

Grounding Natural Language Prompts in Expressive Super Mario Level Generation

Anonymous ACL submission

Abstract

Natural-language-controllable procedural content generation depends critically on how linguistic concepts are grounded in structured representations. We show that widely used benchmarks rely on coarse semantic encodings that collapse distinct concepts, obscure grounding failures, and systematically inflate apparent instruction-following performance. Focusing on Super Mario level generation, we introduce MARIOPCG, a higher-fidelity dataset with expanded semantic coverage, and evaluate multiple decoder-only language models under controlled conditions. Increasing representational granularity exposes severe controllability failures in limited-capacity models that remain invisible under coarser benchmarks, while larger models exhibit stable behavior only when the representation supports meaningful grounding. These findings establish dataset semantic granularity as a necessary condition for valid evaluation of grounded language control and suggest that prior conclusions drawn from semantically collapsed benchmarks reflect representational artifacts rather than model capability. We will publicly release the dataset, prompts, and evaluation code to support reproducibility and further research.

1 Introduction

Procedural content generation (PCG) supports personalization, replayability, and reduced authoring cost by adapting game content to targeted player experiences. A common framing treats generation as the problem of producing content that satisfies experience-level objectives, such as difficulty, engagement, or frustration, often inferred from player behavior or preferences (Yannakakis and Togelius, 2011; Shaker et al., 2012; Yu and Trawick, 2011). Under this perspective, generators are evaluated by their ability to satisfy intended experiential constraints rather than by realism or novelty alone.

To enable controllability, prior PCG systems expose control through tunable parameters, optimiza-

tion objectives, mixed-initiative tools, or latent control variables learned by generative models (Togelius et al., 2011; Liapis et al., 2013; Yannakakis et al., 2014). While effective, these interfaces are often indirect, require domain expertise, or rely on representations that are not easily interpretable by human designers.

Natural language offers a more expressive and accessible control interface. Recent advances in instruction-following language models enable free-form textual descriptions to be mapped to structured outputs, making language a practical mechanism for controllable generation. In games, this approach has been applied to level generation and world synthesis, including text-conditioned generation of functional puzzle and platformer levels (Todd et al., 2023; Sudhakaran et al., 2023; Schrum et al., 2025).

Using language as a control signal introduces a grounding problem. Natural-language prompts specify abstract concepts and relations, while generators must produce concrete structured representations under hard functional constraints. Grounding quality therefore depends not only on model capacity, but also on the semantic fidelity of the underlying representation and dataset. This issue is especially pronounced in tile-based platformer domains such as Super Mario Bros, where commonly used datasets collapse semantically distinct entities and omit gameplay-relevant tiles. As a result, many concepts that are natural to express in language cannot be explicitly represented or evaluated, conflating representational limitations with model limitations.

In this work, we empirically demonstrate that dataset semantics and representation design fundamentally constrain observable instruction-following behavior. Rather than proposing a new generation architecture, we show how representational granularity shapes the space of expressible instructions and the failure modes that emerge dur-

ing expressive Super Mario level generation.

Our main contributions and findings are summarized as follows:

- We establish **dataset semantic granularity** as a structural constraint on natural-language instruction following in procedural content generation, demonstrating that representational choices fundamentally bound prompt expressivity and the grounding failure modes that can be observed during evaluation.
- We introduce **MarioPCG**, a higher-fidelity Super Mario Bros level dataset with expanded semantic coverage, enabling controlled study of grounded language conditioning while remaining compatible with standard PCGML pipelines and evaluation tooling.
- Through controlled experiments across multiple decoder-only language models under matched training and inference conditions, we show that increasing representational expressivity exposes controllability failures in limited-capacity models that remain hidden under coarser benchmarks, while larger models exhibit more stable behavior with expressive supervision.

2 Task Formulation and Problem Definition

Procedural content generation via machine learning (PCGML) learns generators from existing game artifacts and supports intent-driven design workflows such as mixed-initiative and co-creative generation (Summerville et al., 2018). A persistent constraint in PCGML is limited domain data, as game-specific corpora are typically orders of magnitude smaller than standard language or vision datasets (Awiszus et al., 2020). This limitation is compounded by the heterogeneity of game representations, which are often specific to individual games or franchises and constrain dataset scale and semantic expressivity (Summerville et al., 2018). Throughout this work, we use expressivity to denote the number, granularity, and compositionality of semantic concepts that can be referenced in prompts and evaluated through grounded generation.

We focus on tile-based 2D platformer levels represented as grids of discrete symbols corresponding to tiles and objects. Such grids are commonly linearized into sequences to enable sequence modeling while preserving local spatial dependencies

(Summerville and Mateas, 2016). This representation underpins a wide range of PCGML approaches and is used in public datasets such as the Video Game Level Corpus (VGLC) (Summerville et al., 2016).

2.1 Task Definition

We formulate natural-language-controllable level generation as conditional sequence modeling with grounded linguistic inputs. Each training example consists of a natural-language prompt x , describing desired properties of a level, and a corresponding level y , represented as a serialized sequence obtained by linearizing a tile grid into discrete tokens for tiles and objects (Summerville and Mateas, 2016). Given a dataset $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^N$, we fine-tune a decoder-only language model using standard causal language modeling loss, where the prompt x_i provides conditioning context and the loss is applied to the level tokens y_i .

This formulation treats controllable PCG as an instruction-following problem under hard structural and functional constraints. The model must generate sequences that are not only distributionally plausible, but also structurally coherent and aligned with the semantics expressed in the natural-language prompt.

2.2 Natural-Language Controllable PCG as a Grounding Problem

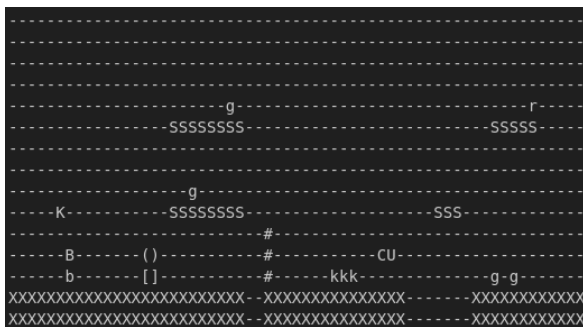
A central challenge in controllable PCG is specifying user intent in a form that is both expressive and operational. Prior work commonly relies on non-linguistic control channels, such as latent vectors optimized to satisfy target properties or structured constraints exposed through interactive editing. For example, Volz et al. evolve latent vectors in the space of a generative adversarial network to steer level generation toward desired characteristics, using continuous control signals that are not directly interpretable by human designers and typically require iterative search (Volz et al., 2018). Mixed-initiative tools similarly expose control through explicit constraints and editor interactions rather than free-form language.

