# DREAMCRAFT: Text-Guided Generation of Functional 3D Environments in Minecraft

Sam Earle
New York University
Broklyn, USA

Filippos Kokkinos
Meta
London, UK

Yuhe Nie
New York University
Brooklyn, USA

Julian Togelius
New York University
Brooklyn, USA

Roberta Raileanu
Meta
London, UK

## Abstract

Procedural Content Generation (PCG) algorithms enable the automatic generation of complex and diverse artifacts. However, they don't provide high-level control over the generated content and typically require domain expertise. In contrast, text-to-3D methods allow users to specify desired characteristics in natural language, offering a high amount of flexibility and expressivity. But unlike PCG, such approaches cannot guarantee functionality, which is crucial for certain applications like game design. In this paper, we present a method for generating functional 3D artifacts from free-form text prompts in the open-world game Minecraft. Our method, DREAMCRAFT, trains quantized Neural Radiance Fields (NeRFs) to represent artifacts that, when viewed in-game, match given text descriptions. We find that DREAM-CRAFT produces more aligned in-game artifacts than a baseline that post-processes the output of an unconstrained NeRF. Thanks to the quantized representation of the environment, functional constraints can be integrated using specialized loss terms. We show how this can be leveraged to generate 3D structures that match a target distribution or obey certain adjacency rules over the block types. DREAMCRAFT inherits a high degree of expressivity and controllability from the NeRF, while still being able to incorporate functional constraints through domain-specific objectives.

*CCS Concepts:* • **Applied computing** → **Media arts**; • **Computing methodologies** → *Neural networks*.

*Keywords:* Procedural Content Generation, Neural Radiance Fields, Minecraft

## 1 Introduction

Procedural Content Generation (PCG) refers to a class of algorithms that can automatically create content such as video game levels [10, 31, 41, 48, 63, 70], 2D or 3D visual assets [6, 45, 48, 58, 78, 89], game rules or mechanics [20, 59, 80, 85, 93], or reinforcement learning (RL) environments [4, 12, 18, 33, 34, 36, 37, 42, 60, 63, 67, 82, 83]. PCG allows for compression of information [79, 85], increased replayability via endless variation [5, 76, 90], expression of particular aesthetics [1, 8, 26, 46], and reduction of human labour otherwise required to manually produce artifacts [13, 17, 71, 72]. These methods are procedural in the sense that they outline sets of procedures or rules for generating artifacts, such as adjacency constraints in Wave Function Collapse, local update rules in cellular automata, or heuristic search in constraint satisfaction. These procedures often leverage domain-specific knowledge in order to guarantee that generated artifacts are functional; that a game environment does not contain structures that violate physics, or that a player is able to navigate between key points within them. However, users cannot generally control such methods via free-form language, and control is limited to those metrics explicitly defined by designers.

In contrast, recent generative models have shown impressive abilities in generating diverse images, videos, or 3D scenes from text prompts describing the desired output in natural language [61, 64, 73]. These advances allow users to create high-quality content even if they are not domain experts. While these models can produce controllable and open-ended generations, the created content is not guaranteed to be functional. Functionality is particularly important for certain applications such as game design or the creation of RL environments. Some recent efforts leverage language
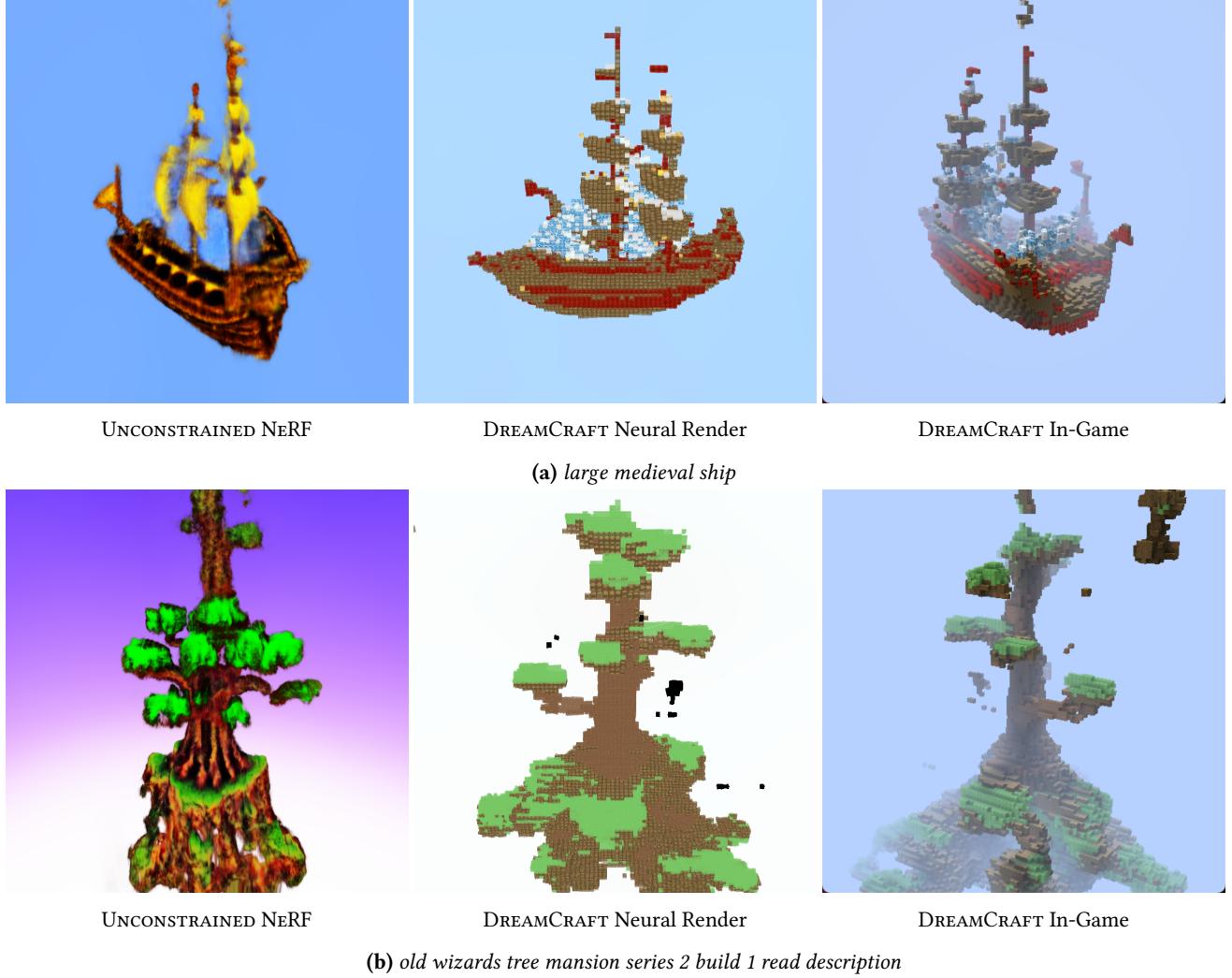
Sam Earle, Filippos Kokkinos, Yuhe Nie, Julian Togelius, and Roberta Raileanu



| Unconstrained NeRF | DreamCraft Neural Render | DreamCraft In-Game |

**(a)** *large medieval ship*



| Unconstrained NeRF | DreamCraft Neural Render | DreamCraft In-Game |

**(b)** *old wizards tree mansion series 2 build 1 read description*

**Figure 1.** Comparison of unconstrained NeRF (left), DreamCraft neural render, and DreamCraft in-game generation for a Planet Minecraft caption. Despite its lower resolution, DreamCraft's generated structures are similar in quality to those of the unconstrained NeRF, both closely matching the corresponding textual description.

models to generate level representations [77, 84], but effectively reduce text-based controls to a series of scalar values. Other methods train generative text-image models to produce levels that are made of discrete tiles [53] or controllable by actions [7], but they do not bring any functional guarantees: houses may spawn in disjointed, and birds may turn into bumblebees after a few frames.

In this work, we pursue a hybrid approach, adapting a generative model to operate on discrete 3D assets and incorporate functional constraints into its loss function. We propose a new method for generating functional 3D environments from free-form text prompts in the open-world game Minecraft. Our method, DreamCraft, trains a quantized NeRF to produce an environment layout that, when viewed in-game, matches a given text description (see Figures 1 and 2

for some examples). Experimenting with various quantization schemes, we find that using soft air blocks or annealing them from continuous (soft) to discrete is crucial for learning stability, while using discrete block types leads to the most recognizable structures. We evaluate the fidelity of our approach in matching generated artifacts to descriptions of both generic and domain-specific scenes and objects. We find that DreamCraft produces in-game artifacts that align with inputs more consistently than a baseline that post-processes the output of an unconstrained NeRF.

Thanks to its quantized representation of the game world, DreamCraft can jointly optimize loss terms that enforce local functional constraints on patterns of blocks. We show how this can be instantiated to, for example, generate 3D structures that match a target distribution or obey certain adjacency rules over the block types. By inheriting a high

degree of expressivity and controllability from the NeRF, while still being able to incorporate functional constraints through domain-specific objectives, DREAMCRAFT combines the strengths of both PCG and generative AI approaches, representing a first step towards democratizing flexible yet functional content creation. Our method has potential applications in the development of AI assistants for game design, as well as in the production of diverse and controllable environments for training and evaluating RL agents.

To summarize, our paper makes the following contributions:

1. introduces DREAMCRAFT, a new method for training a quantized NeRF to produce 3D structures that match a given textual description using a set of discrete Minecraft blocks,
2. studies different quantization schemes such as whether to use discrete or continuous block densities and types,
3. shows that the quantized NeRF produces more accurate Minecraft artifacts than an unconstrained NeRF, and
4. demonstrates how to incorporate functional constraints such as obeying certain target block distributions or adjacency rules.

## 2 Related Work

**Procedural Content Generation (PCG)** is becoming increasingly more popular for training and evaluating robust RL agents that can generalize across a wide range of settings [12, 33, 36, 37, 42, 67, 82, 83]. Generative models like ours provide a way of biasing the environment generations towards human-relevant ones, thus enabling to search more efficiently through vast environment spaces. Existing works indicate that environment generation can be controlled via computable metrics [14, 15, 22, 34, 41, 57, 68, 69, 72]. Novel environments can also be generated by learning on human datasets [28, 29, 48, 50, 74, 80], sometimes with additional functional constraints or post-processing [27, 40, 44, 86, 92]. More recently, Sudhakaran et al. [77], Todd et al. [84] use large language models to generate Sokoban and Mario 2D levels. However, our work is first to show how multi-modal models can be leveraged to guide the generation of 3D game environments. One of the most popular PCG algorithms is wave function collapse [23] which generates structurally consistent content from a single sample, such that its output matches tile-frequency and adjacency constraints. In this paper, we show how such constraints can be used in conjunction with text-guidance to generate environments where both high-level (e.g. via natural language) and low-level (e.g. via block target distributions or adjacency rules) aspects can be controlled.

