tile map transformations.

We show its syntax in Fig. 11, and give an informal description of its language semantics. Instead of writing new grammar rules, a SAnR level specification is a set of declarative properties, which refer to names used in the grammar (line 1). Given a level history as a sequence of rule-based model transformations, e.g., Fig. 7, properties can be evaluated at each point in time, yielding either *true* or *false*. Properties work on tile locations, places on tile maps specified by x and y coordinates denoted as $@(x,y)$, the top left tile being $@(0,0)$. A property is a condition on a set of tile locations visible on a tile map (line 2), which must either be empty (line 4) or of a specific size (line 5). The set is built by collecting tile locations using names from the grammar alphabet, e.g., “door” retrieves a set containing each location of a door. On the example of Fig. 2b this yields $\{ @(2,0), @(4,2) \}$, which means “2x door” is true and “1x door” is false. Locations can be filtered in two optional ways.

- (1) **Adjacency.** The **adjacent to** keyword (lines 8-10), filters locations that do not share at least one side with tiles of another kind, e.g., “door **adjacent to** pillar”, denotes a set of locations of door tiles next to at least one pillar.
- (2) **Topography.** The **in** keyword (lines 11-13), filters out locations that were never affected by a rule rewrite. In other

words, we use rule names to collect sets of tile locations from the level generation history as “*topographical regions*”. The resulting set is the intersection between the left and right hand operands. For example “door **in** walls” gives the set of *door* locations in the region affected by rule *walls*.

4.5 Analyzing Level Generation Histories

The SAnR analysis uses properties for generating level generation reports that show when properties were valid, and when they became invalid. For example, given the level generation history of Fig. 7, and the properties of Fig. 10a, SAnR evaluates the properties after each transformation step, yielding the report of Fig. 10b. From the report we read that at step a_{24} transformation r_4 : $\square \rightarrow \blacksquare @ (3,1)$ places a pillar in front of the north door, which invalidates the property “no pillar adjacent to door”.

4.6 Analyzing Rule Impact

SAnR can also be used to analyze the impact of new rules on existing levels with respect to level properties. For instance, we can spot problems at alternative steps a_{27} in the report of Fig. 10 caused by modules $m4a$ and $m4b$ intended as fixes, shown in Fig. 3a and Fig. 4a. On the one hand, Fig. 10c shows that when module $m4a$ removes the pillar with transformation r_6 : $\blacksquare \rightarrow \square @ (3,0)$ this breaks the property “3x pillar”. On the other hand, Fig. 10d shows that when module $m4b$ moves the pillar to the east with transformation r_7 : $\blacksquare \rightarrow \blacksquare @ (3,0)$ this breaks the property “no pillar adjacent to water”.

5 PRELIMINARY EVALUATION

Here, we report on a preliminary evaluation of the use of MAD and SAnR in the implementation of a prototype level generator called Ludoscope Lite.

5.1 Implementation of Ludoscope Lite

Ludoscope Lite (LL) is a light weight version of Ludoscope intended for rapid prototyping, research and experimentation with analysis and generation techniques for making better grammar-based game levels and generators. Its focus is initially on designing and validating approaches for tile maps, which are later implemented and applied in Ludoscope. We use language work bench [5] and meta-programming language RASCAL¹ [8] to implement MAD and SAnR as separate reusable modules and integrate both in LL².