In contrast, text-conditioned generation treats control as a grounding problem, where natural-language descriptions must be mapped to structured, non-linguistic outputs under domain-specific constraints. Prior work such as MarioGPT demonstrates that linguistic prompts can be used to condition level generation by mapping textual de-

184 descriptions directly to serialized level representa- 214
 185 tions (Sudhakaran et al., 2023). Related work 215
 186 further shows that large language models trained 216
 187 for NLP tasks can be adapted to generate struc- 217
 188 tured game content when appropriately conditioned 218
 189 (Todd et al., 2023). This framing aligns control- 219
 190 lable PCG with established problems in instruc- 220
 191 tion following and semantic grounding, providing a 221
 192 principled basis for analyzing how representational 222
 193 choices shape observable instruction-following be- 223
 194 havior. 224

195 **Semantic Granularity.** Within this framing, we 225
 196 define *semantic granularity* as the resolution at 226
 197 which distinct functional entities and structural re- 227
 198 lations are explicitly represented and supervised 228
 199 in the dataset. Coarser representations collapse se- 229
 200 mantically distinct elements into shared symbols, 230
 201 limiting the set of concepts that can be expressed 231
 202 in natural-language prompts and obscuring ground- 232
 203 ing failures during evaluation. Finer-grained rep- 233
 204 resentations increase expressivity by preserving 234
 205 meaningful distinctions, but simultaneously im- 235
 206 pose stricter grounding requirements on instruc- 236
 207 tion-following models. 237

208 3 MarioPCG Dataset 238



(a) String Representation



(b) Actual Level

Figure 1: String representation and corresponding level layout in the MarioPCG dataset.

209 We introduce MARIOPCG, an open-source 260
 210 dataset for Super Mario Bros level generation that 261
 211 increases semantic granularity in tile-based repre- 262
 212 sentations while remaining compatible with stan- 263
 213 dard PCGML pipelines. MarioPCG builds on the

widely used symbolic grid encoding for Super 214
 Mario Bros, refining it to preserve fine-grained 215
 distinctions among enemies, objects, and structural 216
 elements. 217

MarioPCG addresses two limitations of existing 218
 datasets: limited dataset scale and reduced seman- 219
 tic coverage. The dataset contains a larger number 220
 of levels than commonly used benchmarks and em- 221
 ploys an expanded tile vocabulary that more accu- 222
 rately reflects gameplay-relevant entities present in 223
 the original games. 224

The dataset is constructed from levels recreated 225
 in Super Mario Maker 2, which provides a uni- 226
 fied authoring environment capable of reproducing 227
 levels from multiple Super Mario Bros titles us- 228
 ing a consistent internal representation. MarioPCG 229
 consists primarily of faithful recreations of origi- 230
 nal levels, complemented by a small subset of 231
 original user-created levels. Level selection prior- 232
 itizes adherence to established Super Mario level 233
 design patterns, as characterized by Dahlskog et 234
 al. (Dahlskog and Togelius, 2012) and Thompson 235
 (Thompson, 2015). The released dataset contains 236
 exclusively symbolic, ASCII-based tile grids en- 237
 coding abstract level structure and does not include 238
 original visual, audio, or executable game assets. 239
 Source levels were obtained from a publicly re- 240
 leased Super Mario Maker 2 API dataset by The 241
 Great Rambler (TheGreatRambler, 2022). 242

Because the source data originates from Super 243
 Mario Maker 2, preprocessing is required to con- 244
 vert levels into the standard tile-based representa- 245
 tion used in Super Mario Bros PCG research. This 246
 includes minor tile placement adjustments to bet- 247
 ter reflect original game mechanics. The resulting 248
 dataset supports a wide range of enemies, items, 249
 and blocks compatible with the Mario AI Frame- 250
 work. In total, MarioPCG includes 26 distinct tile 251
 types, compared to 13 in the VGLC, enabling a 252
 broader range of concepts to be expressed and con- 253
 trolled through natural-language prompts. 254

The final dataset contains 112 levels, including 255
 77 faithful recreations from Super Mario Bros, Su- 256
 per Mario Bros: The Lost Levels, Super Mario 257
 Land, and New Super Mario Bros, as well as 35 258
 original user-created levels. While additional recre- 259
 ations exist in Super Mario Maker 2, many could 260
 not be included due to mechanics or enemies not 261
 supported by the Mario AI Framework (Khalifa, 262
 2019). By providing higher-fidelity level repre- 263
 sentations with expanded semantic coverage, Mari- 264
 oPCG enables controlled study of grounded natural- 265

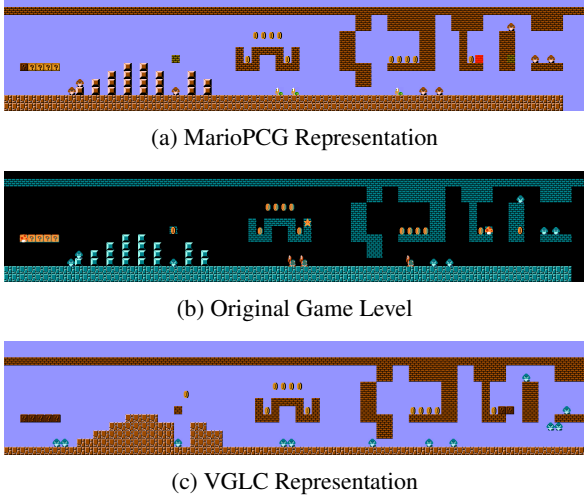


Figure 2: Comparison of level 1–2 in Super Mario Bros across MarioPCG, the original game, and the VGLC representation.

language control in procedural content generation while remaining independent of any specific modeling approach.

4 Experiments

4.1 Setup

Levels are linearized into token sequences using column-wise traversal. Training examples are constructed by extracting fixed-length slices of 800 characters using a sliding window that advances by one column at a time without crossing level boundaries, yielding 17,100 training examples for MarioPCG. Each slice is paired with an automatically generated natural-language prompt following the MarioGPT procedure, where prompts describe the presence and quantity of tiles, enemies, and structures. Prompts and level slices are formatted as user–assistant pairs, with prompts provided as input and serialized levels as output.

To isolate the effect of dataset composition, we construct a parallel training set from the original MarioGPT dataset derived from VGLC resulting in 7,250 training examples. The DistilGPT2 VGLC baseline is reproduced using the official MarioGPT implementation and training recipe.

For evaluation, we generate 500 samples per model and dataset using a fixed sampling temperature of 2.0, and evaluate all models under identical sampling conditions.

4.2 Evaluation Metrics

We evaluate generated levels along four complementary axes that capture different aspects of

grounded instruction following: functional validity, prompt adherence, structural integrity, and memorization.