**Text-to-3D Generation** Our work builds upon the many recent advances in text-to-3D generation [9, 55, 61]. For example, DVGO [81] is a supervised NeRF method that, instead of training an multi-layer perceptron (MLP), directly optimizes a voxel grid over the 3D space. Similarly, Pure-CLIPNeRF [43] uses a CLIP loss to guide a NeRF using both direct and implicit voxel grids. Our approach, DREAMCRAFT, resembles an implicit voxel grid approach, in that it uses MLPs to parameterize activations over discrete grids. But instead of outputting continuous RGB and density values and interpolating between nearest grid vertices (to determine activation at a given point during ray tracing), our approach uses MLPs to produce predictions over block types, and considers only the single nearest grid vertex during ray sampling (to determine within which block a given point resides).

**Text-to-Environment Generation** [53] train a text-conditioned decoder model on a dataset of hand-made levels in a 2D tile domain. Compared to NeRFs, which train a model to represent a single artifact (here, a level), the "5 dollar" decoder model can represent a distribution of artifacts. It can also potentially achieve a certain amount of generalization thanks to the pre-trained LLM which is used to encode text prompts.

**Minecraft Environment Generation** Several prior works have sought to generate Minecraft environments using both supervised and self-supervised methods. Awiszus et al. [2] use a 3D GAN [19] architecture to generate arbitrarily sized world snippets from a single example. Meanwhile, Hao et al. [30] envision Minecraft as a potential "sketchbook" for designing more photorealistic 3D landscapes, using an unsupervised neural rendering approach to generate the latter from preexisting Minecraft landscapes. To assist Minecraft players, Merino et al. [52] introduce a tool for interactive evolution using both a 3D generative model to generate the structure design and an encoding model for applying Minecraft-specific textures to the structure's voxels. Sudhakaran et al. [78] have shown that neural cellular automata can be used to grow complex 3D Minecraft artifacts made out of thousands of blocks such as castles, apartment blocks, and trees. Other works have leveraged search-based methods [91], evolutionary algorithms [51, 75], or even reinforcement learning [34] approaches to generate Minecraft structures. The game has also been a testbed for PCG algorithms [65, 66], open-endedness [21], artificial life [78], RL agents [25, 35, 38, 54], or foundation models for decision-making [3, 16, 39, 87]. However, our work is first to generate functional Minecraft environments directly from text prompts, enhancing the high-level controllability of the creation process.
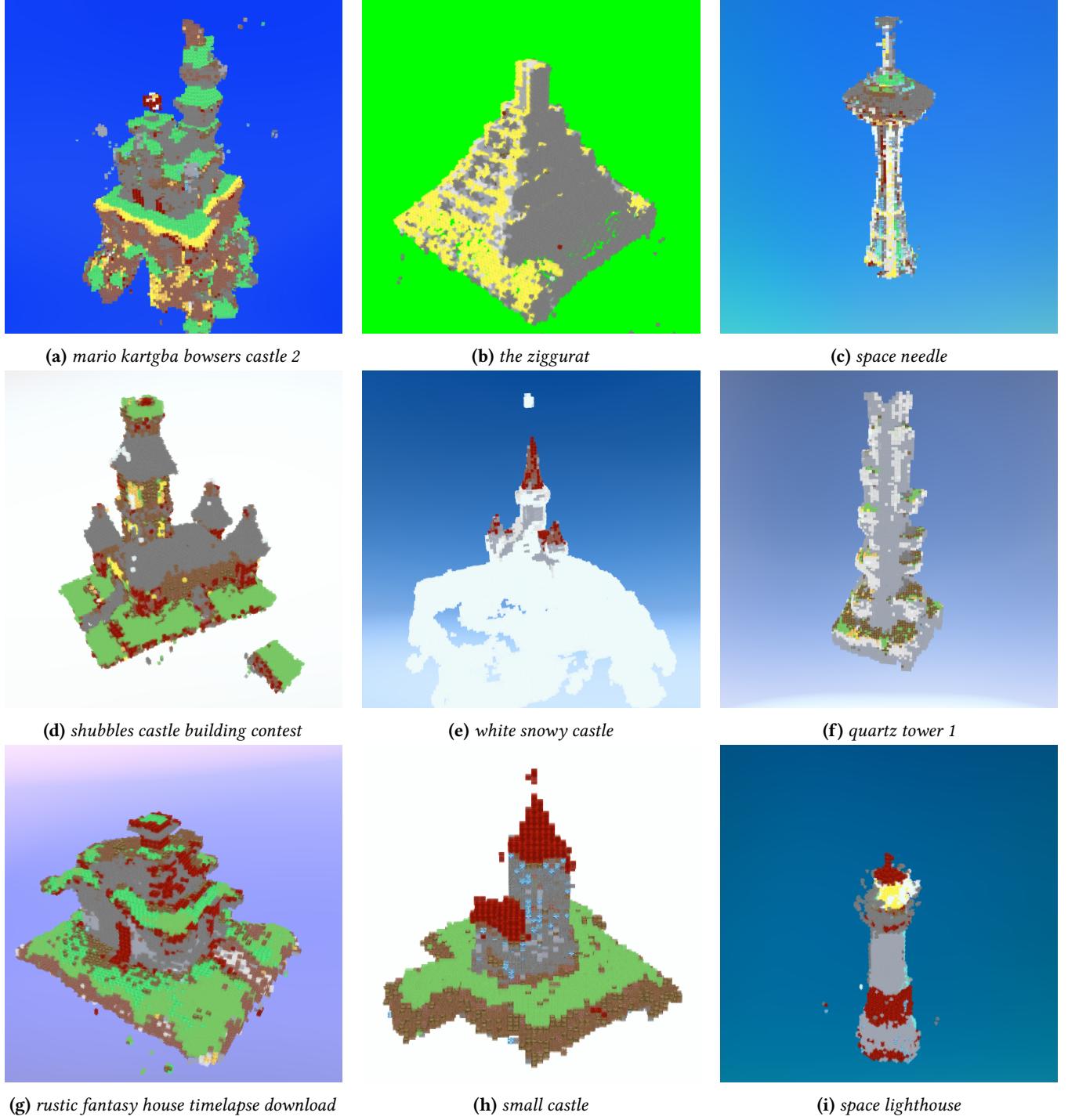
Sam Earle, Filippos Kokkinos, Yuhe Nie, Julian Togelius, and Roberta Raileanu



**(a)** *mario kartgba bowsers castle 2*

**(b)** *the ziggurat*

**(c)** *space needle*

**(d)** *shubbles castle building contest*

**(e)** *white snowy castle*

**(f)** *quartz tower 1*

**(g)** *rustic fantasy house timelapse download*

**(h)** *small castle*

**(i)** *space lighthouse*

**Figure 2.** DreamCraft neural rendered generations for a set of Planet Minecraft captions.

## 3 DreamCraft: Text-Guided Minecraft Environment Generation

### 3.1 Quantized NeRFs for Environment Generation

In this section, we introduce DreamCraft, a quantized NeRF which learns to arrange in-game assets during training (see

Figure 3 for an overview of our approach). Our text-guided NeRF implementation uses score distillation sampling from a pre-trained image generation model to provide a loss function for the optimization of the NeRF. We use a preliminary version of Emu [11], trained only on Shutterstock data. Note
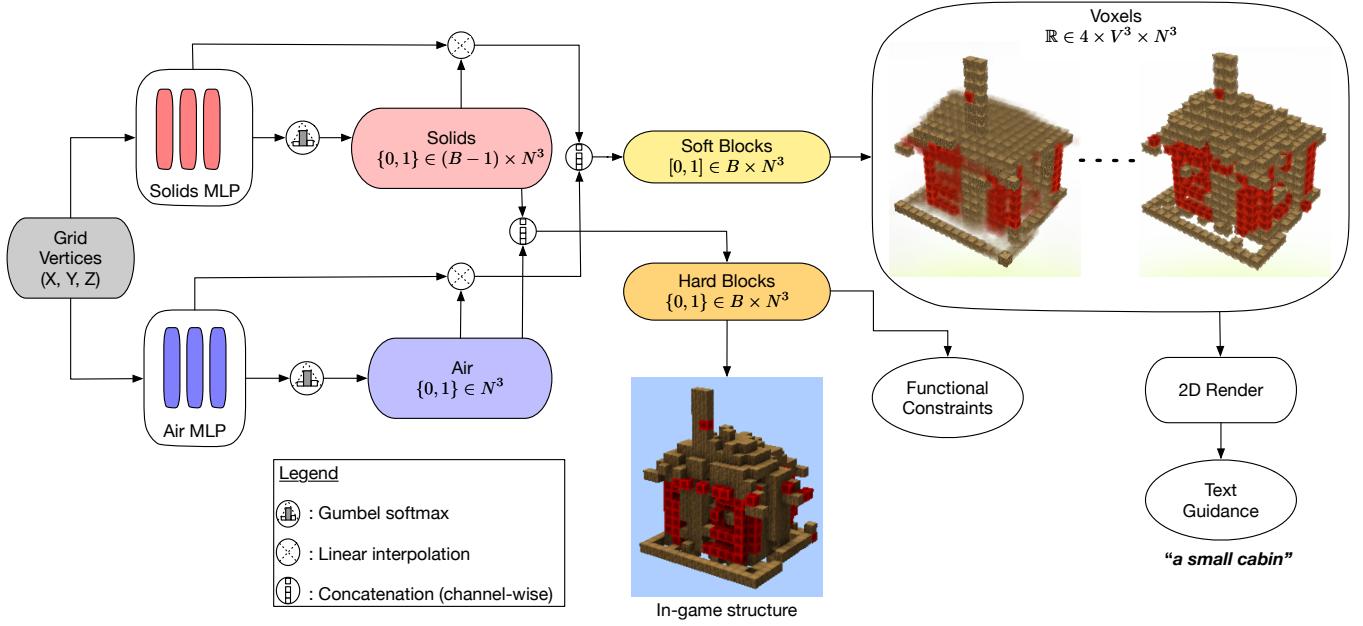
**Figure 3.** DREAMCRAFT produces a distribution of discrete blocks over a grid. Unquantized or "soft" versions of this distribution can be visualized for text-image guidance, while a fully discrete representation can be used to determine functional constraints, and exported to an in-game structure.

that our approach is agnostic to the text-to-3D model used, so it can be applied to any other text-to-3D model architecture and thus benefit from future advances in that field. Also note that we do not use any additional or domain-specific data to train our model apart from the block textures in Table 6.

**From Continuous Points to Discrete Blocks** As in PureCLIPNeRF, we use MLPs to predict continuous activations over a grid by feeding them $X, Y, Z$ coordinates which are encoded using a series of sine waves, as in the original NeRF [55]. We sample $X, Y, Z$ vertices along a cubic $3D$ grid of width $N$, generating predictions over solid block types at each cell in the grid. We refer to the resulting tensor as the soft *solid block grid* $\mathbf{B}_{soft} \in [0, 1]^{M \times N \times N \times N}$, of width $N$, comprising of $M$ different non-air block types. We take gumbel softmax [32] over these predictions to yield a discrete grid of solid blocks $\mathbf{B}_{hard}$ with values in $[0, 1]$. We feed the same $X, Y, Z$ coordinates to a separate MLP to generate a soft *air block grid* $\mathbf{A}_{soft} \in [0, 1]^{N \times N \times N}$, and again apply gumbel softmax to generate its discrete counterpart $\mathbf{A}_{hard} \in \{0, 1\}^{N \times N \times N}$. The distinction between air and solid blocks is analogous to that between the air and albedo MLPs in standard NeRFs.