Table. 1 gives an overview of the components of LL and their size in Lines of Code (LOC) relative to Ludoscope. Of course, the user-friendly IDE of Ludoscope has many features LL lacks, explaining the size difference. LL integrates a grammar-based parser that reads the storage format of Ludoscope. The ultimate goal is compatibility, sharing syntax

¹<https://www.rascal-mpl.org>

²<https://github.com/visknut/LudoscopeLite>

Component	Ludoscope (KLOC)	LL (KLOC)
IDE (features differ)	10.5	0.3
Parser + execution	10	1.7 + 0.4
Test + test data	?	1.5 + 0.7
Metric of Added Detail	not yet	0.1
Level Property Language	not yet	0.3
Extension wrappers	-	0.4
Total	20.5	5.5

Table 1: Source code size of Ludoscope and Ludoscope Lite

Data	Example	+ $m4a$	+ $m4b$
Unique histories	9846	9858	9844
Unique tile maps	9171	9014	8775
Broken tile maps	6254	6132	4613
Bugs found	2	2	4

Table 2: SAnR data on the example the pipeline of Fig. 1 and its two extensions modules $m4a$ and $m4b$

Property	Example	+ $m4a$	+ $m4b$
2x door in walls	-	-	-
1x water	-	-	-
3x pillar	-	r6 (3226x)	r7 (111x) r8 (112x)
no pillar adjacent to door	r5 (3164x)	-	r5 (438x)
no water adjacent to pillar	r5 (5686x)	r5 (5209x)	r5 (5482x)

Table 3: SAnR level generation reports for 10K random executions. The rules r_n refer to Fig. 1, Fig. 3a and Fig. 4a

and semantics for generating and analyzing rules. We apply test-driven development, encoding expected behaviors for most of its features in a combination of unit and integration tests for regression testing. The histories and reports shown in this paper are generated by LL, which currently still generates them as strings. A more user friendly visualization is work in progress.

5.2 Test Automation

We use LL to evaluate SAnR on the running example of Section 3.1³ We wish to learn if LL and SAnR can help automate tests, and run 10K random executions (or simulations) on the pipeline Fig. 1, and its extensions, shown in Fig. 3a and Fig. 4a, which makes 30K executions total. For each execution, we record the model transformation history and use SAnR and the properties of Fig. 10a to obtain a report.

³There is one difference, LL implements Ludoscope recipes for limiting the amount of times that rules are applied. As a side-effect, this limits sequences and reduces the level generation space.

Table. 2 displays an overview of the results, which were obtained in about 10 minutes of run time. The unique number of histories is lower than 10K because some executions yielded the same transformations. In addition, different transformation sequences can produce the same tile map, which explains why there are fewer unique tile maps. We consider a tile map *broken* when not all SAnR property are satisfied. In addition, Table. 3 shows which rules break properties (in how many histories) for each pipeline version, which helps designers compare and analyze causes.

We gain the following insights. The test automation approach is feasible, and issues can be found in seconds. In addition, by relating the number of unique outputs to the number of broken outputs we can get an idea how serious issues are. Naturally, 10K random executions says nothing about test coverage, but it improves upon random manual testing. We confirm that module *m4a* is a bad fix. We note that although extension *m4b* increases the number bugs, it also generates fewer broken tile maps. Clearly, the pipeline still requires fixes. Of course, the example is small and not representative of the size and complexity of transformation pipelines of games such as Unexplored. However, our test automation setup is reusable, and enables testing other grammars with larger pipelines too.

6 DISCUSSION

MAD and SAnR provide a means for answering designer questions of Section 3.2. Here we discuss the benefits and limitations of the approaches and threats to validity.

6.1 MAD Level Design

MAD gives a partial answer to the question if rules generate levels efficiently. The metric helps designers identify rules that remove detail, and possibly waste time on generating cause dead content. It supports the single responsibility principle, exposing modules add many details at once. However, MAD does not address the challenge of analyzing the state space. At best, it can help identify rules that may lead to longer level generation traces. In addition, we do not know if MAD can be used for data structures other than tile maps, e.g., for grammars that work on graphs. Finally, MAD is not yet empirically validated.

6.2 SAnR Level Design

SAnR properties enable analyzing how effectively rules generate intended levels, e.g. for simple tile adjacency, counting, missing tiles, and topographical inclusion. Properties depend only on the names of rules and tiles, which separates concerns but complicates refactoring grammar rules. SAnR analyzes levels by checking properties against generation histories, and assumes these are correctly generated. Therefore, SAnR reports are only as good as the grammar engine,

which may also contain bugs. Of course, our approach is not the first that checks simple invariant conditions. However, to the best of our knowledge, checking properties that use level generation histories and grammar rule names to collect topographical regions of tile locations is new.