Playability. We measure playability using the Mario AI Framework by simulating each generated level with Robin Baumgarten’s A* agent. A level is considered playable if the agent completes it in a single attempt under default settings. We interpret solvability as a conservative functional-validity signal rather than a definitive measure of human playability, due to rare agent failure cases.

Prompt Adherence. To measure how well generated levels satisfy the natural-language prompt, we adopt the caption adherence score (c-score) introduced by Schrum et al. (Schrum et al., 2025), with a modification to account for under-specified prompts. In practice, prompts typically reference only a subset of available concepts. We therefore compute adherence only over concepts explicitly mentioned in the prompt, rather than treating unmentioned concepts as absent.

For example, given the prompt “many goombas and coins”, a valid and playable level necessarily contains ground and supporting structures that are not mentioned. Under the original c-score definition, such levels would be penalized for containing unspecified concepts, while degenerate levels consisting primarily of empty space populated only by goombas and coins could receive higher adherence scores despite being unplayable. This failure mode becomes more pronounced as semantic vocabularies grow more granular. Our modification avoids this issue by scoring adherence only over explicitly referenced concepts.

Given a prompt p and an automatically generated caption c for the produced level, we compute:

$$\text{c-score}(p, c) = \frac{\sum_{t \in T(p)} \text{match}(\text{phrase}(t, p), \text{phrase}(t, c))}{|T(p)|} \quad (1)$$

$$\text{match}(p_t, c_t) = \begin{cases} 1.0 & \text{if } p_t = c_t \\ 1.0 - \frac{|\text{qu}(p_t) - \text{qu}(c_t)|}{|Q| - 1} & \text{if } \text{co}(p_t) \wedge \text{co}(c_t) \\ 0.1 & \text{if } p_t \neq \emptyset \wedge c_t \neq \emptyset \\ -1.0 & \text{otherwise} \end{cases} \quad (2)$$

where $T(p) \subseteq T$ denotes the set of concepts referenced in the prompt. The function $\text{phrase}(t, s)$ extracts the quantity phrase for concept t from text s , and $\text{match}(\cdot, \cdot)$ assigns exact or graded partial credit based on ordinal distance over an ordered quantity scale (e.g., $\{no, few, some, many\}$). This

341 formulation isolates adherence to stated constraints
342 while allowing models to complete unspecified
343 structure necessary for coherent, playable levels.

344 **Structural Integrity.** We additionally evaluate
345 domain-specific integrity constraints for multi-tile
346 structures. Pipes and cannons require consistent
347 top and body arrangements, and generated outputs
348 sometimes violate these constraints. We report the
349 percentage of generated levels containing at least
350 one malformed pipe and at least one malformed
351 cannon.

352 Memorization.

353 To measure direct copying from the training set,
354 we adopt an n-gram-based plagiarism metric over
355 level columns. For a given n , we extract all ordered
356 column n-grams from the training levels and
357 store them as a reference set T_n . For each gener-
358 ated level, we extract the corresponding ordered
359 n-grams L_n and compare each $l \in L_n$ against T_n .
360 An n-gram is counted as matched if it appears ex-
361 actly in T_n , or if it matches a training n-gram with
362 normalized Hamming similarity of at least 95%, al-
363 lowing for small tile-level variations that preserve
364 overall structure.

365 We compute memorization as the fraction of
366 ordered n-grams in L_n that match the training set.
367 A generated level is considered copied if more than
368 50% of its ordered n-grams are matched. Following
369 Dahlskog and Togelius’ definition of level “beats”
370 as short recurring challenge segments separated
371 by lower-intensity regions (Dahlskog and Togelius,
372 2012), we use a small n-gram size ($n=2$) together
373 with a majority-overlap threshold, so that isolated
374 beat reuse is permitted while extended structural
375 copying results in more than 50% n-gram overlap.

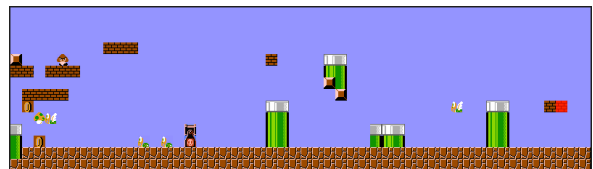
376 5 Results

377 We evaluate controllable level generation across
378 multiple decoder-only language models, including
379 Qwen 2.5, Qwen 3, Gemma 3, and Llama 3.1, all
380 fine-tuned using standard causal language model-
381 ing loss applied to assistant tokens only. Models
382 are trained on user–assistant formatted data with-
383 out architectural modifications, external prompt en-
384 coders, or cross-attention mechanisms. This stands
385 in contrast to MarioGPT’s DistilGPT2-based archi-
386 tecture, which injects prompt information through
387 cross-attention layers and relies on a frozen encoder
388 for prompt representations.

389 5.1 Effect of Prompt and Dataset Granularity

390 Expanding the prompt space reveals severe control-
391 lability failures in the MarioGPT base model that
392 are not observable under coarser prompt schemas.
393 As shown in Figure 3, generation collapses toward
394 overuse of specific structural elements, most nota-
395 bly pipes and, to a lesser extent, cannons, at the
396 expense of structural diversity and coherence. We
397 refer to this failure mode as *Structural Overcom-
398 pensation*, where a small subset of structures domi-
399 nates generation regardless of the prompt’s relative
400 balance. This degradation becomes increasingly
401 pronounced as the prompt space becomes more
402 granular and natural. In contrast, larger and more
403 recent decoder-only models accommodate the ex-
404 panded semantic vocabulary without similar col-
405 lapse, motivating their use for controllable PCG
406 under richer natural-language prompts.

407 Figure 3 further shows that increasing the num-
408 ber of promptable classes causes severe degrada-
409 tion in the MarioGPT DistilGPT2 base model. No-
410 tably, a substantial adherence drop is observed even
411 when holding the prompt schema fixed to the origi-
412 nal four classes and changing only the dataset,
413 indicating that controllability failures cannot be
414 attributed solely to the prompt interface. This be-
415 havior reflects *Prompt Signal Dilution*, where ex-
416 plicit linguistic constraints are overwhelmed by
417 dataset-level regularities under increased semantic
418 expressivity.



(a) MarioGPT (DistilGPT2) fine-tuned on MarioPCG with four prompt classes.



(b) MarioGPT (DistilGPT2) fine-tuned on MarioPCG with eleven prompt classes.

Figure 3: Generation degradation of the MarioGPT base model as the number of promptable classes increases. Prompt: many enemies, some pipes, some blocks, high elevation.