During training, we may interpolate hard and soft variants to interpolate final air and block grids:

$$\mathbf{A} = \alpha \cdot \mathbf{A}_{hard} + (1 - \alpha) \cdot \mathbf{A}_{soft} \quad \mathbf{B} = \beta \cdot \mathbf{B}_{hard} + (1 - \beta) \cdot \mathbf{B}_{soft},$$

where $\alpha$ and $\beta$ in $[0, 1]$ control the hardness of air and solid blocks respectively. In our experiments, we set these to 1 and 0 when air/solids are hard (discrete) and soft (continuous)

respectively, and linearly scale between 0 and 1 to "anneal" the blocks from continuous to discrete. We then mask solid block activations with air activations to yield the final block grid $\mathbf{C} = \mathbf{A} \odot \mathbf{B}$, with $\mathbf{C} \in [0, 1]^{M \times N \times N \times N}$, where $\odot$ denotes the element-wise product, broadcasting over the block type dimension in $\mathbf{B}$. Thus, $\mathbf{C}$ can be seen as a 3D image with as many channels as there are solid block types. When a cell in $\mathbf{C}$ has a value of 1 in a given block channel, and 0s elsewhere, it contains only this block type. When it has values in $(0, 1)$, a block is "partially" present. And when a cell in $\mathbf{C}$ is all zeros, it represents the presence of an air block. When exporting structures to the game engine or computing loss from functional constraints, we set $\alpha = \beta = 1$, producing an entirely discrete block grid $\mathbf{C}_{hard}$.

**From 2D Textures to 3D Structures**

The block grid $\mathbf{C}$ is combined with statically-generated voxel grids to differentiably generate artifacts that visually resemble in-game structures. First, we pre-fabricate $16 \times 16 \times 16$ voxel grids for each block type using in-game textures, applying the game's $16 \times 16$ RGB textures to the appropriate block faces (picking an arbitrary priority order among faces where voxels overlap along the edges of the cube). These grids are frozen and do not pass gradients during learning. During the forward pass, we project our $M \times N \times N \times N$ block grid into a $16N \times 16N \times 16N$ voxelated block grid, where blocks appear in the same arrangement as in the low-resolution block grid, but in their voxelated form. We then apply neural ray tracing to this structure to generate 2D images.
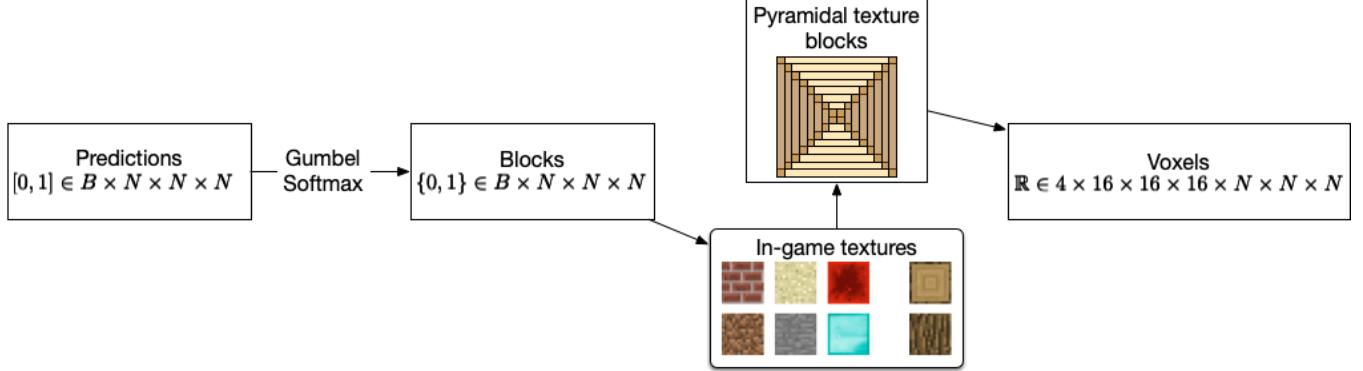
**Figure 4. From 2D textures to 3D structures** Once discretized, block types are mapped to voxel grids corresponding to the appearance of objects in-game. Here, we assemble each block into a 3D voxel grid using 2D textures from Minecraft. We approximate solid blocks by repeating the surface texture in the space inside the block.

When sampling points along a ray during rendering, each point takes the color and density values of the voxel in which it falls, avoiding interpolation between cells to mimic the sharp, pixelated appearance of textures in game. To avoid cases where the inside—but not the surface—of a block's voxel grid is sampled during ray tracing, we repeatedly stack face textures within the block, in a pyramidal pattern, to approximate solid objects (Figure 4).

Given a vertex $\mathbf{x} \in \mathbb{R}^3$ on the solid block grid, we compute a discrete block type, first using an MLP to generate a prediction $\mathbf{b}_{\text{soft}} \in \mathbb{R}^M$ over $M$ block types, then obtaining a onehot vector $\mathbf{b}_{\text{hard}} \in \{e_1, e_2, \ldots, e_M\}$ using the gumbel softmax function:

$$\mathbf{b}_{\text{soft}} = \text{MLP}(\mathbf{x}; \theta_B) \qquad \mathbf{b}_{\text{hard}} = \text{gumbel\_softmax}(\mathbf{b}_{\text{soft}}),$$

where $\theta_B$ denotes the parameters of the solid block type MLP. These vectors serve as the elements of $\mathbf{B}_{\text{soft}}$ and $\mathbf{B}_{\text{hard}}$ respectively.

Again given the block grid vertex $\mathbf{x}$, we compute soft and hard air values, where high values correspond to solid blocks and low values approaching 0 correspond to air.

The soft air grid is computed in the same way as density in previous works, passing the output of an MLP (parameterized by $\theta_A$) through an exponential activation function[1]:

$$\sigma_{\text{soft}} = \exp(\text{MLP}(\mathbf{x}; \theta_A)),$$

where $\theta_A$ denotes the parameters of the air block MLP. To quantize the air grid (effectively computing the presence of air blocks), we first use the soft air values to derive the respective probability of an air/solid block appearing at $\mathbf{x}$:

$$\sigma'_{\text{soft}} = \sigma_{\text{soft}} - 10 \qquad p_{\text{air}} = -\sigma'_{\text{soft}} \qquad p_{\text{solid}} = \sigma'_{\text{soft}}.$$

We then use the gumbel softmax function to discretize these predictions, obtaining a density value of 0 in case of

air, and 1 in case of a solid block:

$$\sigma_{\text{hard}} = \sum^{2} \text{gumbel\_softmax}([p_{\text{air}}, p_{\text{solid}}]).$$

These values serve as the elements of $\mathbf{A}_{\text{soft}}$ and $\mathbf{A}_{\text{hard}}$, respectively.

We use only opaque solid Minecraft blocks (excluding transparent blocks like glass or ice, and porous ones like grass or flowers). We use a separate, shallow background MLP, which takes as input a viewing angle, to model color at the end of each ray, allowing the model to learn a low-resolution background texture (effectively projected onto the inside of a sphere).

### 3.2 Functional Constraints

**Distributional Constraints** The discrete block grid resulting from the quantized NeRF allows us to optimize DREAM-CRAFT to produce a level satisfying a target distribution of block types, producing text-guided objects comprising of particular block mixtures. The user sets a target proportion for each block type, and we apply a loss equal to the difference between this target and the actual proportion of (non-air) grid cells in the NeRF's output (after quantization), which we compute by taking the sum over the relevant channel of the discrete onehot block grid and dividing by the size of the grid. Formally, we define the distributional loss as

$$L_D = \sum_{t \in T} |G(t) - P(t)|,$$

where $t$ is a block type among the set of blocks $T$, $G(t)$ is the target number of occurrences of this block as specified by the user, and $P(t)$ is the number of actual occurrence in the quantized block grid $\mathbf{C}_{\text{hard}}$.

**Adjacency Constraints** We also introduce a loss term corresponding to a penalty or reward incurred whenever a particular configuration of blocks appear in the generated structure. To this end, we construct a convolutional layer
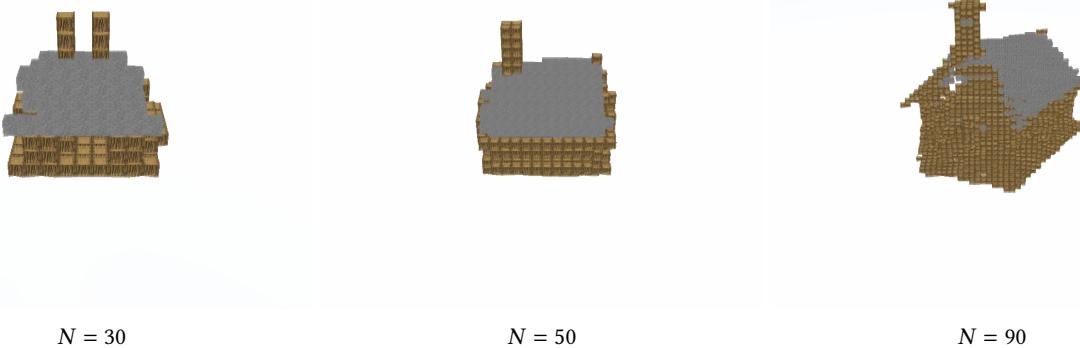
---

[1]During the backward pass, the exponential function is clipped to avoid exploding gradients: $\sigma_{\text{soft}} = \exp(\text{clamp}(y, 15))$.

$N = 30$     $N = 50$     $N = 90$

**Figure 5.** *"medium medieval home"* with block grids of various widths.

that outputs 1 over any matching patch of blocks, and 0 everywhere else, then sum the result, yielding the number of occurences of the relevant pattern in $C_{hard}$. We multiply this sum by a user-specified loss coefficient (negative when the pattern is desired, positive when prohibited). Suppose the user wants to apply a loss/reward of $l_p$ to the pattern of blocks $b_0, b_1, \ldots, b_{j_p}$ occupying a patch of size $K^3$ (with the number of blocks of interest $j_p \leq K^3$), where each block $b_i$ has relative coordinates $x_i, y_i, z_i$ in the patch. We construct a 3D convolutional weight matrix $W_p$ consisting of the one-hot vectors $e_i$ corresponding to each block type, and placing them at position $x_i, y_i, z_i$ in the weight matrix. We apply the resulting convolutional layer, $conv_{W_p}$ to the quantized block grid, subtract $j_p - 1$ from the output, and apply a ReLU activation to obtain the binary pattern-activation tensor. Formally,

$$L_P = \sum_{p \in P} w_p \sum_{i}^{K^3} \left( \text{ReLU} \left( \text{conv}_{W_p} \left( C_{hard} \right) - j_p + 1 \right) \right)_i,$$

where an inner sum is taken over elements $i$ in the binary activation tensor for pattern $p$.