SAnR can help designers analyze quality and remove unwanted situations from the level generation space by identifying transformations and rules that break properties. However, those rules may not be the root cause of the problem, which can originate earlier in the pipeline. In addition, it is hard for developers to analyze the history, since it is not clear where the branch points in the generation process are, and how alternatives would have played out. Finally, the expressive range of properties is currently still rather limited, and a formal semantics relating properties and histories is not yet defined.

7 CONCLUSION

This paper proposes two novel techniques that aim to improve the quality of grammar-based procedural level generation for grammars that work on tile maps. The first, is the Metric of Added Detail (MAD), a novel metric that indicates if a grammar rule adds or removes detail to a tile map. The second, is Specification Analysis Reporting (SAnR), a technique that offers level property language for expressing level qualities. SAnR analyzes and reports how these properties evolve over time in level generation histories. We demonstrated the feasibility of MAD and SAnR with LudoScope Lite, a light-weight version of Ludoscope intended to study level quality. Our preliminary evaluation shows that SAnR can express and analyze simple level properties, and that MAD is intuitive and raises flags for rules that remove detail. In addition, SAnR can be used in test automation. MAD and SAnR augment existing approaches by supporting gradually adding detail and analyzing level generation histories, which ultimately helps designers make better levels and level generators. Of course, LL is an academic research prototype that is not yet extensively validated in practice.

7.1 Future Work

Future work includes the following.

- **Validation.** A case study on Boulder Dash is current work. We also plan to study Unexplored to identify which additional SAnR property features are needed to express design intent more fully, e.g., better filters, validity ranges, and for shapes, paths and relative positions. We hope to identify bugs that would otherwise be hard or impossible to find.
- **Analyses.** Additional analyses on rule dependencies, and partial orderings may be identified of different rule orders generating the same levels, e.g., for increasing test coverage and level generation variety.

For assessing the variety of generated content, existing metrics can be reused. For instance, Smith and Whitehead assess the expressive range of a generator by comparing metrics for linearity and leniency of platform levels [14].

- **Generation.** Here we use SAnR for analyzing level generation histories after they are generated. However, by integrating SAnR into a level generator we could also prune the search space and filter out potential unwanted levels before they are ever produced. A feasibility study can assess the impact on efficiency and scalability of this approach.
- **Formal semantics.** Reproducible dynamic analyses require a formal semantics for the execution of generative grammars, separate from tools and games that interpret them.
- **Parsing.** We observe that ambiguous grammars for parsing and level grammars generating the same tile map with different rule orderings are related. Given a bugged tile map, how many different rule orderings can reproduce it? When changing the rules, can the new rules produce the tile map with a different generation history?
- **Debugging.** Debugging level grammars requires an interactive debugger, in particular for back in time debugging, exploring what-if scenarios and saving and replaying generated levels while testing new rules. Additional visualizations are needed to see how the generation space unfolds.

ACKNOWLEDGMENTS

We thank Paul Klint, Anders Bouwer, Rafael Bidarra and the anonymous reviewers for their insightful comments that helped improve this paper. We thank Joris Dormans for collaborating with us and answering our many questions about the design of Ludoscope.