Table 1 summarizes results across datasets and model capacities. To isolate the effect of dataset

Model	Adherence	Playable	Broken Pipes	Broken Cannons	Memorization
VGLC (Summerville et al., 2016)					
DistilGPT2 (MarioGPT)	0.559	0.872	0.255	0.006	0.824
Qwen 2.5 (7B)	0.488	0.684	0.180	0.003	0.194
Qwen 2.5 (14B)	0.537	0.718	0.188	0.004	0.288
Qwen 3 (8B)	0.548	0.720	0.179	0.002	0.192
Qwen 3 (14B)	0.496	0.738	0.171	0.002	0.314
Gemma 3 (12B)	0.361	0.612	0.073	0.012	0.183
Llama 3.1 (8B)	0.511	0.662	0.146	0.002	0.178
MarioPCG (Ours)					
Qwen 2.5 (7B)	0.690	0.780	0.018	0.004	0.196
Qwen 2.5 (14B)	0.704	0.784	0.006	0.004	0.272
Qwen 3 (8B)	0.698	0.750	0.014	0.000	0.190
Qwen 3 (14B)	0.691	0.733	0.004	0.000	0.226
Gemma 3 (12B)	0.449	0.612	0.088	0.060	0.234
Llama 3.1 (8B)	0.677	0.733	0.016	0.021	0.264

Table 1: Main results across datasets and model capacities. All results use the same 4-class prompt interface at inference time. Models trained on MarioPCG receive 11-class semantic supervision during training.

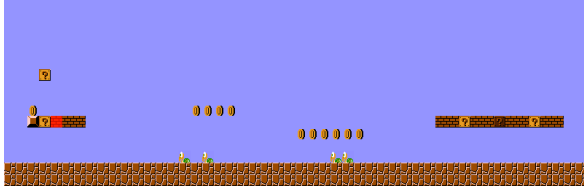
semantic granularity from the prompt interface, all results use the same four-class prompt schema at inference time. For the VGLC dataset, models are trained and evaluated using this schema. For MarioPCG, models are trained with the full eleven-class semantic schema but evaluated using prompts that mention only four randomly selected classes.

On the VGLC dataset, the MarioGPT baseline exhibits both high playability and high memorization. Because many VGLC levels are fully playable by construction, reproducing large portions of the training data can yield high solvability under automated agents. As a result, playability alone can overestimate generalization performance in settings where memorization is prevalent, since functional validity may arise from copying rather than grounded instruction following. Modern decoder-only models generally exhibit lower copying on VGLC but do not outperform MarioGPT in prompt adherence under the same coarse prompt schema. We additionally observed that Gemma 3 trained on VGLC often collapses toward generating overly simple and repetitive level structures that are largely invariant to the prompt, producing outputs that are technically playable and satisfy coarse prompt constraints but are qualitatively poor and weakly responsive to prompt variation.

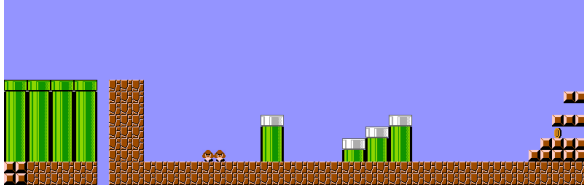
When trained on MarioPCG but evaluated using the four-class prompt interface, the MarioGPT

base model exhibits essentially unchanged prompt adherence and memorization. We hypothesize that this behavior arises from implicit attribute interference. Prior work shows that unannotated attributes in datasets can provide strong learning signals that degrade control of explicitly specified attributes, a phenomenon referred to as attribute transfer (Ma et al., 2023), which manifests here where unprompted structural and compositional attributes dominate generation in limited-capacity models. In this setting, the four-class prompt constrains only coarse tile categories, while many structural and compositional properties of level slices remain unspecified. Because MarioPCG covers a broader range of realizations under the same prompt labels, limited-capacity models may prioritize modeling dataset-level regularities over preserving the prompt signal, weakening controllability. The absence of this degradation in larger decoder-only models supports this interpretation.

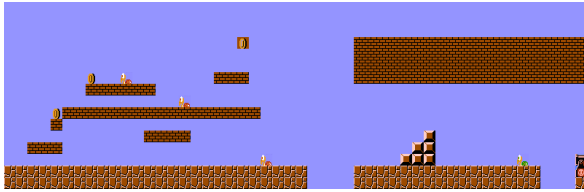
Figure 4 presents qualitative examples of levels generated under distinct natural-language prompts, illustrating how entity types and quantity specifications are grounded in the generated structure. Taken together, these effects highlight recurring patterns of grounding failure as semantic granularity increases, which can help interpret aggregate controllability metrics.



(a) Make a level with many coins, some powerups and some koopas



(b) I want a level that has many pipes and some goombas



(c) Generate a level with many breakable blocks, no pipes, some koopas, no goombas, a few special enemies and a few coin blocks

Figure 4: Examples of levels generated under distinct natural-language prompts.

Model	Adherence	Playable	Broken Pipes	Broken Cannons	Memorization
Qwen 2.5 (7B)	0.725	0.707	0.017	0.002	0.116
Qwen 2.5 (14B)	0.753	0.712	0.016	0.002	0.087
Qwen 3 (8B)	0.723	0.596	0.014	0.000	0.108
Qwen 3 (14B)	0.724	0.716	0.022	0.000	0.090
Gemma 3 (12B)	0.659	0.587	0.044	0.002	0.091
Llama 3.1 (8B)	0.672	0.701	0.018	0.002	0.170

Table 2: Results on MarioPCG using exact numerical quantity specification in prompts.

5.2 Effect of Quantity Specification in Prompts

We compare semantic quantity descriptors, such as *some* or *many*, with exact numerical quantity specifications in natural-language prompts. Aside from the form of quantity expression, all other factors, including dataset, models, and training procedure, are held constant.

Numerical prompts specify exact target counts during generation but are evaluated using the same ordinal quantity bins as semantic prompts. After generation, realized quantities are mapped to bins and compared against the bins implied by the prompt, ensuring both formulations are scored under the same adherence metric.

Exact numerical quantity specification yields modest improvements in prompt adherence for

the strongest models, particularly Qwen 2.5 and Qwen 3, while Gemma 3 and Llama 3.1 exhibit slight degradation. Playability generally decreases under numerical prompting, indicating a trade-off between precise quantity control and structural validity. These results suggest that enforcing exact numerical constraints can interfere with maintaining global level coherence, especially in smaller or less robust models.

This behavior aligns with prior work on numeracy in language models. Earlier studies show that standard language models treat numerals as discrete symbols without explicit numerical structure, leading to poor handling of exact values (Spithourakis and Riedel, 2018). More recent work demonstrates that large language models encode numerical magnitude in hidden representations in a linearly decodable manner, indicating that numeric value information is implicitly present despite token-level limitations (Zhu et al., 2025). These results further illustrate that increased constraint precision does not guarantee better grounding under insufficient representational support.