## 4 Experiments

**Baselines** We compare the performance of DREAMCRAFT, our quantized NeRF which learns to arrange in-game assets *during training*, to **UNCONSTRAINED NERF**, a baseline that maps the continuous outputs to game assets *after training*. UNCONSTRAINED NERF trains a text-guided NeRF and then maps its output to Minecraft blocks using nearest neighbors. For text-guidance, we use a variant of Emu [11] trained only on Shutterstock, a model using text-conditioned latent diffusion to generate images. The nearest RGB value to each Minecraft block is determined by taking the average color over each of the (normally repeated) $16 \times 16$ textures covering each of its 6 faces. We define a width-$N$ grid over the 3D

output space of the unconstrained model and query the RGB and density values at the center of each cell in the grid. We calculate the $L_2$ distance between each centerpoint and average Minecraft block color, mapping each cell to the closest block. We then select a density threshold $s = 10$, and place air blocks wherever $\sigma < s$.

**Ablations** We study different quantization schemes in order to understand what is the best way to map the continuous outputs of the air and solid block MLPs into discrete grids of Minecraft blocks. The output of either MLP can be passed through the gumbel softmax function to produce a discrete grid of air or solid blocks (see Figure 3). If these values are not discretized *i.e.,* $\alpha < 1$ or $\beta < 1$, meaning that $\mathbf{b}_{soft}$ or $\sigma_{soft}$ are used instead of their "hard" counterparts, then the resulting voxel grid can include solid blocks interpolated with air blocks (*i.e.,* semi-transparent) or with one another (*i.e.,* multi-texture). We also experiment with linearly annealing these values from their soft to hard counterparts over the course of training.

**Evaluation Datasets** We evaluate our method on both generic and domain-specific text prompts, using the **COCO** [47] and **Planet Minecraft** [49] datasets, respectively. For COCO, we use the same 153 prompts as in prior text-to-3D works [43, 56, 61]. For Planet Minecraft, we take the names of the top 150 most downloaded assets uploaded by users in 2016 under the "Maps" category. Some examples include "a desperate and lonely wizards tower pmc chunk challenge entry lore", "mario kartgba bowsers castle 2", and "icarly set and nickelodeon studio". A full list of prompts can be found in Section D.

**Evaluation Metrics** To quantitatively evaluate the performance of our model, we measure its fidelity using **R-precision**. More specifically, we query other pre-trained joint text-image encoders, namely CLIP ViT-B/16 and CLIP

**Table 1. Fidelity of Neural Renders on COCO and Planet Minecraft** DreamCraft's quantized generations are compared to unquantized generations from an Unconstrained NeRF, which can be considered an upper bound for neural renders given their higher resolution. DreamCraft's relative performance increases when moving to a set of domain-relevant prompts.

| | COCO | | Planet Minecraft | |
| Model | CLIP ViT-B/16 | CLIP ViT-B/32 | CLIP ViT B/16 | CLiP ViT B/32 |
|---|---|---|---|---|
| Unconstrained NeRF | **61.44** | **66.67** | **25.17** | **31.29** |
| DreamCraft | 19.74 | 21.05 | 11.56 | 17.01 |
| ratio | 0.32 | 0.32 | 0.46 | 0.54 |

**Table 2. Fidelity of In-Game Renders on Planet Minecraft** DreamCraft, which learns to use discrete blocks during training, outperforms a baseline whose output is discretized only after training.

| Model | CLIP ViT-B/16 | CLIP ViT-B/32 |
|---|---|---|
| Unconstrained NeRF | 2.72 | 2.72 |
| DreamCraft | **5.44** | **6.12** |

**Table 3. Fidelity of Different Quantization Schemes for Planet Minecraft**. Using a hard block type and soft block density achieves the best performance.

| block type | block density | CLIP ViT-B/16 | | CLIP ViT-B/32 | |
| | | RGB | depth | RGB | depth |
|---|---|---|---|---|---|
| anneal | anneal | 7.73 | 5.07 | 9.07 | 5.73 |
| | hard | 2.67 | 1.20 | 5.60 | 2.00 |
| | soft | 7.33 | 7.33 | 10.40 | 5.20 |
| hard | anneal | 8.53 | 6.00 | 8.40 | 7.33 |
| | hard | 2.27 | 1.33 | 2.80 | 1.87 |
| | soft | **10.40** | **8.93** | **13.20** | **8.40** |
| soft | anneal | 8.40 | 3.60 | 10.00 | 4.80 |
| | hard | 4.00 | 0.80 | 7.60 | 2.80 |
| | soft | 7.87 | 6.53 | 11.07 | 5.47 |

**Table 4.** Fidelity of generated structures given different quantization schemes for block density and type. COCO dataset, neural renders. Both RGB screenshots and depth-only greyscale heatmaps of the generated 3D structure are evaluated. A combination of fully discrete block types and "soft" block density leads to best performance in terms of depth (i.e. structure topology), while soft block types and density during training lead to the best RGB images. R-precision averaged over 5 poses for each experiment.

| block type | block density | CLIP ViT-B/16 | | CLIP ViT-B/32 | |
| | | RGB | depth | RGB | depth |
|---|---|---|---|---|---|
| anneal | anneal | 12.68 | 4.31 | 11.37 | 3.01 |
| | hard | 5.36 | 0.65 | 6.80 | 0.92 |
| | soft | 22.88 | 8.89 | 22.88 | 8.24 |
| hard | anneal | 12.68 | 5.88 | 14.12 | 4.71 |
| | hard | 4.05 | 2.22 | 4.97 | 1.18 |
| | soft | 19.61 | **12.03** | 21.83 | **11.11** |
| soft | anneal | 17.65 | 4.84 | 15.69 | 4.44 |
| | hard | 9.80 | 1.05 | 7.32 | 0.92 |
| | soft | **24.31** | 8.76 | **24.97** | 7.32 |

ViT-B/32 [62], and test whether they can recognize the caption responsible for a given NeRF rendering from a set of distractors (other, randomly selected captions from the dataset). For each caption, we repeat the process for 5 different test images and average the result.

## 4.1 Quality of the Generations

In Figures 1 and 8, we can see that, using only Minecraft blocks, our model produces structures that are visually similar to those of the Unconstrained NeRF, for both generic

and domain-specific text prompts. Note that the Unconstrained NeRF model is a strong upper bound because it has a continuous and thus much larger output space (i.e., higher resolution) than our discretized DreamCraft model. DreamCraft's relative performance increases when moving to a set of domain-relevant prompts.

In Table 1, we compare the fidelity of DreamCraft with that of the Unconstrained NeRF on COCO and Planet Minecraft. Note that here, the generations are evaluated using the renders from the neural ray tracing engine rather than in-game renders. As expected, limiting the NeRF's output to a specific set of discretely assembled blocks drastically reduces its space of generations. This is reflected in DreamCraft's lower fidelity with respect to the Unconstrained NeRF, when generating objects from both the COCO dataset

and Planet Minecraft datasets. However, the performance gap between DʀᴇᴀᴍCʀᴀғᴛ and the Uɴᴄᴏɴsᴛʀᴀɪɴᴇᴅ NᴇRF is reduced when moving from the generic COCO dataset to the domain-specific Planet Minecraft dataset. This suggests that despite its restricted output space, DʀᴇᴀᴍCʀᴀғᴛ is particularly effective at generating high-quality structures when the input prompts are relevant to the (discrete) domain at hand.

In Table 2, we evaluate the R-precision using 2D captures of generated Minecraft block layouts in the game engine itself. Note that here, the generations are evaluated using the in-game renders rather than the neural renders. For game design applications, in-game render fidelity is more relevant than neural render fidelity, so this is ultimately the metric we care about in our study. In this case, the fidelity of the Unconstrained NeRF is lower than that of DʀᴇᴀᴍCʀᴀғᴛ for both COCO and Planet Minecraft. This indicates that postprocessing the output of a NeRF in order to discretize it using nearest-neighbor leads to worse results than learning to use discrete blocks during the generation process. In Figure 9, we see that mapping a discrete set of grid vertices to nearest neighbor block types via average color leads to sub-optimal results that are particularly bad at maintaining consistency in terms of texture and color. This result demonstrates the difficulty of translating unconstrained generations to a constrained repertoire of domain-specific assets. We conclude that by incorporating these game assets in the learning process, we can generate more faithful in-game structures.

## 4.2 Quantization Schemes

In Table 3, we compare the effect of applying soft, hard, and annealed quantization schemes to solid and air blocks. Maintaining *soft air blocks* (continuous-valued block transparency), in combination with *hard solid blocks* (discrete block types), leads to the highest R-precision at test time. This suggests that learning the topology (in contrast to the color/texture) of a generated structure is a sensitive process in which relaxing the quantization scheme (and presumably simplifying the loss landscape) is crucial.

Forcing block density to be discrete throughout training leads to poor performance (as can be seen in the first row of Figure 6), with the poorest performance coming from models in which both block type and density are fully discrete. This may be because of noisier learning dynamics resulting from the quantized output space of the model.

Conversely, using soft block density can lead to situations in which apparently solid surfaces are emulated by layering a number of semi-transparent blocks. At test time, when rendering the fully discrete block grid, such surfaces can suddenly be culled from the image, as none of the individual blocks of which they consist are enough to result in a "solid" binary output after quantization.
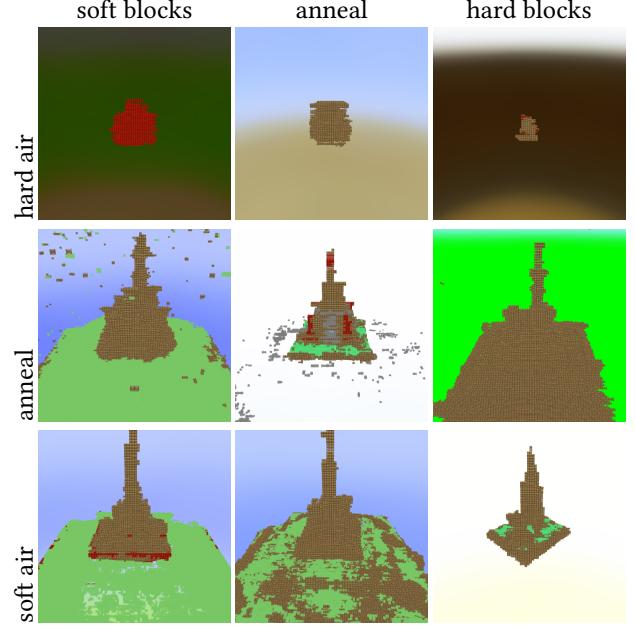


**Figure 6. Qualitative Effect of Different Quantization Schemes** for *"dwarven entrance"*. The best results obtained using hard block types and soft air blocks. Meaningful representations are learned only when the discreteness of air blocks (or by analogy, the density or topology of the generated structure) is relaxed or annealed throughout training (bottom two rows).
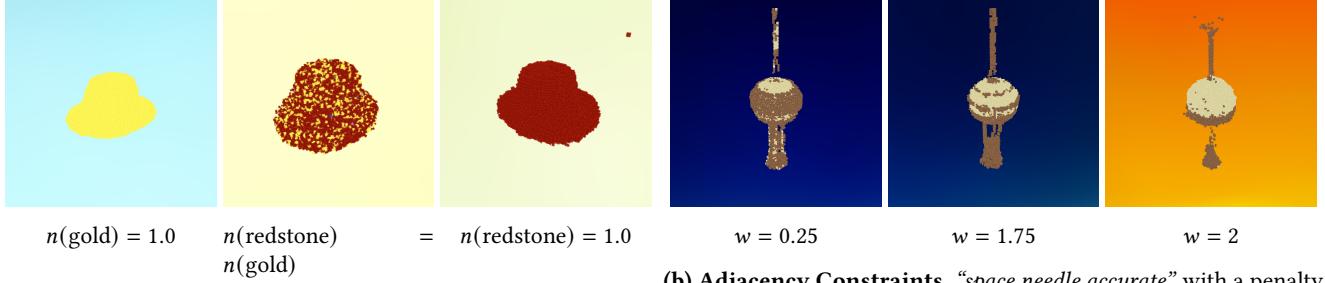
.

