# BlenderAlchemy: Editing 3D Graphics with Vision-Language Models

Ian Huang, Guandao Yang, and Leonidas Guibas

Department of Computer Science, Stanford University

**Abstract.** Graphics design is important for various applications, including movie production and game design. To create a high-quality scene, designers usually need to spend hours in software like Blender, in which they might need to interleave and repeat operations, such as connecting material nodes, hundreds of times. Moreover, slightly different design goals may require completely different sequences, making automation difficult. In this paper, we propose a system that leverages Vision-Language Models (VLMs), like GPT-4V, to intelligently search the design action space to arrive at an answer that can satisfy a user's intent. Specifically, we design a vision-based edit generator and state evaluator to work together to find the correct sequence of actions to achieve the goal. Inspired by the role of visual imagination in the human design process, we supplement the visual reasoning capabilities of VLMs with "imagined" reference images from image-generation models, providing visual grounding of abstract language descriptions. In this paper, we provide empirical evidence suggesting our system can produce simple but tedious Blender editing sequences for tasks such as editing procedural materials from text and/or reference images, as well as adjusting lighting configurations for product renderings in complex scenes. [1]
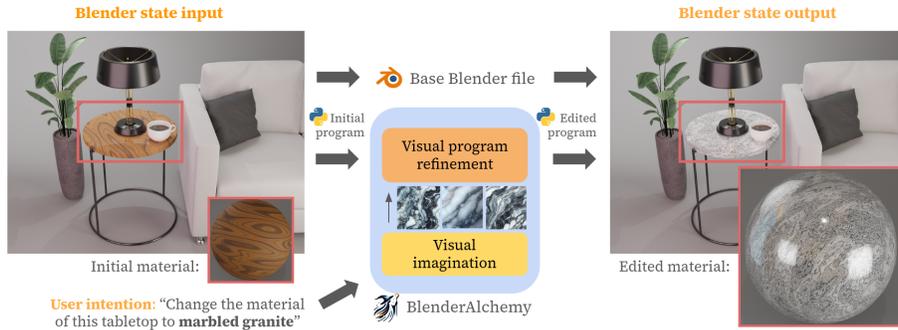
## 1 Introduction

To produce the compelling graphics content we see in movies or video games, 3D artists usually need to spend hours in software like Blender to find appropriate surface materials, object placements, and lighting arrangements. These operations require the artist to create a mental picture of the target, experiment with different parameters, and visually examine whether their edits get closer to the end goal. One can imagine automating these processes by converting language or visual descriptions of user intent into edits that achieve a design goal. Such a system can improve the productivity of millions of 3D designers and impact various industries that depend on 3D graphic design.

Graphic design is very challenging because even a small design goal requires performing a variety of different tasks. For instance, modeling of a game environment requires the 3D artist to cycle between performing modeling, material

---

[1] For project website and code, please go to: https://ianhuang0630.github.io/BlenderAlchemyWeb/
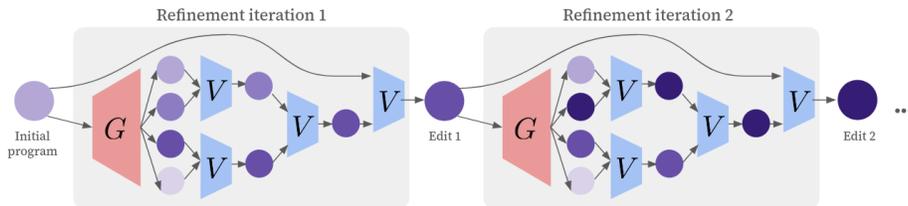
**Fig. 1: Overview of BlenderAlchemy.** Given an input Blender state and a user intention specified using either language or reference images, BlenderAlchemy edits the Blender state to satisfy that intention by *iteratively refining a Blender python program* that executes to produce the final Blender state. Our system additionally leverages text-to-image generation for *visual imagination*, a step that expands a text-only user intention to a concrete visual target to improve program refinement.

design, texture painting, animation, lighting, and scene composition. Prior attempts usually focus on specific editing tasks, like material synthesis [11, 65]. While these approaches show promising performance in the tasks they are designed for, it is non-trivial to combine these task-specific methods to satisfy an intended design goal. An alternative is to leverage Large Language Models (LLMs) [25, 30, 48] to digest user intent and suggest design actions by proposing intelligent combination of existing task-specific tools [39] or predicting editing programs step-by-step [52, 60]. While LLMs have excellent abilities to understand user intentions and suggest sequences of actions to satisfy them, applying LLMs to graphical design remains challenging largely because language cannot capture the visual consequences of actions performed in software like Blender.

One promising alternative is to leverage vision language models (VLM), such as LLaVA [26], GPT-4V [31], Gemini [47], and DallE-3 [10]. These VLMs have shown to be highly capable of understanding detailed visual information [16, 59, 61, 62] and generating compelling images [10]. We posit that these VLMs can be leveraged to complete different kinds of design tasks within the Blender design environment, like material editing and the design of lighting setups.

In this paper, we provide a proof-of-concept system using the vision foundation model GPT-4V to generate and edit programs that modify the state of a Blender workspace to satisfy a user intention. Specifically, we first initialize the state of a workspace within Blender. The Blender state is parameterized as a short Python program, and a base Blender file. The user will then input a text description and potentially a reference image to communicate the desirable design outcome. The system is tasked to edit the program so that, when executed on the base Blender file, the rendered image can satisfy the user's intention. Figure 1 provides an illustration of the problem setup.