6 Segment-Chained Generation

Decoder-only models trained on fixed-width level slices do not naturally generate arbitrary-length levels in a single forward pass. We therefore introduce *segment-chained generation*, an inference-time continuation procedure that composes multiple fixed-width generations into a longer level while leaving the underlying model and training objective unchanged.

Fixed-size scene generation is common in PCGML, and longer levels are often constructed by composing multiple scenes or segments (Volz et al., 2018; Green et al., 2020; Summerville and Mateas, 2016). Similar distinctions between scene-level generation and full-level construction appear in recent text-conditioned diffusion approaches (Schrum et al., 2025). Segment-chained generation extends slice-trained decoder-only generators to arbitrary length while preserving local grounding behavior. We do not evaluate long-horizon grounding in this work.

6.1 Overlap-Based Continuation

Models are trained to generate fixed-width segments of 50 columns and 16 tiles in height. At inference time, we repeatedly decode segments to reach a target level length. To encourage continu-

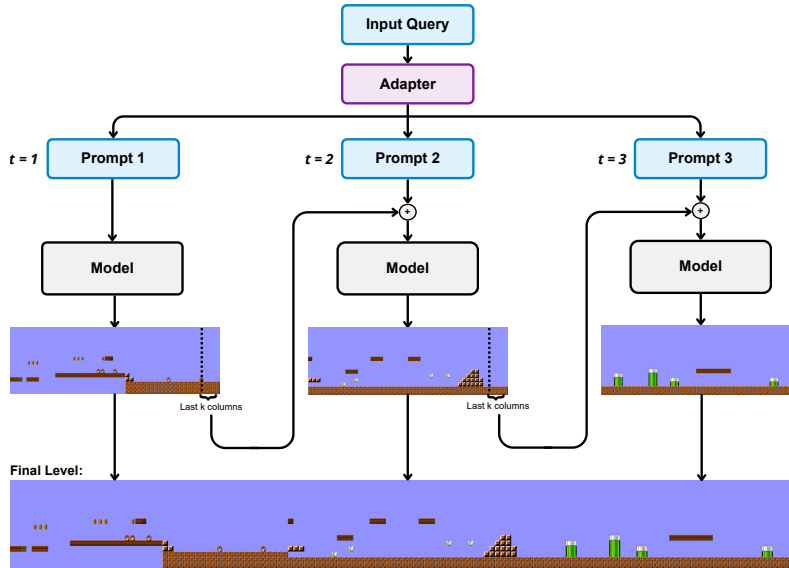


Figure 5: Temporal prompt decomposition for segment-chained generation. An input query is decomposed into ordered prompts. Each segment is generated conditioned on the current prompt and overlapping context, and only the non-overlapping suffix is appended to the growing level.

ity, we carry forward an overlap of o columns by inserting the final o columns of segment t as initial assistant-side context for segment $t+1$. A new segment of width W is then decoded, and only its non-overlapping suffix is appended to the growing level. This rolling-window procedure stays within the model’s trained segment format and requires no architectural modifications.

6.2 Temporal Prompt Decomposition

Segment-chained generation also supports temporally varying intent. Given a prompt describing ordered phases, we apply a prompt adapter that decomposes the input into a sequence of prompts. Segments are generated sequentially, each conditioned on the current sub-prompt and the overlapping context from the previous segment. The prompt is updated when advancing to the next phase. Figure 5 illustrates this procedure.

Conceptually, this approach aligns with long-form generation methods in NLP that decompose generation into ordered steps or phases (Yang et al., 2022, 2023). Unlike revision-based reprompting pipelines, segment-chained generation uses a fixed-window continuation loop with scheduled prompt changes, chosen to match the fixed-segment training regime while enabling temporal control.

7 Conclusion

This work establishes dataset semantic granularity as a decisive factor in revealing grounding and

controllability in natural-language-conditioned procedural content generation. By varying representation design while holding model architectures and training fixed, we show that coarse benchmarks systematically overestimate instruction-following ability by collapsing distinct semantic concepts and obscuring failure modes. Through the introduction of MARIOPCG, a higher-fidelity Super Mario level dataset with expanded semantic coverage, we demonstrate that increased representational expressivity exposes grounding failures that remain hidden under existing evaluation protocols. Limited-capacity models degrade sharply as the prompt space becomes richer and more natural, while larger decoder-only models exhibit stable behavior only when the representation supports meaningful grounding. These results indicate that prior conclusions about controllable generation are often shaped by representational constraints in benchmark design rather than model capability alone. More broadly, this work reframes controllable procedural generation as a grounding and representation problem, providing a principled foundation for future benchmarks and evaluations of instruction-following behavior.

8 Limitations

This study is conducted under a deliberately constrained experimental setting. The evaluation is limited to tile-based Super Mario Bros level generation using symbolic grid representations, which

603 may not directly generalize to other procedural con-
604 tent generation domains such as three-dimensional
605 environments, continuous control spaces, or graph-
606 structured representations.

607 Natural-language prompts are automatically de-
608 rived from level content rather than authored by hu-
609 mans. While this enables controlled manipulation
610 of semantic granularity, it restricts linguistic diver-
611 sity and evaluates grounding under structured, syn-
612 thetic prompt distributions rather than open-ended
613 human instructions.

614 Grounding quality is measured using automated
615 proxy metrics, including prompt adherence, playa-
616 bility, structural validity, and memorization indica-
617 tors. These metrics capture complementary aspects
618 of grounded generation but do not fully reflect hu-
619 man judgments of controllability or design intent.

620 Finally, models generate fixed-width level seg-
621 ments, with longer levels produced via inference-
622 time composition. This design facilitates local
623 grounding analysis but does not address native long-
624 horizon generation or global structural dependency
625 modeling.

626 Although grounded in PCG, the observed failure
627 modes arise from representational collapse and may
628 therefore extend beyond this domain.

629 9 Ethical Considerations

630 This work uses level layout data derived from recre-
631 ations in Super Mario Maker 2. The released
632 MarioPCG dataset contains only symbolic, ASCII-
633 based representations of level layouts and does not
634 include original visual, audio, or executable game
635 assets. Levels are represented as abstract character
636 grids commonly used in prior Super Mario Bros.
637 procedural content generation research.

638 The symbolic grid representations are a long-
639 established abstraction in procedural content gen-
640 eration research for Super Mario Bros. and related
641 domains. They encode functional tile categories
642 and structural relationships rather than expressive
643 audiovisual elements, and are used as analytical
644 representations for modeling and evaluation rather
645 than as substitutes for the original game content.