**Table 5.** Fidelity of generated structures given block grids of varying resolutions. Planet Minecraft dataset. Neural renders. Performance increases with the resolution of the block grid.

| | CLIP ViT-B/16 | | CLIP ViT-B/32 | |
|---|---|---|---|---|
| | RGB | depth | RGB | depth |
| block grid | | | | |
| 20 | 5.47 | 2.40 | 6.00 | 2.93 |
| 40 | 6.53 | 6.13 | 8.67 | 6.93 |
| 60 | 7.33 | 6.80 | 9.33 | 4.40 |
| 80 | 11.73 | 8.27 | 12.27 | 6.67 |
| 100 | **10.80** | **9.07** | **13.33** | 8.67 |

## 5 Block Grid Resolution

We experiment with the resolution of the block grid, learning a $N \times N \times N$-block representation of text prompts with $N \in \{10, 20, \ldots, 100\}$.

In Table 5, we see that increasing block grid resolution leads to an increase in R-precision. We note that the more blocks in the grid, the closer each block comes to being represented by only a single pixel in each 2D render of the

Sam Earle, Filippos Kokkinos, Yuhe Nie, Julian Togelius, and Roberta Raileanu



| | | | | | |
|---|---|---|---|---|---|
| $n(\text{gold}) = 1.0$ | $n(\text{redstone}) = n(\text{gold})$ | $n(\text{redstone}) = 1.0$ | $w = 0.25$ | $w = 1.75$ | $w = 2$ |

**(a) Distributional Constraints.** *"a stylish hat"* jointly optimized to satisfy a given target distribution over the block types.

**(b) Adjacency Constraints.** *"space needle accurate"* with a penalty for floating sand, weighted more heavily from left to right. The block set is limited to sand (light tan) and dirt (brown).

**Figure 7.** Examples of incorporation of functional constraints via auxiliary loss terms.

3D block layout. In other words, we can expect the output of these higher-resolution quantized NeRFs to approximate that of their unconstrained counterparts with increasing accuracy. By analogy with visual art, we can say that the model uses blocks in an increasingly pointillistic fashion.

### 5.1 Functional Constraints

In Figure 7a, we illustrate the effect of adding distributional constraints to a prompt asking for "a stylish hat". We can see that the model produces a similar structure using either entirely gold, redstone, or an even split of both when adding the distributional loss term.

In Figure 7b, we illustrate the effect of adding an adjacency constraint prohibiting sand from "floating", i.e. being placed directly above an air tile, and ask for "space needle accurate" (a Planet Minecraft prompt). As the weight of the adjacency loss term is increased, the use of sand blocks becomes restricted to the central "bulb" of the tower, where it is more likely to be supported by a the circular base of dirt blocks. When the adjacency loss is weighted less heavily, sand is often incorporated into the underside of the bulb (where it is unsupported and will fall to the ground in-game).

These experiments demonstrate that functional constraints can be easily integrated with our text-to-3D model to create controllable Minecraft structures that obey both high-level and low-level user specifications and can be rendered in-game.

## 6 Conclusion

In this work, we develop a new approach for generating functional game environments in Minecraft from free-form text descriptions. DreamCraft quantizes the output of a text-to-3D NeRF to predict discrete block types which are then mapped to game assets (*i.e.,* voxel-grids corresponding to in-game blocks). This allows the NeRF to use the game assets to represent the content described by a text prompt. We demonstrate that our approach has higher fidelity to the text prompt than a baseline that discretizes the output of an Unconstrained NeRF *after* learning. DreamCraft is,

to our knowledge, the first generator capable of generating diverse, functional, and controllable 3D game environments directly from free-form text. Since our model can adapt to the unique appearance of user-supplied modular game assets to produce environments with high-level aesthetic properties, it may be particularly useful for game designers working in new domains that don't yet have have large datasets of game layouts.

One limitation of DreamCraft is that it takes a few hours to generate a single structure. However, it could benefit from recent and future speed improvements in NeRFs [24, 81, 88, 92]. Another promising direction for future work is to model lightning and shadow, in addition to color and density, which could be achieved using an auxiliary loss to model the in-game lightning effects.

While we focus on MineCraft, DreamCraft could be extended beyond cube-based environments, to any 3D environment involving discrete assets that can be approximated by voxel grids. It could also incorporate more complex functional constraints, assuming these could be implemented to be differentiable. For example, one could compute the path length between key blocks (e.g. a player spawn block or a treasure chest) using convolutions, and use the difference between the target achieved path lengths as part of the loss, similar to the reward used in Earle et al. [14], Khalifa et al. [41]. It may also be beneficial—especially where functioial constraints cannot be made differentiable—to use RL methods to approximate the gradient from such additional functional scores. Vis-a-vis embodied player agents, environments could be generated to result in certain rewards, dynamics, regret or learnability with respect to these agents.

## References

[1] Alberto Alvarez, Steve Dahlskog, Jose Font, Johan Holmberg, and Simon Johansson. 2018. Assessing aesthetic criteria in the evolutionary dungeon designer. In *Proceedings of the 13th International Conference on the Foundations of Digital Games*. 1–4.

[2] Maren Awiszus, Frederik Schubert, and Bodo Rosenhahn. 2021. World-gan: a generative model for minecraft worlds. In *2021 IEEE Conference on Games (CoG)*. IEEE, 1–8.

[3] Bowen Baker, Ilge Akkaya, Peter Zhokov, Joost Huizinga, Jie Tang, Adrien Ecoffet, Brandon Houghton, Raul Sampedro, and Jeff Clune. 2022. Video pretraining (vpt): Learning to act by watching unlabeled online videos. *Advances in Neural Information Processing Systems* 35 (2022), 24639–24654.

[4] Philip Bontrager and Julian Togelius. 2021. Learning to generate levels from nothing. In *2021 IEEE Conference on Games (CoG)*. IEEE, 1–8.

[5] Nathan Brewer. 2017. Computerized Dungeons and Randomly Generated Worlds: From Rogue to Minecraft [Scanning Our Past]. *Proc. IEEE* 105, 5 (2017), 970–977.

[6] Michele Brocchini, Marco Mameli, Emanuele Balloni, Laura Della Sciucca, Luca Rossi, Marina Paolanti, Emanuele Frontoni, and Primo Zingaretti. 2022. MONstEr: A Deep Learning-Based System for the Automatic Generation of Gaming Assets. In *Image Analysis and Processing. ICIAP 2022 Workshops: ICIAP International Workshops, Lecce, Italy, May 23–27, 2022, Revised Selected Papers, Part I*. Springer, 280–290.

[7] Jake Bruce, Michael Dennis, Ashley Edwards, Jack Parker-Holder, Yuge Shi, Edward Hughes, Matthew Lai, Aditi Mavalankar, Richie Steigerwald, Chris Apps, et al. 2024. Genie: Generative Interactive Environments. *arXiv preprint arXiv:2402.15391* (2024).

[8] Alessandro Canossa and Gillian Smith. 2015. Towards a procedural evaluation technique: Metrics for level design. In *The 10th International Conference on the Foundations of Digital Games*. sn, 8.

[9] Anpei Chen, Zexiang Xu, Andreas Geiger, Jingyi Yu, and Hao Su. 2022. Tensorf: Tensorial radiance fields. In *Computer Vision–ECCV 2022: 17th European Conference, Tel Aviv, Israel, October 23–27, 2022, Proceedings, Part XXXII*. Springer, 333–350.

[10] Steve Dahlskog and Julian Togelius. 2014. Procedural content generation using patterns as objectives. In *Applications of Evolutionary Computation: 17th European Conference, EvoApplications 2014, Granada, Spain, April 23-25, 2014, Revised Selected Papers 17*. Springer, 325–336.

[11] Xiaoliang Dai, Ji Hou, Chih-Yao Ma, Sam Tsai, Jialiang Wang, Rui Wang, Peizhao Zhang, Simon Vandenhende, Xiaofang Wang, Abhimanyu Dubey, et al. 2023. Emu: Enhancing image generation models using photogenic needles in a haystack. *arXiv preprint arXiv:2309.15807* (2023).

[12] Michael Dennis, Natasha Jaques, Eugene Vinitsky, Alexandre Bayen, Stuart Russell, Andrew Critch, and Sergey Levine. 2020. Emergent complexity and zero-shot transfer via unsupervised environment design. *Advances in neural information processing systems* 33 (2020), 13049–13061.

[13] Robert Ota Dieterich. 2017. *Using Proof-Of-Concept Feedback to Explore the Relationship Between Artists and Procedural Content Generation in Computer Game Development Tools*. Ph. D. Dissertation.

[14] Sam Earle, Maria Edwards, Ahmed Khalifa, Philip Bontrager, and Julian Togelius. 2021. Learning controllable content generators. In *2021 IEEE Conference on Games (CoG)*. IEEE, 1–9.

[15] Sam Earle, Justin Snider, Matthew C Fontaine, Stefanos Nikolaidis, and Julian Togelius. 2022. Illuminating diverse neural cellular automata for level generation. In *Proceedings of the Genetic and Evolutionary Computation Conference*. 68–76.

[16] Linxi Fan, Guanzhi Wang, Yunfan Jiang, Ajay Mandlekar, Yuncong Yang, Haoyi Zhu, Andrew Tang, De-An Huang, Yuke Zhu, and Anima Anandkumar. 2022. Minedojo: Building open-ended embodied agents with internet-scale knowledge. *arXiv preprint arXiv:2206.08853* (2022).

[17] Tianhan Gao, Jin Zhang, and Qingwei Mi. 2022. Procedural Generation of Game Levels and Maps: A Review. In *2022 International Conference on Artificial Intelligence in Information and Communication (ICAIIC)*. IEEE, 050–055.

[18] Linus Gisslén, Andy Eakins, Camilo Gordillo, Joakim Bergdahl, and Konrad Tollmar. 2021. Adversarial reinforcement learning for procedural content generation. In *2021 IEEE Conference on Games (CoG)*. IEEE, 1–8.