**Fig. 2: Iterative visual program editing** employs a edit generator $G$ and a state evaluator $V$ in each iteration to explore and prune different potential program edits, where $G$ generates plausible variants of an input program and $V$ picks between two programs based on the consequences they have to the Blender visual state and their alignment to the user intention. Each iteration of the refinement explores variations of the most promising program from the previous iteration. See Algorithm 1 for details.

Naively applying VLMs to this editing setting gives rise to many failure cases, possibly due to the fact that out-of-the-box VLMs have a poor understanding of the visual consequences of Blender program edits. To counter this, we propose a visually-guided program search procedure that combines a vision-aware edit generator and a visual state evaluator to iteratively search for a suitable program edit (Figure 2). Inspired by the human design process, our system performs guided trial-and-error, capped by some computation budget. Within each iteration, the visual program generator will propose several possible edits on the current program. These edits will be applied and executed in Blender to produce a rendered image. These rendered images will be provided to visual evaluator, which will select the best renders via pairwise comparison by assessing which choice better satisfies the design goal specified by text and reference image. The program achieving the best render will replace the current program as a starting point for the visual program generator. This iterative search process, however, has very low success rate because of the sparsity of correct program edits in the vast program space. To improve the success rate, we further propose two techniques. First, when the proposed program of the new iteration does not contain a viable candidate, we revert to the best candidate of the prior program. This reversion mechanism make sure the search procedure will not diverge when facing a batch of bad edit candidates. Second, to facilitate our visual evaluator and generator to better understand user intent, we leverage the "visual imagination" of text-to-image generative model to imagine a reference image. We show that our method is capable of accomplishing graphical design tasks within Blender, guided by user intention in the form of text and images. We demonstrate the effectiveness of our system on material and lighting design tasks, both parts of the 3D design process where artists spend a significant amount of time, ranging from 20 hrs to 4-6 workdays *per* model [3, 4]. We show that our method can outperform prior works designed for similar problem settings, such as BlenderGPT [2]. In summary, the main contributions of this paper include:

1. We propose BlenderAlchemy, a system that's able to edit visual programs based on user input in the form of text or images.

2. We identify key components that make the system work: a visual state evaluator, a visual edit generator, a search algorithm with an edit reversion mechanism, and a visual imagination module to facilitate the search.
3. We provide evidence showing that our system can outperform prior works in text-based procedural material editing, as well as its applicability to other design tasks including editing lighting configurations.

## 2  Related works

***Task-specific tools for Material Design.*** Large bodies of works have been dedicated to using learning-based approaches to generate materials. Prior works exploit 2D diffusion models to generate texture maps either in a zero-shot way [11, 24, 36, 53, 57] or through fine-tuning [40, 63]. Though these works open the possibility of generating and modifying textures of objects using natural language for 3D meshes, these works do not model the material properties in a way that allows such objects to be relit. Other methods directly predict the physical properties of the surface of a material through learning, using diffusion [49, 50] or learning good latent representations from data [13, 20, 23, 28, 42, 65]. However, for all of the work mentioned so far, their fundamentally image-based or latent-based representations make the output materials difficult to edit in existing 3D graphics pipelines. People have also explored combining learning-based approaches with symbolic representations of materials [37, 46]. These works often involve creating differentiable representations of the procedural material graph often used in 3D graphical design pipelines, and backpropagating gradients throughout the graph to produce an image that can match the target [19, 22, 41]. Other works like Infinigen [35] use rule-based procedural generation over a library of procedural materials. However, no prior works in this direction have demonstrated the ability to edit procedural material graphs using *user intention specified by language* [37], a task that we are particularly interested in. Though the aforementioned works excel at material design, they aren't generalizable to other 3D graphics design task settings. BlenderAlchemy, on the other hand, aims to produce a system that can perform various design tasks according to user intents. This usually requires combining different methods together in a non-trivial way.

***LLM as general problem solvers.*** Large language models (LLMs) like GPT-4 [30,31], Llama [48], and Mistral [25] have in recent years demonstrated unprecedented results in a variety of problem settings, like robotics [7,15,21,27,44,55,64], program synthesis [29, 38, 43], and graphic design [24, 53, 58]. Other works have shown that by extending such models with an external process like Chain-of-Thought [52], Tree-of-Thought [60] or memory/skill database [14, 32, 33, 51], or by embedding such systems within environments where it can perceive and act [14, 32, 39, 51], a range of new problems that require iterative refinement can be solved. Their application to visual problem settings, however, has mostly been limited due to the nonexistent visual perception capabilities of the LLMs [14, 18, 51, 53, 56, 58]. While this could be sidestepped by fully condensing the visual state of the environment using text [32] or some symbolic representation [14, 51], doing so for 3D graphic design works poorly. For instance, the task

of editing a Blender material graph to create a desired material requires many trial-and-error cycles and an accurate understanding of the consequences certain design actions can have on the visual output. Recent works that apply LLMs to graphical design settings [14,18] and ones that more specifically do so within Blender (like BlenderGPT [2], 3D-GPT [45] and L3GO [56]) do not use visual information to inform or refine their system's outputs, leading to unsatisfactory results. BlenderAlchemy borrows ideas from existing LLM literature and tries to address this issue by inputting visual perception into the system.

***Vision-Language Models.*** State-of-the-art vision language models, such as LLaVA [26], GPT-4V [31], and Gemini [47] have demonstrated impressive understanding of the visual world and its connections to language and semantics, enabling many computer vision tasks like scene understanding, visual question answering, and object detection to be one API call away [16,17,59,61,62]. Works such as [9,54] suggest that such models can also be used as a replacement for human evaluators for a lot of tasks, positioning them as tools for guiding planning and search by acting as flexible reward functions. BlenderAlchemy takes the first steps to apply VLMs to solve 3D graphic design tasks, a novel yet challenging application rather unexplored by existing works.