646 The source data was obtained from a publicly
647 released Super Mario Maker 2 API dataset by The
648 Great Rambler under a CC BY-NC-SA 4.0 license
649 (TheGreatRambler, 2022). MarioPCG does not
650 include player identifiers, metadata, or other user-
651 related information and is distributed solely for
652 non-commercial academic research.

Acknowledgments

Anonymous Acknowledgments

References

- Maren Awiszus, Frederik Schubert, and Bodo Rosenhahn. 2020. *Toad-gan: Coherent style level generation from a single example*. In *Proceedings of the Sixteenth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*. 656-659.
- Yi-Chun Chen and Arnav Jhala. 2025. *Gametilenet: a semantic dataset for low-resolution game art in procedural content generation*. In *Proceedings of the Twenty-First AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment, AIIDE '25*. AAAI Press. 661-666.
- Steve Dahlskog and Julian Togelius. 2012. *Patterns and procedural content generation: revisiting mario in world 1 level 1*. In *Proceedings of the First Workshop on Design Patterns in Games, DPG '12*, page 1-8. ACM. 667-671.
- Michael Cerny Green, Luvneesh Mugrai, Ahmed Khalifa, and Julian Togelius. 2020. *Mario level generation from mechanics using scene stitching*. In *Proceedings of the IEEE Conference on Games (CoG)*. arXiv. 672-675.
- Sergey Karakovskiy and Julian Togelius. 2012. *The mario AI benchmark and competitions*. *IEEE Transactions on Computational Intelligence and AI in Games*, 4(1):55-67. 676-679.
- Ahmed Khalifa. 2019. *Mario-ai-framework*. <https://github.com/amidos2006/Mario-AI-Framework>. GitHub repository. 680-682.
- Ahmed Khalifa, Philip Bontrager, Sam Earle, and Julian Togelius. 2020. *Pcgrl: Procedural content generation via reinforcement learning*. In *Proceedings of the Sixteenth AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment (AIIDE)*. 683-686.
- Antonios Liapis, Georgios N. Yannakakis, and Julian Togelius. 2013. *Sentient sketchbook: Computer-aided game level authoring*. In *International Conference on Foundations of Digital Games*. 688-691.
- Congda Ma, Tianyu Zhao, Makoto Shing, Kei Sawada, and Manabu Okumura. 2023. *Focused prefix tuning for controllable text generation*. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 1116-1127, Toronto, Canada. Association for Computational Linguistics. 692-698.
- Jacob Schrum, Olivia Kilday, Emilio Salas, Bess Hagan, and Reid Williams. 2025. *Text-to-level diffusion models with various text encoders for super mario bros*. In *Proceedings of the 21st AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*. 699-703.

705	Noor Shaker, Julian Togelius, and Mark J. Nelson. 2011.	Vanessa Volz, Jacob Schrum, Jialin Liu, Simon M. Lu-	760
706	The 2010 mario AI championship: Level generation	cas, Adam Smith, and Sebastian Risi. 2018. Evolving	761
707	track. <i>IEEE Transactions on Computational Intelli-</i>	mario levels in the latent space of a deep convolu-	762
708	<i>gence and AI in Games</i> , 3(4):332–347.	tional generative adversarial network. In <i>Proceed-</i>	763
		<i>ings of the Genetic and Evolutionary Computation</i>	764
709	Noor Shaker, Georgios N. Yannakakis, and Julian To-	<i>Conference</i> , GECCO '18, page 221–228. ACM.	765
710	gelius. 2012. Towards player-driven procedural con-		
711	tent generation . In <i>Proceedings of the 9th Conference</i>	Kevin Yang, Dan Klein, Nanyun Peng, and Yuandong	766
712	<i>on Computing Frontiers</i> , CF '12, page 237–240, New	Tian. 2023. DOC: Improving long story coherence	767
713	York, NY, USA. Association for Computing Machin-	with detailed outline control . In <i>Proceedings of the</i>	768
714	ery.	<i>61st Annual Meeting of the Association for Compu-</i>	769
		<i>tational Linguistics (Volume 1: Long Papers)</i> , pages	770
715	Kristin Siu, Matthew Guzdial, and Mark Riedl. 2017.	3378–3465, Toronto, Canada. Association for Com-	771
716	Evaluating singleplayer and multiplayer in human	putational Linguistics.	772
717	computation games . pages 1–10.		
		Kevin Yang, Yuandong Tian, Nanyun Peng, and Dan	773
718	Georgios Spithourakis and Sebastian Riedel. 2018. Nu-	Klein. 2022. Re3: Generating longer stories with	774
719	meracy for language models: Evaluating and improv-	recursive reprompting and revision . In <i>Proceedings</i>	775
720	ing their ability to predict numbers . In <i>Proceedings</i>	<i>of the 2022 Conference on Empirical Methods in Nat-</i>	776
721	<i>of the 56th Annual Meeting of the Association for</i>	<i>ural Language Processing</i> , pages 4393–4479, Abu	777
722	<i>Computational Linguistics (Volume 1: Long Papers)</i> ,	Dhabi, United Arab Emirates. Association for Com-	778
723	pages 2104–2115, Melbourne, Australia. Association	putational Linguistics.	779
724	for Computational Linguistics.		
		Georgios Yannakakis and Julian Togelius. 2011.	780
725	Shyam Sudhakaran, Miguel Gonzalez-Duque, Claire	Experience-driven procedural content generation . <i>Af-</i>	781
726	Glanois, Matthias Freiberger, Elias Najarro, and Se-	<i>fective Computing, IEEE Transactions on</i> , 2:147–	782
727	bastian Risi. 2023. MarioGPT: Open-ended text2level	161.	783
728	generation through large language models . In <i>Ad-</i>		
729	<i>vances in Neural Information Processing Systems</i> .	Georgios N. Yannakakis, Antonios Liapis, and Con-	784
730	Curran Associates, Inc.	stantine Alexopoulos. 2014. Mixed-initiative co-	785
		creativity . In <i>International Conference on Founda-</i>	786
731	Adam Summerville and Michael Mateas. 2016. Super	<i>tions of Digital Games</i> .	787
732	mario as a string: Platformer level generation via		
733	lstms . In <i>Proceedings of the First International Joint</i>	Hong Yu and Tyler Trawick. 2011. Personalized proce-	788
734	<i>Conference of DiGRA and FDG</i> .	dural content generation to minimize frustration and	789
		boredom based on ranking algorithm . <i>Proceedings</i>	790
735	Adam Summerville, Sam Snodgrass, Matthew Guz-	<i>of the AAAI Conference on Artificial Intelligence and</i>	791
736	dial, Christoffer Holmgard, Amy K. Hoover, Aaron	<i>Interactive Digital Entertainment</i> , 7(1):208–213.	792
737	Isaksen, Andy Nealen, and Julian Togelius. 2018.		
738	Procedural content generation via machine learning	Fangwei Zhu, Damai Dai, and Zhifang Sui. 2025. Lang-	793
739	(pcgml) . <i>IEEE Transactions on Games</i> , 10(3):257–	uage models encode the value of numbers linearly .	794
740	270.	In <i>Proceedings of the 31st International Conference</i>	795
		<i>on Computational Linguistics</i> , pages 693–709, Abu	796
741	Adam James Summerville, Sam Snodgrass, Michael	Dhabi, UAE. Association for Computational Linguis-	797
742	Mateas, and Santiago Ontañón Villar. 2016. The	tics.	798
743	vglc: The video game level corpus . <i>Proceedings of</i>		
744	<i>the 7th Workshop on Procedural Content Generation</i> .	A Related Work	799
745	TheGreatRambler. 2022. Mario maker 2 datasets .		
		A.1 Procedural Content Generation	800
746	T. Thompson. 2015. The fine line between rehash and	Procedural Content Generation (PCG) studies al-	801
747	sequel: Design patterns of the super mario series.	gorithmic methods for producing game content	802
748	Design Patterns in Games Workshop.	such as levels, rules, and quests, often with the	803
		goal of improving replayability, personalization,	804
749	Graham Todd, Sam Earle, Muhammad Umair Nasir,	or authoring efficiency. Early PCG systems were	805
750	Michael Cerny Green, and Julian Togelius. 2023.	typically rule-based or search-based, while more	806
751	Level generation through large language models . In	recent work increasingly leverages machine learn-	807
752	<i>Proceedings of the 18th International Conference on</i>	ing to imitate designer distributions and capture	808
753	<i>the Foundations of Digital Games</i> , FDG 2023, page	implicit design regularities, an approach com-	809
754	1–8. ACM.	monly referred to as PCG via Machine Learning	810
		(PCGML) (Summerville et al., 2018). A recurring	811
755	Julian Togelius, Georgios Yannakakis, Kenneth Stanley,	challenge for PCGML is data scarcity and dataset	812
756	and Cameron Browne. 2011. Search-based proce-		
757	dural content generation: A taxonomy and survey .		
758	<i>Computational Intelligence and AI in Games, IEEE</i>		
759	<i>Transactions on</i> , 3:172 – 186.		