REFERENCES

- [1] Kate Compton. 2015. Tracery. <http://tracery.io>. (2015). Visited May 19th 2018.
- [2] Kate Compton, Ben Kybartas, and Michael Mateas. 2015. Tracery: An Author-Focused Generative Text Tool. In *Interactive Storytelling*, Henrik Schoenau-Fog, Luis Emilio Bruni, Sandy Louchart, and Sarune Baceviciute (Eds.). Springer International Publishing, 154–161.
- [3] Joris Dormans. 2010. Adventures in Level Design: Generating Missions and Spaces for Action Adventure Games. In *Proceedings of the 2010 workshop on procedural content generation in games*. ACM, 1.
- [4] Joris Dormans. 2011. Level Design as Model Transformation: a Strategy for Automated Content Generation. In *Proceedings of the 2nd International Workshop on Procedural Content Generation in Games*. ACM, 2.
- [5] Sebastian Erdweg, Tijs Van Der Storm, Markus Völter, Meinte Boersma, Remi Bosman, William R Cook, Albert Gerritsen, Angelo Hulshout, Steven Kelly, Alex Loh, and others. 2013. The State of the Art in Language Workbenches. In *International Conference on Software Language Engineering*. Springer, 197–217.
- [6] Ilja Heitlager, Tobias Kuipers, and Joost Visser. 2007. A Practical Model for Measuring Maintainability. In *Quality of Information and Communications Technology, 2007. QUATIC 2007. 6th International Conference on the*. IEEE, 30–39.
- [7] Daniël Karavolos, Anders Bouwer, and Rafael Bidarra. 2015. Mixed-Initiative Design of Game Levels: Integrating Mission and Space into Level Generation. In *Foundations of Digital Games*.
- [8] Paul Klint, Tijs Van Der Storm, and Jurgen Vinju. 2009. Rascal: A Domain Specific Language for Source Code Analysis and Manipulation. In *Source Code Analysis and Manipulation, 2009. SCAM'09. Ninth IEEE International Working Conference on*. IEEE, 168–177.
- [9] Ralf Lämmel. 2018. *Software Languages: Syntax, Semantics, and Metaprogramming*. Springer.
- [10] Stephen Lavelle. 2015. PuzzleScript. <https://www.puzzlescript.net>. (2015). Visited May 20th 2018.
- [11] Tom Mens, Michel Wermelinger, Stéphane Ducasse, Serge Demeyer, Robert Hirschfeld, and Mehdi Jazayeri. 2005. Challenges in Software Evolution. In *Principles of Software Evolution, Eighth International Workshop on*. IEEE, 13–22.
- [12] Noor Shaker, Julian Togelius, and Mark J. Nelson. 2016. *Procedural Content Generation in Games: A Textbook and an Overview of Current Research*. Springer.
- [13] Adam M Smith and Michael Mateas. 2011. Answer set programming for procedural content generation: A design space approach. *IEEE Transactions on Computational Intelligence and AI in Games* 3, 3 (2011), 187–200.
- [14] Gillian Smith and Jim Whitehead. 2010. Analyzing the Expressive Range of a Level Generator. In *Proceedings of the 2010 Workshop on Procedural Content Generation in Games*. ACM, 4.
- [15] Adam Summerville, Julian R. H. Mariño, Sam Snodgrass, Santiago Ontañón, and Levi H. S. Lelis. 2017. Understanding Mario: An Evaluation of Design Metrics for Platformers. In *Proceedings of the 12th International Conference on the Foundations of Digital Games (FDG '17)*. Article 8, 10 pages.
- [16] Roland Van der Linden, Ricardo Lopes, and Rafael Bidarra. 2013. Designing Procedurally Generated Levels. In *Proceedings of the the second workshop on Artificial Intelligence in the Game Design Process*.
- [17] Roland van der Linden, Ricardo Lopes, and Rafael Bidarra. 2014. Procedural Generation of Dungeons. *IEEE Transactions on Computational Intelligence and AI in Games* 6, 1 (2014), 78–89.
- [18] Arie Van Deursen, Paul Klint, and Joost Visser. 2000. Domain-Specific Languages: An Annotated Bibliography. *ACM Sigplan Notices* 35, 6 (2000), 26–36.
- [19] Georgios N. Yannakakis and Julian Togelius. 2018. *Artificial Intelligence and Games*. Springer. <http://gameaibook.org>.