## 3    Method

The goal of our system is to perform edits within the Blender 3D graphic design environment through iteratively refining programs that define a sequence of edits in Blender. This requires us to (1) decompose the input initial Blender input into a combination of programs and a "base" Blender state (Section 3.1) and (2) develop a procedure to edit each program to produce the desired visual state within Blender to match a user intention (Sections 3.2).

### 3.1    Representation of the Blender Visual State

The state of the initial Blender design environment can be decomposed into an "base" Blender state $S_{\text{base}}$ and a set of initial programs $\{p_0^{(1)}, p_0^{(2)}, ..., p_0^{(k)}\}$ that acts on state $S_{\text{base}}$ to produce the initial Blender environment through a dynamics function $F$ that transitions from one state to another based on a set of programmatic actions:

$$S_{\text{init}} = F\left(\left\{p_0^{(i)}\right\}_{i=1...k}, S_{\text{base}}\right)$$

In our problem setting, $F$ is the python code executor within the Blender environment that executes $\{p_0^{(i)}\}_{i=1...k}$ in sequence. We set each initial program $p_0^{(i)}$ to be a program that concerns a single part of the 3D graphical design workflow – for instance, $p_0^{(1)}$ is in charge of the material on one mesh within the scene, and $p_0^{(2)}$ is in charge of the lighting setup of the entire scene. The decomposition of $S_{\text{init}}$ into $S_{\text{base}}$ and $p_0^1...p_0^k$ can be done using techniques like the

"node transpiler" from Infinigen [35], which converts entities within the Blender instance into lines of Python code that can recreate a node graph, like a material shader graph. We develop a suite of tools to do this in our own problem setting.

Though it's possible for *all* edits to be encompassed in a *single* program instead of $k$ programs, this is limiting in practice – either because the VLM's output length isn't large enough for the code necessary or because the VLM has a low success rate, due to the program search space exploding in size. Although it is possible that future VLMs will substantially mitigate this problem, splitting the program into $k$ task-specific programs may still be desirable, given the possibility of querying $k$ task-specific fine-tuned/expert VLMs in parallel.

### 3.2 Iterative Refinement of Individual Visual Programs

Suppose that to complete a task like material editing for a single object, it suffices to decompose the initial state into $S_{\text{base}}$ and a single script, $p_0$ – that is, $S_{\text{init}} = F(\{p_0\}, S_{\text{base}})$. Then our goal is to discover some edited version of $p_0$, called $p_1$, such that $F(\{p_1\}, S_{\text{base}})$ produces a visual state better aligned with some user intention $I$. Our system assumes that the user intention $I$ is provided in the form of language and/or image references, leveraging the visual understanding of the latest VLM models to understand user intention.

---

**Algorithm 1** Iterative Refinement of Visual Programs

---

1: **procedure** Tournament(State candidates $\{S_1, S_2, ...S_k\}$, Visual state evaluator $V$, User intention $I$)
2:     **if** $k > 2$ **then**
3:         $w_1 \leftarrow$ Tournament($\{S_1, S_2, ...S_{k/2}\}$, $V$, $I$),
4:         $w_2 \leftarrow$ Tournament($\{S_{k/2}, S_{k+1}, ...S_k\}$, $V$, $I$)
5:     **else**
6:         $w_1 \leftarrow S_1, w_2 \leftarrow S_2$
7:     **end if**
8:     **return** $V(w_1, w_2, I)$
9: **end procedure**
10: **procedure** Refine(Depth $d$, Breadth $b$, Intention $I$, Edit Generator $G$, State Evaluator $V$, Base state $S_{\text{base}}$, Initial program $p_0$, Dynamics Function $F$)
11:     $S_0 \leftarrow F(p_0, S_{\text{base}})$, $S_{\text{best}} \leftarrow S_0$, $p_{\text{best}} \leftarrow p_0$
12:     **for** $i \leftarrow 1$ to $d$ **do**
13:         $\mathcal{P}_i = \mathcal{N}(p_{\text{best}})$ if $i$ is odd else $\mathcal{P}$          ▷ "Tweak" or "Leap" edits
14:         $p_i^1, p_i^2...p_i^b \leftarrow G(p_{\text{best}}, S_{\text{best}}, I, b, \mathcal{P}_i)$          ▷ Generate $b$ options
15:         $S_i^1 \leftarrow F(p_i^1, S_0), ..., S_i^b \leftarrow F(p_i^b, S_0)$          ▷ Observe the visual states
16:         $S_i^* \leftarrow$ Tournament($\{S_i^1, S_i^2...S_i^b\} \cup \{S_{\text{best}}\}$, $V$, $I$)          ▷ Choose the best
17:         $S_{\text{best}} \leftarrow S_i^*$, $p_{\text{best}} \leftarrow p_i^*$          ▷ Best visual state and programs so far
18:     **end for**
19:     **return** $S_{\text{best}}, p_{\text{best}}$
20: **end procedure**

---

To discover a good edit to $p_0$, we introduce the procedure outlined in Algorithm 1, an iterative refinement loop that repeatedly uses a visual state evaluator $V$ to select among the hypotheses from an edit generator $G$. A single "agent" for a certain task like procedural material design can be fully described by $(G, V)$.