813 bias, since level representations and design con- 864
814 ventions are highly game-specific. Public corpora 865
815 such as the Video Game Level Corpus (VGLC) par- 866
816 tially address this by providing standardized, tile- 867
817 based level encodings across multiple games (Sum- 868
818 merville et al., 2016).

819 Super Mario Bros (SMB) has become a canon- 869
820 ical benchmark for PCG research due to its well- 870
821 studied design principles, strong functional con- 871
822 straints (playability), and widely used evaluation 872
823 tooling. The Mario AI Benchmark and its associ- 873
824 ated competitions helped standardize experimen- 874
825 tal protocols for SMB agents and level generation, 875
826 enabling reproducible evaluation and comparison 876
827 across generators (Karakovskiy and Togelius, 2012; 877
828 Shaker et al., 2011). Building on this ecosystem, 878
829 a large body of work has explored both learning- 879
830 based and search-based generators for SMB, in- 880
831 cluding sequential modeling with string representa- 881
832 tions (Summerville and Mateas, 2016), adversarial 882
833 approaches that evolve levels in a learned latent 883
834 space (Volz et al., 2018), one-shot or few-shot gen- 884
835 erative methods designed to reduce data require- 885
836 ments (Awiszus et al., 2020), and reinforcement 886
837 learning formulations that frame level design as a 887
838 sequential decision process (Khalifa et al., 2020).

839 A.2 Prompt Guided Level Generation

840 Recent work has explored natural language as a 890
841 high-level control interface for level generators, 891
842 motivated by the broader trend of using prompts 892
843 to steer generative models. In this setting, the 893
844 generator must map a free-form or templated tex- 894
845 tual request to a structured level representation 895
846 while preserving both global coherence and lo- 896
847 cal constraints. Todd et al. (Todd et al., 2023) 897
848 show that large language models can generate func- 898
849 tional levels (in Sokoban) and that generation qual- 899
850 ity scales strongly with dataset size, highlighting 900
851 both the promise of language-based generators and 901
852 the importance of training data construction. In 902
853 the SMB domain, MarioGPT (Sudhakaran et al., 903
854 2023) demonstrates open-ended text-to-level gen- 904
855 eration by fine-tuning a GPT-style model on tile se- 905
856 quences paired with automatically derived prompts, 906
857 enabling controllable generation through descrip- 907
858 tive text. More recently, text-conditioned diffu- 908
859 sion models have also been explored for tile-based 909
860 SMB generation, emphasizing captioning strate- 910
861 gies, text encoders, and designer-facing workflows 911
862 for assembling longer levels from generated seg- 912
863 ments (Schrum et al., 2025).

A.3 PCG Datasets

864 PCGML methods depend strongly on the availabil- 865
866 ity and structure of training data, yet game level 867
868 datasets are often game-specific and encode strong 868
869 representational assumptions. As noted by Sum- 869
870 merville et al. (Summerville et al., 2018), even 870
871 games within the same genre frequently differ in 871
872 underlying data structures and semantic interpreta- 872
873 tions, limiting dataset reuse across titles and con- 873
874 straining dataset scale. Public datasets therefore 874
875 play a central role in shaping what types of genera- 875
876 tive models and controls are feasible.

876 The Video Game Level Corpus (VGLC) (Sum- 876
877 merville et al., 2016) is the most widely used 877
878 dataset for level generation research. It pro- 878
879 vides standardized, parseable representations of 879
880 levels from twelve games, including multiple Su- 880
881 per Mario Bros. titles, using tile-based, graph- 881
882 based, and vector encodings. The tile represen- 882
883 tation in particular has enabled a large body of 883
884 PCGML work and effectively established SMB 884
885 as a canonical benchmark. However, VGLC re- 885
886 lies on simplified abstractions introduced during 886
887 automated extraction, such as collapsing multiple 887
888 enemy types into a single symbol and omitting 888
889 certain structural details. As a result, several re- 889
890 cent works rely on cleaned or modified versions of 890
891 VGLC, treating its levels as noisy or incomplete 891
892 data rather than ground truth. Text-conditioned dif- 892
893 fusion approaches for Super Mario Bros. (Schrum 893
894 et al., 2025) follow this strategy by correcting er- 894
895 rors in the corpus and pairing level segments with 895
896 structured textual descriptions, enabling generation 896
897 and repair conditioned on semantic intent and im- 897
898 proving alignment beyond direct imitation of tile 898
899 encodings.