[19] Ian Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2020. Generative adversarial networks. *Commun. ACM* 63, 11 (2020), 139–144.

[20] Daniele Gravina, Ahmed Khalifa, Antonios Liapis, Julian Togelius, and Georgios N Yannakakis. 2019. Procedural content generation through quality diversity. In *2019 IEEE Conference on Games (CoG)*. IEEE, 1–8.

[21] Djordje Grbic, Rasmus Berg Palm, Elias Najarro, Claire Glanois, and Sebastian Risi. 2021. Evocraft: A new challenge for open-endedness. In *Applications of Evolutionary Computation: 24th International Conference, EvoApplications 2021, Held as Part of EvoStar 2021, Virtual Event, April 7–9, 2021, Proceedings 24*. Springer, 325–340.

[22] Michael Cerny Green, Luvneesh Mugrai, Ahmed Khalifa, and Julian Togelius. 2020. Mario level generation from mechanics using scene stitching. In *2020 IEEE Conference on Games (CoG)*. IEEE, 49–56.

[23] Maxim Gumin. 2016. Wave Function Collapse Algorithm. https://github.com/mxgmn/WaveFunctionCollapse

[24] Xiang Guo, Guanying Chen, Yuchao Dai, Xiaoqing Ye, Jiadai Sun, Xiao Tan, and Errui Ding. 2022. Neural Deformable Voxel Grid for Fast Optimization of Dynamic View Synthesis. In *Proceedings of the Asian Conference on Computer Vision*. 3757–3775.

[25] William H Guss, Mario Ynocente Castro, Sam Devlin, Brandon Houghton, Noboru Sean Kuno, Crissman Loomis, Stephanie Milani, Sharada Mohanty, Keisuke Nakata, Ruslan Salakhutdinov, et al. 2021. The minerl 2020 competition on sample efficient reinforcement learning using human priors. *arXiv preprint arXiv:2101.11071* (2021).

[26] Matthew Guzdial, Duri Long, Christopher Cassion, and Abhishek Das. 2017. Visual procedural content generation with an artificial abstract artist. In *Proceedings of ICCC computational creativity and games workshop*.

[27] Matthew Guzdial, Sam Snodgrass, and Adam J Summerville. 2022. Constraint-Based PCGML Approaches. In *Procedural Content Generation via Machine Learning: An Overview*. Springer, 51–66.

[28] Matthew Guzdial, Sam Snodgrass, and Adam J Summerville. 2022. PCGML Process Overview. In *Procedural Content Generation via Machine Learning: An Overview*. Springer, 35–49.

[29] Matthew Guzdial, Sam Snodgrass, and Adam J Summerville. 2022. *Procedural Content Generation Via Machine Learning: An Overview*. Springer.

[30] Zekun Hao, Arun Mallya, Serge Belongie, and Ming-Yu Liu. 2021. Gancraft: Unsupervised 3d neural rendering of minecraft worlds. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 14072–14082.

[31] Mark Hendrikx, Sebastiaan Meijer, Joeri Van Der Velden, and Alexandru Iosup. 2013. Procedural content generation for games: A survey. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)* 9, 1 (2013), 1–22.

[32] Eric Jang, Shixiang Gu, and Ben Poole. 2017. Categorical Reparametrization with Gumble-Softmax. In *International Conference on Learning Representations (ICLR 2017)*. OpenReview. net.

[33] Minqi Jiang, Edward Grefenstette, and Tim Rocktäschel. 2021. Prioritized level replay. In *International Conference on Machine Learning*. PMLR, 4940–4950.

[34] Zehua Jiang, Sam Earle, Michael Green, and Julian Togelius. 2022. Learning Controllable 3D Level Generators. In *Proceedings of the 17th International Conference on the Foundations of Digital Games*. 1–9.

[35] Matthew Johnson, Katja Hofmann, Tim Hutton, and David Bignell. 2016. The Malmo Platform for Artificial Intelligence Experimentation.. In *Ijcai*. 4246–4247.

[36] Arthur Juliani, Ahmed Khalifa, Vincent-Pierre Berges, Jonathan Harper, Ervin Teng, Hunter Henry, Adam Crespi, Julian Togelius, and Danny Lange. 2019. Obstacle tower: A generalization challenge in vision, control, and planning. *arXiv preprint arXiv:1902.01378* (2019).

[37] Niels Justesen, Ruben Rodriguez Torrado, Philip Bontrager, Ahmed Khalifa, Julian Togelius, and Sebastian Risi. 2018. Illuminating generalization in deep reinforcement learning through procedural level generation. In *NeurIPS Workshop on Deep Reinforcement Learning*.

[38] Anssi Kanervisto, Stephanie Milani, Karolis Ramanauskas, Nicholay Topin, Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, Wei Yang, et al. 2022. Minerl diamond 2021 competition: Overview, results, and lessons learned. *NeurIPS 2021 Competitions and Demonstrations Track* (2022), 13–28.

[39] Anssi Kanervisto, Stephanie Milani, Karolis Ramanauskas, Nicholay Topin, Zichuan Lin, Junyou Li, Jianing Shi, Deheng Ye, Qiang Fu, Wei Yang, Weijun Hong, Zhongyue Huang, Haicheng Chen, Guangjun Zeng, Yue Lin, Vincent Micheli, Eloi Alonso, François Fleuret, Alexander Nikulin, Yury Belousov, Oleg Svidchenko, and Aleksei Shpilman. 2022. MineRL Diamond 2021 Competition: Overview, Results, and Lessons Learned. In *Proceedings of the NeurIPS 2021 Competitions and Demonstrations Track (Proceedings of Machine Learning Research, Vol. 176)*, Douwe Kiela, Marco Ciccone, and Barbara Caputo (Eds.). PMLR, 13–28. https://proceedings.mlr.press/v176/kanervisto22a.html

[40] Isaac Karth and Adam M Smith. 2019. Addressing the fundamental tension of PCGML with discriminative learning. In *Proceedings of the 14th International Conference on the Foundations of Digital Games*. 1–9.

[41] Ahmed Khalifa, Philip Bontrager, Sam Earle, and Julian Togelius. 2020. Pcgrl: Procedural content generation via reinforcement learning. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 16. 95–101.

[42] Heinrich Küttler, Nantas Nardelli, Alexander Miller, Roberta Raileanu, Marco Selvatici, Edward Grefenstette, and Tim Rocktäschel. 2020. The nethack learning environment. *Advances in Neural Information Processing Systems* 33 (2020), 7671–7684.

[43] Han-Hung Lee and Angel X Chang. 2022. Understanding pure clip guidance for voxel grid nerf models. *arXiv preprint arXiv:2209.15172* (2022).

[44] Vivian Lee, Nathan Partlan, and Seth Cooper. 2020. Precomputing Player Movement in Platformers for Level Generation with Reachability Constraints.. In *AIIDE Workshops*.

[45] Antonios Liapis, Georgios Yannakakis, and Julian Togelius. 2013. Designer modeling for personalized game content creation tools. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 9. 11–16.

[46] Antonios Liapis, Georgios N Yannakakis, and Julian Togelius. 2012. Adapting models of visual aesthetics for personalized content creation. *IEEE Transactions on Computational Intelligence and AI in Games* 4, 3 (2012), 213–228.

[47] Tsung-Yi Lin, Michael Maire, Serge Belongie, James Hays, Pietro Perona, Deva Ramanan, Piotr Dollár, and C Lawrence Zitnick. 2014. Microsoft coco: Common objects in context. In *European conference on computer vision*. Springer, 740–755.

[48] Jialin Liu, Sam Snodgrass, Ahmed Khalifa, Sebastian Risi, Georgios N Yannakakis, and Julian Togelius. 2021. Deep learning for procedural content generation. *Neural Computing and Applications* 33, 1 (2021), 19–37.

[49] Cyprezz LLC. [n. d.]. Planet Minecraft Community: Creative fansite for everything minecraft! https://www.planetminecraft.com/.

[50] Christian E López, James Cunningham, Omar Ashour, and Conrad S Tucker. 2020. Deep reinforcement learning for procedural content generation of 3d virtual environments. *Journal of Computing and Information Science in Engineering* 20, 5 (2020).

[51] Alejandro Medina, Melanie Richey, Mark Mueller, and Jacob Schrum. 2023. Evolving Flying Machines in Minecraft Using Quality Diversity. *arXiv preprint arXiv:2302.00782* (2023).

[52] Timothy Merino, M Charity, and Julian Togelius. 2023. Interactive Latent Variable Evolution for the Generation of Minecraft Structures. In *Proceedings of the 18th International Conference on the Foundations*

[53] Timothy Merino, Roman Negri, Dipika Rajesh, M Charity, and Julian Togelius. 2023. The Five-Dollar Model: Generating Game Maps and Sprites from Sentence Embeddings. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 19. 107–115.

[54] Stephanie Milani, Nicholay Topin, Brandon Houghton, William H Guss, Sharada P Mohanty, Keisuke Nakata, Oriol Vinyals, and Noboru Sean Kuno. 2020. Retrospective analysis of the 2019 MineRL competition on sample efficient reinforcement learning. In *NeurIPS 2019 Competition and Demonstration Track*. PMLR, 203–214.

[55] B Mildenhall, PP Srinivasan, M Tancik, JT Barron, R Ramamoorthi, and R Ng. 2020. Nerf: Representing scenes as neural radiance fields for view synthesis. In *European conference on computer vision*.

[56] Nasir Mohammad Khalid, Tianhao Xie, Eugene Belilovsky, and Tiberiu Popa. 2022. CLIP-Mesh: Generating textured meshes from text using pretrained image-text models. In *SIGGRAPH Asia 2022 Conference Papers*. 1–8.

[57] Justin Mott, Saujas Nandi, and Luke Zeller. 2019. Controllable and coherent level generation: A two-pronged approach. In *Experimental AI in games workshop*.

[58] Rohit Nair. 2020. *Using Raymarched shaders as environments in 3D video games*. Drexel University.

[59] Mark J Nelson, Julian Togelius, Cameron Browne, and Michael Cook. 2016. Rules and mechanics. *Procedural Content Generation in Games* (2016), 99–121.

[60] Jack Parker-Holder, Minqi Jiang, Michael Dennis, Mikayel Samvelyan, Jakob Foerster, Edward Grefenstette, and Tim Rocktäschel. 2022. Evolving curricula with regret-based environment design. In *International Conference on Machine Learning*. PMLR, 17473–17498.

[61] Ben Poole, Ajay Jain, Jonathan T Barron, and Ben Mildenhall. 2022. Dreamfusion: Text-to-3d using 2d diffusion. *arXiv preprint arXiv:2209.14988* (2022).

[62] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, et al. 2021. Learning transferable visual models from natural language supervision. In *International conference on machine learning*. PMLR, 8748–8763.

[63] Sebastian Risi and Julian Togelius. 2020. Increasing generality in machine learning through procedural content generation. *Nature Machine Intelligence* 2, 8 (2020), 428–436.

[64] Robin Rombach, Andreas Blattmann, Dominik Lorenz, Patrick Esser, and Björn Ommer. 2022. High-Resolution Image Synthesis With Latent Diffusion Models. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 10684–10695.