Inspired by works like [54], we propose a visual state evaluator $V(S_1, S_2, I)$, which is tasked with returning whichever of the two visual states ($S_1$ or $S_2$) better matches the user intention $I$. This evaluator is applied recursively to choose the most suitable visual state candidate among $b$ visual state candidates by making $\mathcal{O}(\log(b))$ queries, as done in TOURNAMENT in Algorithm 1.

Though it seems straightforward to ask the same VLM to edit the code in a single pass, this leads to many failure cases (see more in Section 4.2). Due to the VLM's lack of baked-in understanding of the visual consequences of various programs within Blender, a multi-hypothesis and multi-step approach is more appropriate. Extending Tree-of-Thoughts [60] to the visual domain, $G(p, S, I, b, \mathcal{P})$ is a module tasked with generating $b$ different variations of program $p$, conditioned on the current visual state $S$ and user intention $I$, constrained such that the output programs fall within some family of programs $\mathcal{P}$, which can be used to instill useful priors to the edit generator.

Below we describe some additional system design decisions that ensured better alignment of the resultant edited program to the user intention, either by improving the stability of the procedure or by supplementing the visual understanding of VLMs. The effect of each is investigated in Section 4.2.

**Hypothesis Reversion.** To improve the stability of the edit discovery process, we add the visual state of the program being edited at every timestep ($S_{\text{best}}$ for $p_{\text{best}}$) as an additional candidate to the selection process, providing the option for the process to revert to an earlier version if the search at a single iteration was unsuccessful. Line 16 in Algorithm 1 shows this.

**Tweak and Leap Edits.** An important characteristic of visual programs is that continuous values hard-coded within the program can modify the output just as much as structural changes. This is in contrast to non-visual program synthesis tasks based on unit-tests of I/O specs, like [8, 12], where foundation models are mostly tasked to produce the right *structure* of the program with minimal hard-coded values. Given a a single visual program $p \in \mathcal{P}$, the space of visual outputs achievable through *only* tweaking the numerical values to function parameters and variable assignments can cover a wide range, depending on the fields available in the program. In Algorithm 1, we refer this space as the "neighborhood" of $p$ or "tweak" edits, $\mathcal{N}(p)$. Though this can result in a very small change to the program, this can lead to a large visual difference in the final output, and in a few edits can change the "wrong" program into one more aligned with the user intention. On the other hand, more drastic changes (or "leap" edits) may be needed to accomplish a task. Consider for example the task of changing a perfectly smooth material to have a noised level of roughness scattered sparsely across the material surface. This may require the programmatic addition of the

relevant nodes (e.g. color ramps or noise texture nodes), and thus the edited program falls outside of $\mathcal{N}(p)$.

Empirically, we find that the optimal edits are often a mix of tweak and leap edits. As such, our procedure cycles between restricting the edits of $G$ to two different sets: the neighborhood of $p_{\text{best}}$, and the whole program space $\mathcal{P}$ (Line 13 in Algorithm 1). In practice, such restrictions are softly enforced through in-context prompting of VLMs, and though their inputs encourage them to abide by these constraints, the model can still produce more drastic "tweak" edits or conservative "leap" edits as needed.

***Visual Imagination.*** In the case when the user intention is communicated purely textually, their intention may be difficult for the VLM to turn into successful edits. Prior works like [24, 53, 58] have made similar observations of abstract language for 3D scenes. Consider, for example, the prompt "make me a material that resembles a *celestial nebula*". To do this, the VLM must know what a celestial nebula looks like, *and* how it should change the parameters of the material shader nodes of, say, a wooden material. We find that in such cases, it's hard for the VLM to directly go from abstract descriptions to low-level program edits that affect low-level properties of the Blender visual state.

Instead, we propose supplementing the text-to-program understanding of VLM's with the text-to-image understanding in state-of-the-art image generation models. Intermediary visual artifacts (images generated using the user intention) are created and used to guide the refinement process towards a more plausible program edit to match the desired outcome, as shown in Figure 1. The generated images act as image references *in addition to the textual intention provided by the user*. This constitutes a simple visual chain of thought [52] for visual program editing, which not only creates a reference image for $G$ and $V$ to guide their low-level visual comparisons (e.g. color schemes, material roughness... *etc.*), but also provides a user-interpretable intermediary step to confirm the desired goals behind an otherwise vague user intention.

## 4    Experiments

We demonstrate BlenderAlchemy on editing procedural materials and lighting setups within Blender, two of the most tedious parts of the 3D design workflow.

### 4.1    Procedural Material Editing

Procedural material editing has characteristics that make it difficult for the same reason as a lot of other visual program settings: (1) small edit distances of programs may result in very large visual differences, contributing to potential instabilities in the edit discovery process, (2) it naturally requires trial-and-error when human users are editing them as the desired magnitude of edits depends on the visual output and (3) language descriptions of edit intent typically do not

contain low-level information for the editor to know immediately which part of the program to change.