900 Beyond static level corpora, player activity has 900
901 also been explored as a data source. The Gwario 901
902 human computation game (Siu et al., 2017) embeds 902
903 annotation tasks into modified Super Mario Bros. 903
904 levels, illustrating how player behavior can provide 904
905 structured signals during play. Although not de- 905
906 signed for level generation, such playtrace-based 906
907 data captures aspects of difficulty and engagement 907
908 absent from static layouts, but is costly to collect 908
909 and difficult to scale.

910 GameTileNet (Chen and Jhala, 2025) instead 910
911 focuses on semantic annotation of low-resolution 911
912 game art assets. By providing labeled tiles with af- 912
913 fordance and connectivity information, it supports 913
914 narrative-to-visual alignment and vision-language 914

915 modeling at the asset level. While it does not rep-
 916 resent full levels, it complements level corpora
 917 by enriching tile-level semantics within PCGML
 918 datasets.

919 B Dataset Details



























Tile Type	Symbol	Visualization
Empty	-	
Ground Block	X	
Hard Block	#	
Breakable Block	S	
Question Block w/ Powerup	?	
Coin	o	
Left pipe top	<	
Right pipe top	>	
Left pipe lower	[
Right pipe lower]	
Pipe Top Left w/ Plant Enemy	(
Pipe Top Right w/ Plant Enemy)	
Cannon Top	B	
Cannon Body	b	
Coin Brick Block	C	
Mushroom Brick Block	U	
Coin Question Block	!	
Invisible 1-Up Block	1	
Invisible Coin Block	2	
Goomba	g	
Winged Goomba	G	
Koopa	k	
Winged Koopa	K	
Red Koopa	r	
Winged Red Koopa	R	
Spiny	y	

Table 3: All Entities in Dataset

920 Table 3 demonstrates all unique tiles in the Mar-
 921 ioPCG dataset

922 C Quantity Bins

923 Prompts are constructed as combinations of entity
 924 classes and quantity specifications, computed over
 925 fixed-width level segments of 50 columns.

926 **Count-based entities** may be specified using either
 927 ordinal keywords or exact numerical quantities:

- 928 • {*no, few, some, many*} or exact counts special
 929 enemies
- 930 • {*no, few, some, many*} or exact counts pipes

- {*few, some, many*} or exact counts ground
 931 blocks 932
- {*few, some, many*} or exact counts hard blocks 933
- {*few, some, many*} or exact counts breakable
 934 blocks 935
- {*few, some, many*} or exact counts coin blocks 936
- {*no, few, some, many*} or exact counts coins 937
- {*no, few, some, many*} or exact counts
 938 powerups 939
- {*no, few, some, many*} or exact counts goom-
 940 bas 941
- {*no, few, some, many*} or exact counts koopas 942

When ordinal keywords are used, entity counts are mapped to dataset-derived bins specific to each class, as defined in Table 4. The keyword *no* corresponds to zero occurrences. When numerical quantities are used, the specified count is interpreted directly during generation and mapped back to ordinal bins for evaluation.

Ordinal and numerical quantity specifications correspond to separate training conditions; models are trained and evaluated using one prompt formulation at a time and do not jointly support both within a single model.

Elevation is specified independently:

- {*low, high*} elevation 956

The *low* and *high* elevation labels are determined from the height of the highest unbreakable blocks in a segment of the level.

Example prompts:

- many goombas, some pipes, few coin blocks,
 961 high elevation 962
- 5 goombas, 2 pipes, low elevation 963

Table 4 summarizes the entity classes, their associated tile symbols, and the thresholds used to map counts to ordinal bins.

967 D Training Details

968 D.1 Model Initialization

969 All models are initialized from instruction-tuned,
 970 decoder-only checkpoints and fine-tuned using
 971 parameter-efficient adaptation. We use the Unsloth
 972 framework for model loading and training. The

Entity	Tile Symbols	Low	Medium	High
Special Enemy	B, y	1	2	3
Pipe	<>, ()	1	2	5
Ground Block	X	1	87	150
Hard Block	#	1	11	23
Coin Block	Q, !, 2, C	1	2	6
Breakable Block	S	1	19	38
Coin	o	1	5	10
Powerup	1, ?, U, L	1	2	3
Goomba	g, G	1	2	5
Koopa	r, R, k, K	1	3	6

Table 4: Entity-to-tile mappings and ordinal quantity bin thresholds used for prompt generation and evaluation. Each entity aggregates multiple tile symbols from Table 3. Quantity bins correspond to dataset-derived low, medium, and high count thresholds.

maximum sequence length is fixed to 2048 tokens for all experiments. All models are loaded in 4-bit precision using bitsandbytes quantization, with dtype left to the framework default.

The following checkpoints are used:

- unsloth/Qwen3-14B-unsloth-bnb-4bit
- unsloth/Qwen3-8B-bnb-4bit
- unsloth/gemma-3-12b-it-bnb-4bit
- unsloth/Qwen2.5-7B-Instruct-bnb-4bit
- unsloth/Qwen2.5-14B-Instruct-bnb-4bit
- unsloth/Meta-Llama-3.1-8B-Instruct-bnb-4bit

D.2 Fine-Tuning Method

All models are fine-tuned using Low-Rank Adaptation (LoRA), with base model weights frozen throughout training. LoRA adapters are applied to attention and MLP projection layers following standard practice for decoder-only transformers.

The LoRA configuration is held constant across all model families:

- LoRA rank: 128
- LoRA alpha: 128
- LoRA dropout: 0.05
- Bias: none

D.3 Data Formatting

Training data are formatted as chat-style user–assistant interactions. Prompts are provided as user messages and serialized level representations are provided as assistant responses. Model-specific chat templates are applied for Gemma 3, Llama 3.1, Qwen 2.5, and Qwen 3 to ensure tokenizer compatibility. Aside from required formatting differences,

the training data and procedure are identical across models.

D.4 Optimization and Training

All models are trained for a single epoch using identical optimization hyperparameters. One epoch was selected based on internal validation, where additional epochs led to increased memorization without consistent improvements in prompt adherence. During training, 10% of the data is held out as a validation set. The validation split is used only for monitoring training behavior, selecting hyperparameters and selecting the number of epochs.

Training hyperparameters are as follows:

- Per-device batch size: 2
- Gradient accumulation steps: 4
- Effective batch size: 8
- Learning rate: 2×10^{-4}
- Warmup ratio: 0.05
- Optimizer: AdamW (8-bit)
- Weight decay: 0.01
- Learning rate scheduler: linear