[65] Christoph Salge, Claus Aranha, Adrian Brightmoore, Sean Butler, Rodrigo De Moura Canaan, Michael Cook, Michael Green, Hagen Fischer, Christian Guckelsberger, Jupiter Hadley, et al. 2022. Impressions of the GDMC AI Settlement Generation Challenge in Minecraft. In *Proceedings of the 17th International Conference on the Foundations of Digital Games*. 1–16.

[66] Christoph Salge, Michael Cerny Green, Rodrigo Canaan, and Julian Togelius. 2018. Generative design in minecraft (gdmc) settlement generation competition. In *Proceedings of the 13th International Conference on the Foundations of Digital Games*. 1–10.

[67] Mikayel Samvelyan, Robert Kirk, Vitaly Kurin, Jack Parker-Holder, Minqi Jiang, Eric Hambro, Fabio Petroni, Heinrich Küttler, Edward Grefenstette, and Tim Rocktäschel. 2021. Minihack the planet: A sandbox for open-ended reinforcement learning research. *arXiv preprint arXiv:2109.13202* (2021).

[68] Anurag Sarkar and Seth Cooper. 2020. Sequential segment-based level generation and blending using variational autoencoders. In *Proceedings of the 15th International Conference on the Foundations of Digital Games*. 1–9.

[69] Anurag Sarkar, Zhihan Yang, and Seth Cooper. 2020. Conditional level generation and game blending. *arXiv preprint arXiv:2010.07735* (2020).

[70] Noor Shaker, Julian Togelius, and Mark J Nelson. 2016. Procedural content generation in games. (2016).

[71] Noor Shaker, Julian Togelius, Mark J Nelson, Antonios Liapis, Gillian Smith, and Noor Shaker. 2016. Mixed-initiative content creation. *Procedural content generation in games* (2016), 195–214.

[72] Noor Shaker, Georgios Yannakakis, and Julian Togelius. 2010. Towards automatic personalized content generation for platform games. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 6. 63–68.

[73] Uriel Singer, Adam Polyak, Thomas Hayes, Xi Yin, Jie An, Songyang Zhang, Qiyuan Hu, Harry Yang, Oron Ashual, Oran Gafni, et al. 2022. Make-a-video: Text-to-video generation without text-video data. *arXiv preprint arXiv:2209.14792* (2022).

[74] Matthew Siper, Ahmed Khalifa, and Julian Togelius. 2022. Path of Destruction: Learning an Iterative Level Generator Using a Small Dataset. *arXiv preprint arXiv:2202.10184* (2022).

[75] Ole Edvin Skjeltorp. 2022. *3D Neural Cellular Automata-Simulating morphogenesis: Shape, color and behavior of three-dimensional structures.* Master's thesis.

[76] 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.

[77] Shyam Sudhakaran, Miguel González-Duque, Claire Glanois, Matthias Freiberger, Elias Najarro, and Sebastian Risi. 2023. MarioGPT: Open-Ended Text2Level Generation through Large Language Models. arXiv:2302.05981 [cs.AI]

[78] Shyam Sudhakaran, Djordje Grbic, Siyan Li, Adam Katona, Elias Najarro, Claire Glanois, and Sebastian Risi. 2021. Growing 3d artefacts and functional machines with neural cellular automata. *arXiv preprint arXiv:2103.08737* (2021).

[79] Adam Summerville and Michael Mateas. 2015. Sampling hyrule: Multi-technique probabilistic level generation for action role playing games. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 11. 63–67.

[80] Adam Summerville, Sam Snodgrass, Matthew Guzdial, Christoffer Holmgård, Amy K Hoover, Aaron Isaksen, Andy Nealen, and Julian Togelius. 2018. Procedural content generation via machine learning (PCGML). *IEEE Transactions on Games* 10, 3 (2018), 257–270.

[81] Cheng Sun, Min Sun, and Hwann-Tzong Chen. 2022. Direct voxel grid optimization: Super-fast convergence for radiance fields reconstruction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5459–5469.

[82] Adaptive Agent Team, Jakob Bauer, Kate Baumli, Satinder Baveja, Feryal Behbahani, Avishkar Bhoopchand, Nathalie Bradley-Schmieg, Michael Chang, Natalie Clay, Adrian Collister, et al. 2023. Human-Timescale Adaptation in an Open-Ended Task Space. *arXiv preprint arXiv:2301.07608* (2023).

[83] Open Ended Learning Team, Adam Stooke, Anuj Mahajan, Catarina Barros, Charlie Deck, Jakob Bauer, Jakub Sygnowski, Maja Trebacz, Max Jaderberg, Michael Mathieu, et al. 2021. Open-ended learning leads to generally capable agents. *arXiv preprint arXiv:2107.12808* (2021).

[84] Graham Todd, Sam Earle, Muhammad Umair Nasir, Michael Cerny Green, and Julian Togelius. 2023. Level Generation Through Large Language Models. In *Proceedings of the 18th International Conference on the Foundations of Digital Games*. 1–8.

[85] Julian Togelius, Georgios N Yannakakis, Kenneth O Stanley, and Cameron Browne. 2011. Search-based procedural content generation: A taxonomy and survey. *IEEE Transactions on Computational Intelligence and AI in Games* 3, 3 (2011), 172–186.

[86] Ruben Rodriguez Torrado, Ahmed Khalifa, Michael Cerny Green, Niels Justesen, Sebastian Risi, and Julian Togelius. 2020. Bootstrapping conditional gans for video game level generation. In *2020 IEEE Conference on Games (CoG)*. IEEE, 41–48.

[87] Peihao Wang, Xuxi Chen, Tianlong Chen, Subhashini Venugopalan, Zhangyang Wang, et al. 2022. Is Attention All NeRF Needs? *arXiv preprint arXiv:2207.13298* (2022).

[88] Peng Wang, Yuan Liu, Zhaoxi Chen, Lingjie Liu, Ziwei Liu, Taku Komura, Christian Theobalt, and Wenping Wang. 2023. F$^2$-NeRF: Fast Neural Radiance Field Training with Free Camera Trajectories. *arXiv preprint arXiv:2303.15951* (2023).

[89] Ryan Watkins. 2016. *Procedural content generation for unity game development.* Packt Publishing Ltd.

[90] Georgios N Yannakakis and Julian Togelius. 2011. Experience-driven procedural content generation. *IEEE Transactions on Affective Computing* 2, 3 (2011), 147–161.

[91] Cristopher Yates. 2021. *The use of Poisson Disc Distribution and A\* Pathfinding for Procedural Content Generation in Minecraft.* Ph. D. Dissertation. Ph. D. Dissertation. Memorial University.

[92] Hejia Zhang, Matthew Fontaine, Amy Hoover, Julian Togelius, Bistra Dilkina, and Stefanos Nikolaidis. 2020. Video game level repair via mixed integer linear programming. In *Proceedings of the AAAI Conference on Artificial Intelligence and Interactive Digital Entertainment*, Vol. 16. 151–158.

[93] Alexander Zook and Mark O Riedl. 2014. Generating and adapting game mechanics. In *Proceedings of the 2014 Foundations of Digital Games Workshop on Procedural Content Generation in Games*.

## A Minecraft Textures

The set of Minecraft textures used to imitate in-game blocks within the neural rendering engine is displayed in Table 6.
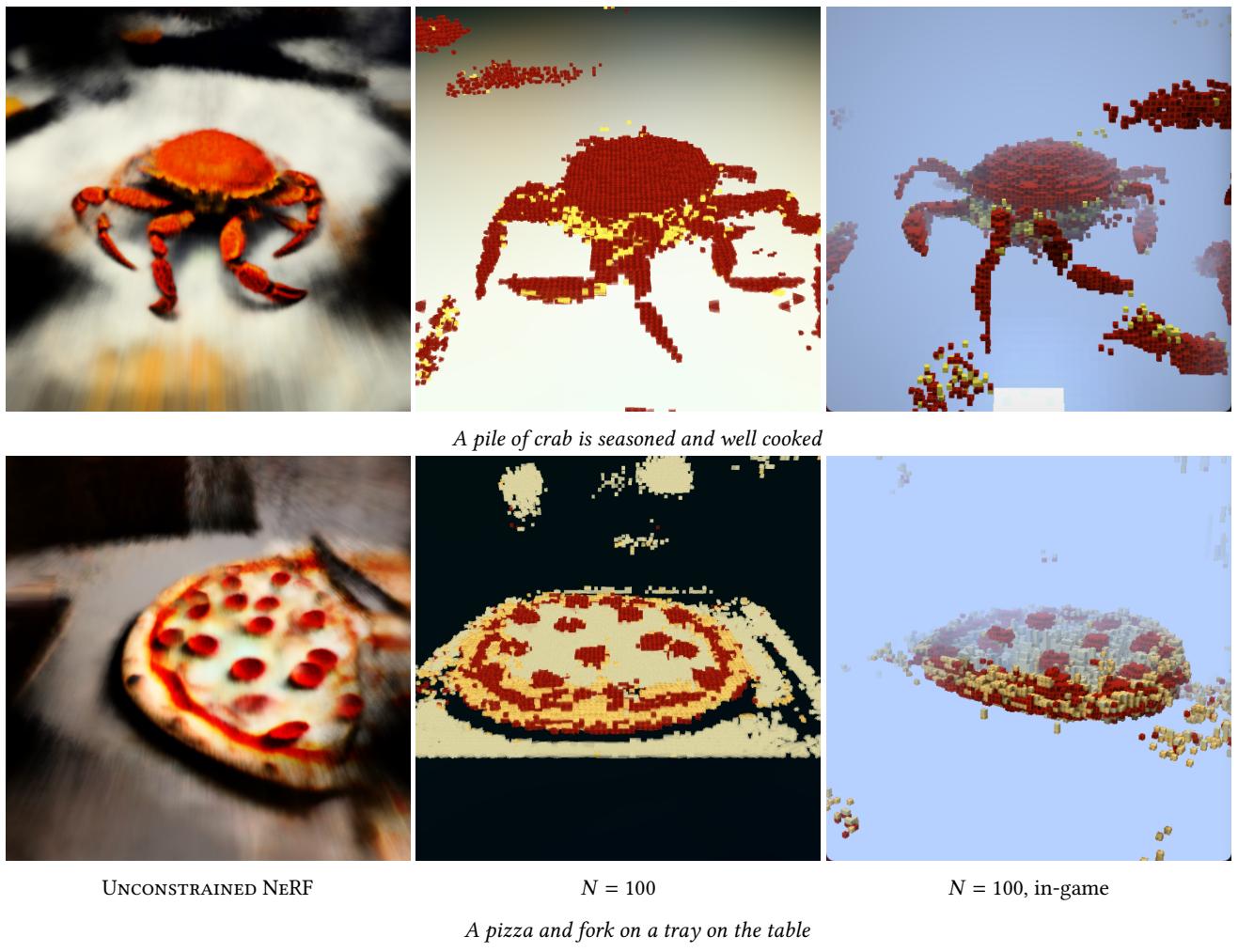
## B DREAMCRAFT Generations

## C Limitations

In some cases, the model uses negative space to represent an object, modulating the background texture to a particular color, then occluding parts of it with foreground blocks/density, to give it an apparent shape. This undesirable swapping of roles between foreground/background models may be more likely to occur in the quantized NeRF: whereas certain colors or textures may be difficult or impossible to replicate using the provided blocks, the background MLP remains unconstrained. To mitigate this, future work could investigate constraining the background MLP to only use 2D projections of "distant" game assets.