We demonstrate this on two different kinds of editing tasks: (1) turning the same initial material (a synthetic wooden material) into other materials described by a list of **language descriptions** and (2) editing many different initial materials to resemble the same target material described by an **image input**. For the starting materials, we use the synthetic materials from Infinigen [35].
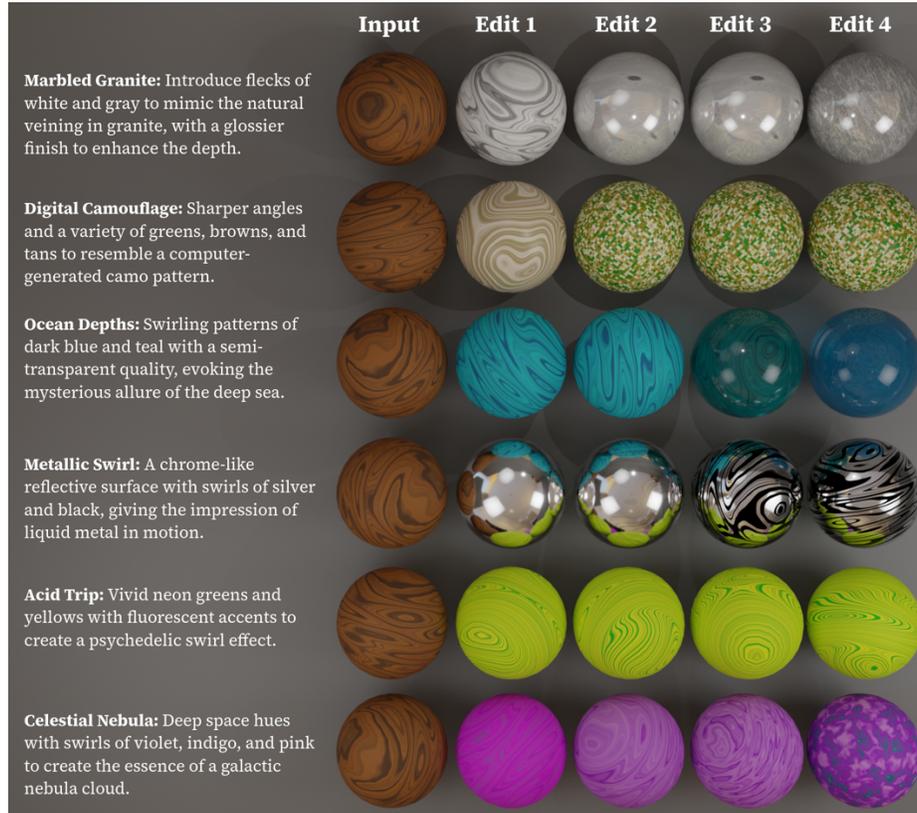
**Text-based material-editing** An attractive application of our system is in the modification of preexisting procedural materials using natural language descriptions that communicate user intent, a desireable but thus far undemonstrated capability of neurosymbolic methods [37]. We demonstrate our system's capabilities by asking it to edit a wooden material from Infinigen [35] into many other materials according to diverse language descriptions of target materials that are, importantly, *not wood*. In reality, this is a very challenging task, since this may require a wide range in the size of edits even if the language describing desired target material may be very similar.

Figure 3 shows examples of edits of the same starter wood material using different language descriptions, and Figure 4 demonstrates the intermediary steps of the problem solving process for a single instance of the problem. Our system is composed of an edit generator that generates 8 hypotheses per iteration, for 4 iterations ($d = 4, b = 8$), cycling between tweak and leap edits. It uses GPT-4V for edit generation and state evaluation, and DallE-3 for visual imagination.

We compare against BlenderGPT, the most recent open-sourced Blender AI agent that use GPT-4 to execute actions within the Blender environment through the Python API. [2] We provide the same target material text prompt to BlenderGPT, as well as the starter code for the initial wood material for reference. We compare the CLIP similarity of their output material to the input text description against our system. BlenderGPT reasons only about how to edit the program using the input text description, doing so in a single pass without state evaluation or multi-hypothesis edit generation. To match the number of edit generator queries we make, we run their method a maximum of 32 times, using the first successful example as its final output. Everything else remains the same, including the starter material program, text description, base Blender state, and lighting setup.

We find that qualitatively, BlenderGPT produces much shallower and more simplistic edits of the input material, resulting in low-quality output materials and poor alignment with the user intention. Examples can be seen in Figure 5. For instance, observe that for the "digital camouflage" example, BlenderAlchemy is able to produce the "sharper angles" that the original description requests (See Figure 3) whereas BlenderGPT produces the right colors but fails to create the sharp, digital look. For "metallic swirl" example, our system's visual state selection process would weed out insufficiently swirly examples such as the one

---

[2] Concurrent works like 3DGPT [45], L3GO [56] have not yet open-sourced their code.
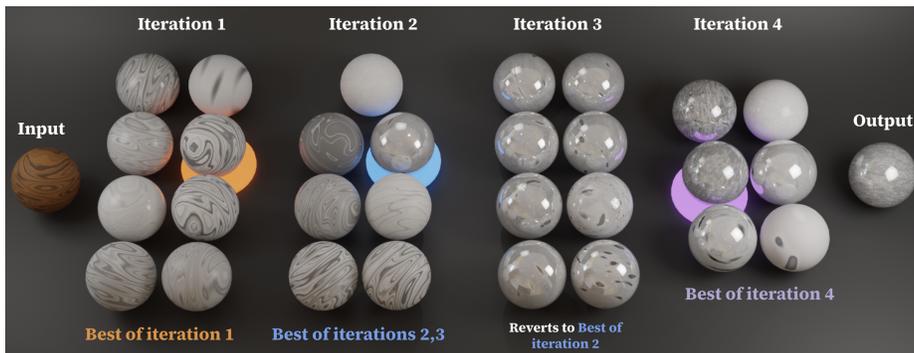
**Fig. 3: Text-based Material Editing Results.** The step-by-step edits of a 4x8 version of BlenderAlchemy to the same wooden material, given the text description on the left as the input user intention.

given by BlenderGPT, enabling our method to produce a material closer to the prompt.

Table 1 demonstrates the comparison in terms of the average ViT-B/32 and ViT-L/14 CLIP similarity scores [34] with respect to the language description. Our system's ability to iteratively refine the edits based on multiple guesses at each step gives it the ability to make more substantive edits over the course of the process. Moreover, visual grounding provided both by the visual state evaluator as well as the output of the visual imagination stage guides the program search procedure to better align with the textual description.

**Image-based material-editing** Given an image of a desired material, the task is to convert the code of the starter material into a material that contains many of the visual attributes of the input image, akin to doing a kind of style transfer for procedural materials.

At each step, the edit generator is first asked to textually enumerate a list of obvious visual differences between the current material and the target, then asked

**Fig. 4: The edit discovery process of turning a wooden material into "marbled granite".** Each column shows the hypotheses generated by $G$, with the most promising candidates chosen by $V$ indicated by the highlights. Note that iteration 3 proved to be unfruitful according to $V$, and the method reverts to the best candidate from iteration 2, before moving onto iteration 4.
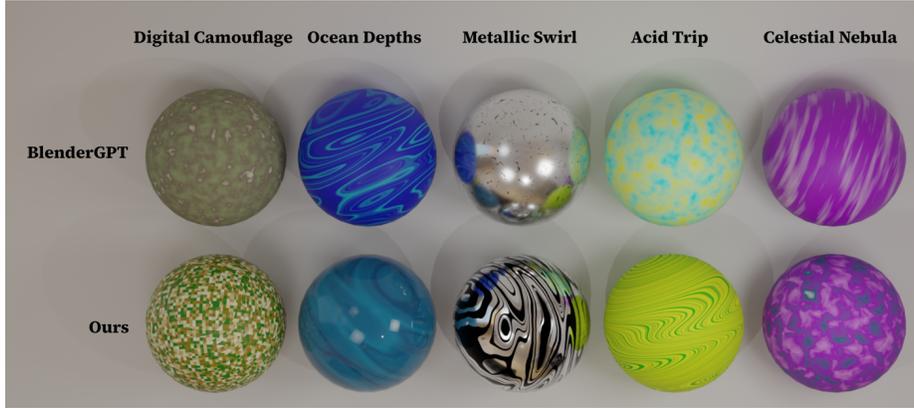
**Table 1: CLIP scores of BlenderAlchemy vs. BlenderGPT** for the text-based material editing task. We find that a version of our system that has *no visual components* (**-Vision**) still outperforms BlenderGPT. By adding vision to the state evaluator alone (**-Vision G**) and not the edit generator, a further improvement is observed.

| Metric | BlenderGPT | -Vision | -Vision G | Ours |
|---|---|---|---|---|
| ViT-B/32 ($\uparrow$) | 25.2 | 25.7 | 27.8 | **28.2** |
| ViT-L/14 ($\uparrow$) | 21.1 | 21.8 | 23.4 | **24.0** |

to locate lines within the code that may be responsible for these visual differences (*e.g.* "the target material looks more rough" $\rightarrow$ "line 23 sets a roughness value, which we should try to increase") before finally suggested an edited version of the program. As such, low-level visual differences (*e.g.* color discrepancies) are semantically compressed first into language descriptions (*e.g.* "the target is more red"), before being fed into the editing process, resulting in the behavior that our system produces variations of the input material that resembles the target image along many attributes, even if its outputs don't perfectly match the target image (See Figure 6). Our system is the same as for text-based material editing, but without the need for visual imagination.

## 4.2 Ablation Experiments on Material Editing

**Value of visually-grounded edit generators and state evaluators** How important is it that the edit generator $G$ and state evaluator $V$ have access to vision? We ablate (1) visual perception of $G$ by prompting it to *only* propose edits based on the text of the initial program $p_0$, without *any visual information* (neither the initial rendering $S_0$ or, if applicable, images of intentions) and (2) additionally the visual perception of $V$, by letting it decide which of two programs

**Fig. 5: Comparisons between our method and BlenderGPT for the text-based material editing task setting.** Note how our materials align better with the original language prompts (See Figure 3 for the original prompts). The input material being edited is the same wooden material.



**Fig. 6: Material editing based on image inputs.** Our edit intention is described by the target image shown on the right. 5 different Infinigen [35] initial materials are shown here, and the final edit. Note how in each case, certain attributes of the target material (metallicness, color, texture) are transferred.

is better purely on the basis of the program code and the textual description of the user intention $I$. Table 3 shows the consequences of such ablations to the alignment to the user intention.

**Value of the multi-hypothesis, multi-step approach** We vary the dimensions of our system ($d$ and $b$ within Algorithm 1) to demonstrate the importance of having a well-balanced dimensions. We keep the total number of calls to the edit generator constant ($d \times b$) at 32 requests. Table 2 shows the result of this on the image-based material editing task. This suggests that there's a sweet-spot in the system dimensions ($8 \times 4$ or $4 \times 8$) where the system is neither too wide (large $b$, in which case the system tends to overly explore the program space, and produce insubstantial edits) or too tall (large $d$, and in which case the system has high risk of overly exploiting suboptimal edit candidates). By the same argument, even if BlenderGPT [2] was equipped with visual perception and a

**Table 2: Different dimensions of BlenderAlchemy and their metrics on the image-based material editing task**. "$1 \times 32$" indicates a setup that uses $d = 1$ (number of iterations) and $b = 32$ (number of hypotheses per iteration) in Algorithm 1. This shows the clear advantage of using a more balanced choice of $d$ and $b$ over sequential iterative refinement method ($32 \times 1$) or querying the language model multiple times without refinement ($1 \times 32$).