Another potential issue is the lack of semantic grounding with respect to block types. For example the model may just as well satisfy the prompt "large medieval ship" by using a combination of dirt and redstone, as with actual wooden logs or planks, so long as these give the *appearance* of wood. Our preliminary work on functional constraints suggests that this particular problem can be addressed by setting per-block-type targets (e.g. requiring 0% dirt and 50% wood blocks), but a more general approach might lie in "demonstrating" what each block type should be used to represent by adding this information in the prompt.

| block | texture | block | texture | block | texture |
|---|---|---|---|---|---|
| log oak | | stone | | dirt | |
| brick | | clay | | snow | |
| glazed terracotta light blue | | glazed terracotta yellow | | redstone block | |
| gold block | | iron block | | diamond block | |
| emerald block | | cobblestone | | slime | |

**Table 6.** In-game blocks and textures used by DREAMCRAFT



*A pile of crab is seasoned and well cooked*



UNCONSTRAINED NERF                    $N = 100$                    $N = 100$, in-game

*A pizza and fork on a tray on the table*

**Figure 8.** DREAMCRAFT output given COCO dataset captions, viewed in-game (right), by neural rendering (middle), and standard text-guided NeRF output (left) given the same prompts.

Whereas traditional NeRFs can model lighting and shadows, this is not the case in DREAMCRAFT, where the color at each point in 3D space is derived directly from a voxel grid corresponding to the in-game appearance of a Minecraft block. When structures are rendered inside the neural engine, they thus appear "flat" in contrast to the kind of shadow and

**(a)** *a japanese temple*



Unconstrained text-guided NeRF

text-guided NeRF with post hoc quantization
using nearest neighbor mapping

DreamCraft

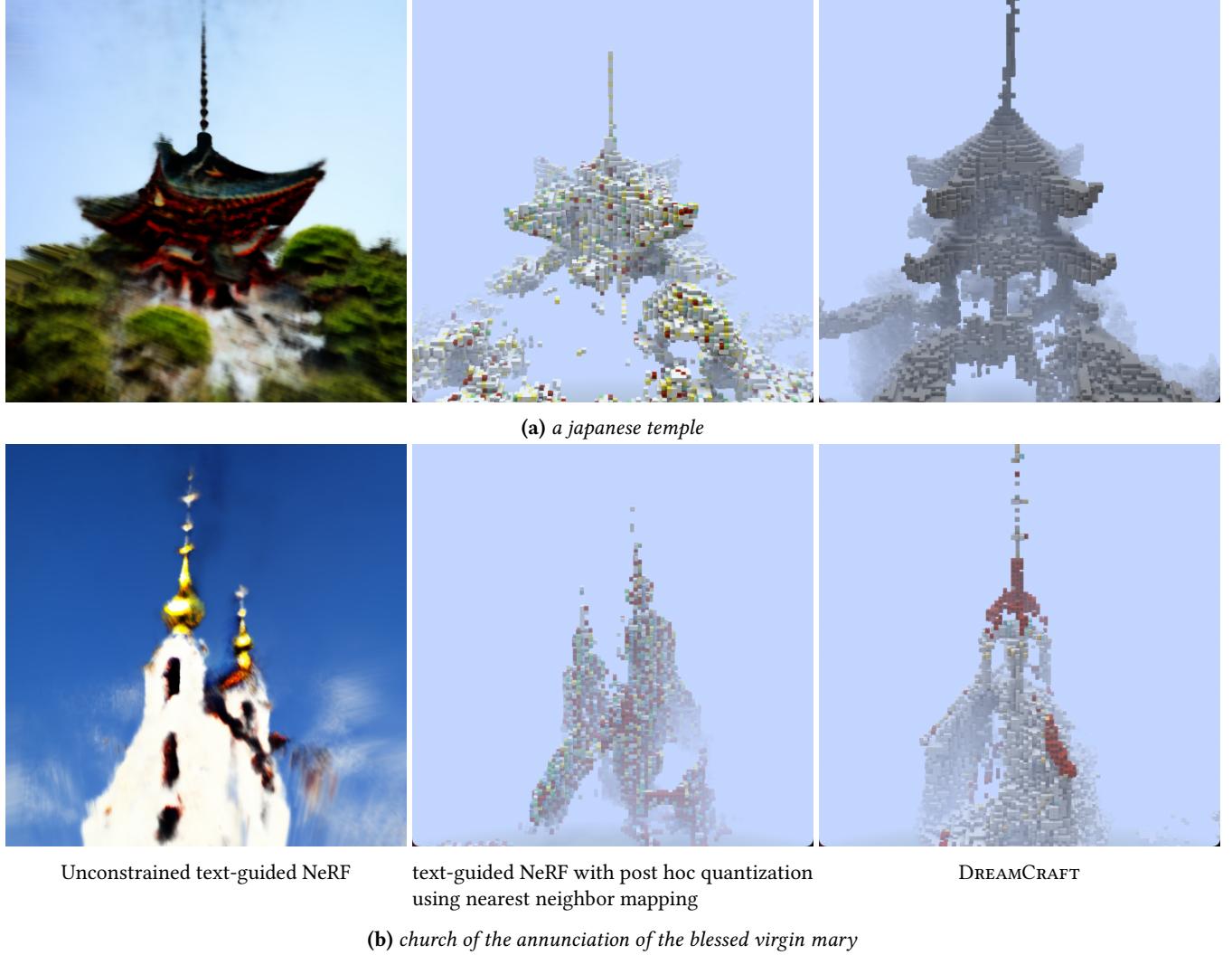**(b)** *church of the annunciation of the blessed virgin mary*

**Figure 9.** DreamCraft directly optimizes a representation using Minecraft blocks, leading to a more faithful reconstruction of the text prompt than results from post-processing the output of an Unconstrained NeRF, which leads to less coherent block type selection and structure topology. Captions from the planet minecraft dataset.

lighting effects that appear in the Minecraft game engine. Ideally, we could train an auxiliary model to mimic the effects of in-game lighting, for example by training it on paired datapoints of 3D block grids, and their appearance in-game at various angles. This learned renderer could replace the differentiable raycasting component of the NeRF pipeline (as in [87]), further sparing us from having to re-implement the rendering of irregular game objects such as plants and glass.

DreamCraft is currently too slow to be feasibly used in an online player-environment generator loop, taking a few hours to generate a single structure. Future versions could benefit from recent and future speed improvements in NeRFs [24, 81, 88, 92]. Alternatively, it could be leveraged to generate a training set for a conditional, guidance-free generative models of game worlds.

## D  Planet Minecraft Dataset

To test DreamCraft's ability to generate environments specific to the domain for which it was designed, we source text prompts from Planet Minecraft, a fan-operated site where users can upload and share custom content. We consider a subset of assets uploaded to the "Maps" category in 2016 (the year in which the most such assets were uploaded), and select the top 150 maps of this subset as measued by the number of user downloads. The prompts correspond to the names of these assets. We do not collect the assets themselves or any further data from the site.

The set of prompts scraped from Planet Minecraft is given below:

| paris eiffel tower | la valle dor by mrbatou download cinematic | reims cathedral | summit creative house | mexican hacjenda | coruscant senate building | from my house to yours merry christmas | the ziggurat | polaris skyscraper 25 | the craftsmans abode pmc solo contest 4 | rustic fantasy house timelapse download | skyscraper 31 ias | chateau de silveberg | fuminsh the city that never sleeps | ontario tower | dirt modern house | farin rocks | elven tower of the wise | distorsion chunk challenge | tours thiers nancy france skyscraper 3 ias | luxury beach house | space lighthouse | central place modern office complex | minecraft is a small world | tiger ii 101 scale | shurwyth snowlands download weareconquest | bridges | jurassic world v2 for jurassicraft 20 minecraft dinosaurs jurassic park isla nublar | brynwalda survival map | icarly set and nickelodeon studio | battle of hoth map echo base star wars hoth map | dream | the little castle nebelburg | steampunk island | land of azorth | jaws ride and amity village | avatar base | quartz tower 1 | small modern house 3 full interior | minecraft disneyland 1965 | 2012 skyrim | greenfield project neoclassical house | 2012 sea dogs village | eternal haven | patronis | japanese temple | star labs the flash cw | hub spawn | church of the annunciation of the blessed virgin mary inowrocaw | space needle accurate | the dark city hokkaido | a watch tower inspired by the game firewatch | undertale | white snowy castle | avengers tower | fantme villa modern house 2 | abandoned wild west | greenfield building vista creek elementary school | fantasy bundle level 25 special | a modern house 1 | server spawn by infro_ | fantasy inspired village danjgames | classic american farm | the builders shrine chunk challange | der eisendrache | modern house by real architect | ahzvels hq download re upload | download a medieval detached farm showcase | wg tower vice city | calypso a modern villa | tf2 egypt | minigames map atlantis | large medieval ship | small cabin | sustainable city | large medieval ship | sequoia valley 30 | kraehenfels survival version | small hospital | panem mc 1st quarter quell arena download | a desperate and lonely wizards tower pmc chunk challenge entry lore | paper mariocolour splash port prisma | skyscraper planus | hollywood residence | a nordic mountain village | old wizards tree mansion series 2 build 1read description | hidden in the sand | five nights at candys 2 roleplay map | castle ardor | epic server spawn | medium medieval home | medeival windmill 1102 | grand stadium pixelmon | prison mine 2 with download | babylon gardens | the new world trade center 11 | 30days day 28 orcish butchers slaughter house | old west home | ark labs outpost 26 | icebornminigames map | grand university | medieval mountain castle | minecraft maps | victorian manor | 30days day 8 dwarven entrance | spherical greenhouse 18 19 110 | huge minecraft server spawn airidale | big cottage | greenfield typical victorian | old fortress | modern condoapartment | small castle | ss tropic a custom by prestogo | two story old west shop | citadel hill fort george | batman arkham asylum | forest cottage | mountain temple | the conjuring | northwich | the amazing word of gumball the wattersons house 18 | five nights at freddys minecraft map 19 and above | fnaf roleplay map | middle eastern farm | modern house build sorry about no music in the video | a medieval farm | woodland mansion | survival map jairus isle | mario kartgba bowsers castle 2 | the scandinavian townhouse | mesa fort | three cool wood structures | brenttwood estates | toy shop 192 | dota 2 | server shop | store fletchers retreat | custom realistic terrain | the piggy sphinx | shubbles castle building contest | mountain housecastle 1 | roman outpost | fallout 4 red rocket | fantasy house | savanna village in the sky cinematic download | kahuai city | minecraft lets build timelapse fantasy update 12 over hanging house | kent regional airport | medieval house | redstone bunker 13 redstone creations version 1

DREAMCRAFT prompts from Planet Minecraft