|  | $1 \times 32$ | $2 \times 16$ | $4 \times 8$ | $8 \times 4$ | $16 \times 2$ | $32 \times 1$ |
|---|---|---|---|---|---|---|
| ViT-B/32 ($\uparrow$) | 81.7 | 82.9 | 81.7 | **84.1** | 83.6 | 81.6 |
| Photometric ($\downarrow$) | 0.066 | 0.066 | **0.049** | 0.050 | 0.056 | 0.087 |
| LPIPS ($\downarrow$) | 0.64 | 0.54 | 0.52 | **0.50** | 0.52 | 0.59 |

**Table 3: Ablating system design decisions.** For the text-based material editing task, we compare against variants in which we remove (1) visual perception from $G$ *and* $V$ (**-Vision**), (2) visual perception from $G$ and *not* from $V$ (**-Vision G**), (3) visual imagination (**-Imagin.**), (4) reversion capabilities (**-Revert**), (5) the option of leap edits (**-Leap**) or (6) tweak edits (**-Tweak**). We use a $4 \times 8$ version of BlenderAlchemy. Edits 1 to 4 correspond to the output at each refinement step of the BlenderAlchemy process. We show the ViT-B/32 CLIP scores here.

|  | -Vision | -Vision G | -Imagin. | -Revert | -Leap | -Tweak | Ours |
|---|---|---|---|---|---|---|---|
| Edit 1 | 27.4 | 27.6 | 26.8 | 27.1 | 27.1 | 27.2 | **27.8** |
| Edit 2 | 26.1 | 27.6 | 27.1 | 26.4 | 27.6 | 27 | **27.9** |
| Edit 3 | 26.5 | 27.6 | 26.8 | 26.6 | 27.7 | 26.9 | **28.4** |
| Edit 4 | 25.7 | 27.8 | 26.9 | 25.8 | 27.8 | 26.6 | **28.2** |

state evaluator to choose among 32 candidates, it would still suffer from the same issues as the $1 \times 32$ version of our method.

**Hypothesis Reversion** The intended effect of hypothesis reversion is to ensure the stability of the procedure, especially when (1) the tree's depth is sufficient for the edit search to go astray and (2) when leap edits can cause large and potentially disruptive edits in a single iteration of our procedure. As seen in Table 3, removing the ability to revert hypotheses causes divergence of the alignment to the text description over several edits, and larger drops in average CLIP-similarity corresponds to when leap edits happen (Edit 2 and Edit 4).

**Tweak/Leap Edits** When we ablate all "tweak" edits ("-Tweak" in Table 3) by making all edits "leap" edits, we observe a strong divergence from the user intention. Conversely, ablating all "leap" edits ("-Leap" in Table 3) leads to slow but steady increase in alignment with the user intention, but too conservative to match the "tweak+leap" variant ("Ours" in Table 3).

**Visual Imagination** Visual imagination is an additional image-generation step before launching the procedure in Algorithm 1, with the intended effect of guiding the edit generator and state evaluator with text-to-image understanding of

**"Hand lotion under disco lights, like at a nightclub."**



Fig. 7: **Optimizing the lighting** of the scene setup to the text-based user intention. We base the initial Blender state input based on a product visualization made by a Blender artist *Nam Nguyen*, downloaded from BlenderKit.

**"Hand cream sitting on a hot comet outside in the night"**



Fig. 8: **Optimizing the lighting and material iteratively** within a product visualization scene, to satisfy the text-based user intention. Note the dimming of the environment lighting for the nighttime lighting, and the glowing-hot material the editing procedure has produced. We base the initial Blender state input based on a product visualization made by an artist *blendervacations Kushwaha*, downloaded from BlenderKit.

state of the art image diffusion models. Without it, user intentions communicated using abstract language descriptions lead to poorer edits due to having limited information to properly guide the low-level visual comparisons (e.g. color, textures, ...*etc.*) by the state evaluator and edit generator (See Table 3).

## 5   Lighting

We show that BlenderAlchemy can be used to adjust the lighting of scenes according to language instructions as well. Figure 7 shows this being done for an input product visualization designed by a Blender artist.

As mentioned in Section 3.1, we can consider iteratively optimizing two separate programs, one controlling the lighting of the whole scene and another controlling the material of an object within the scene. That is, $S_{\text{init}} = F(\{p_0^{(L)}, p_0^{(M)}\}, S_{\text{base}})$ for initial lighting program $p_0^{(L)}$ and material program $p_0^{(M)}$. Figure 8 shows an example of this, where two separate pairs of edit gener-

ators and state evaluators, $(G_M, V_M)$ and $(G_L, V_L)$, are used to achieve an edit to the scene that aligns with the user intention, using Algorithm 2, with $k = 2$.

---

**Algorithm 2** Optimization of many programs using Algorithm 1

---

1: **procedure** MULTISKILLREFINE(Iteration number $N$, Agent collection $\{(G_1, V_1), (G_2, V_2)...(G_k, V_k)\}$, Base state $S_{\text{base}}$ and Initial programs $\{p_0^{(1)}, p_0^{(2)}, ...p_0^{(k)}\}$, User intention $I$, Dynamics function $F$)

2:     $p_{best}^{(1)} \leftarrow p_0^{(1)}, ..., p_{best}^{(k)} \leftarrow p_0^{(k)}$          ▷ Initialize best program edits

3:     **for** $i \leftarrow 1$ to $N$ **do**:

4:         **for** $j \leftarrow 1$ to $k$ **do**:

5:             $F_j \leftarrow \lambda x, S_{\text{base}} : F(\{p_{\text{best}}^{(a)}\}_{a \neq j} \cup \{x\}, S_{base})$          ▷ partial function

6:             $S_{\text{best}}^j, p_{\text{best}}^j \leftarrow \text{REFINE}(d, b, I, G_j, V_j, S_{\text{base}}, p_{\text{best}}^{(j)}, F_j)$

7:         **end for**

8:     **end for**

9:     **return** $\{p_{\text{best}}^{(1)}, ...p_{\text{best}}^{(k)}\}$

10: **end procedure**

---

## 6    Conclusion & Discussion

In this paper, we introduce BlenderAlchemy, a system that performs edits within the Blender 3D design environment by leveraging vision-language models to iteratively refining a program to be more aligned with the user intention, by using visual information to both explore and prune possibilities within the program space. We equip our system with visual imagination by providing it access to text-to-image models, a tool it uses to guide itself towards program edits that better align with user intentions.

***Limitations and Future works.*** We've demonstrated this concept on lighting setup and material editing, and we believe that future work should seek to extend this to other graphical design workflows as well, such as character animation and 3D modeling. Furthermore, as our core focus has been demonstrating the plausibility of the program refinement procedure, we do not build or use a library of tools/skills [14,51,56]. BlenderAlchemy can be further strengthened if a strong workflow-specific library of skills is either provided [14,39,56] or learned [51].

## References

1. Blender github, https://github.com/blender/blender 20
2. Blendergpt, https://github.com/gd3kr/BlenderGPT 3, 5, 12, 20
3. How long does it take to create a 3d model?, https://3d-ace.com/blog/how-long-does-it-take-to-create-a-3d-model/ 3
4. How long does it take to make a 3d model?, https://pixune.com/blog/how-long-does-it-take-to-create-a-3d-model/ 3

5. Ollama, https://ollama.com/ 21
6. Reddit r/blender, https://www.reddit.com/r/blender/ 20
7. Ahn, M., Brohan, A., Brown, N., Chebotar, Y., Cortes, O., David, B., Finn, C., Fu, C., Gopalakrishnan, K., Hausman, K., et al.: Do as i can, not as i say: Grounding language in robotic affordances. arXiv preprint arXiv:2204.01691 (2022) 4
8. Austin, J., Odena, A., Nye, M., Bosma, M., Michalewski, H., Dohan, D., Jiang, E., Cai, C., Terry, M., Le, Q., et al.: Program synthesis with large language models. arXiv preprint arXiv:2108.07732 (2021) 7
9. Baumli, K., Baveja, S., Behbahani, F., Chan, H., Comanici, G., Flennerhag, S., Gazeau, M., Holsheimer, K., Horgan, D., Laskin, M., et al.: Vision-language models as a source of rewards. arXiv preprint arXiv:2312.09187 (2023) 5
10. Betker, J., Goh, G., Jing, L., Brooks, T., Wang, J., Li, L., Ouyang, L., Zhuang, J., Lee, J., Guo, Y., et al.: Improving image generation with better captions. Computer Science. https://cdn. openai. com/papers/dall-e-3. pdf **2**(3), 8 (2023) 2
11. Chen, D.Z., Siddiqui, Y., Lee, H.Y., Tulyakov, S., Nießner, M.: Text2tex: Text-driven texture synthesis via diffusion models. arXiv preprint arXiv:2303.11396 (2023) 2, 4
12. Chen, M., Tworek, J., Jun, H., Yuan, Q., de Oliveira Pinto, H.P., Kaplan, J., Edwards, H., Burda, Y., Joseph, N., Brockman, G., Ray, A., Puri, R., Krueger, G., Petrov, M., Khlaaf, H., Sastry, G., Mishkin, P., Chan, B., Gray, S., Ryder, N., Pavlov, M., Power, A., Kaiser, L., Bavarian, M., Winter, C., Tillet, P., Such, F.P., Cummings, D., Plappert, M., Chantzis, F., Barnes, E., Herbert-Voss, A., Guss, W.H., Nichol, A., Paino, A., Tezak, N., Tang, J., Babuschkin, I., Balaji, S., Jain, S., Saunders, W., Hesse, C., Carr, A.N., Leike, J., Achiam, J., Misra, V., Morikawa, E., Radford, A., Knight, M., Brundage, M., Murati, M., Mayer, K., Welinder, P., McGrew, B., Amodei, D., McCandlish, S., Sutskever, I., Zaremba, W.: Evaluating large language models trained on code (2021) 7
13. Chen, Y., Chen, R., Lei, J., Zhang, Y., Jia, K.: Tango: Text-driven photorealistic and robust 3d stylization via lighting decomposition. Advances in Neural Information Processing Systems **35**, 30923–30936 (2022) 4
14. De La Torre, F., Fang, C.M., Huang, H., Banburski-Fahey, A., Fernandez, J.A., Lanier, J.: Llmr: Real-time prompting of interactive worlds using large language models. arXiv preprint arXiv:2309.12276 (2023) 4, 5, 15, 20, 21
15. Firoozi, R., Tucker, J., Tian, S., Majumdar, A., Sun, J., Liu, W., Zhu, Y., Song, S., Kapoor, A., Hausman, K., et al.: Foundation models in robotics: Applications, challenges, and the future. arXiv preprint arXiv:2312.07843 (2023) 4
16. Fu, C., Chen, P., Shen, Y., Qin, Y., Zhang, M., Lin, X., Yang, J., Zheng, X., Li, K., Sun, X., Wu, Y., Ji, R.: Mme: A comprehensive evaluation benchmark for multimodal large language models. arXiv preprint arXiv:2306.13394 (2023) 2, 5
17. Fu, C., Zhang, R., Wang, Z., Huang, Y., Zhang, Z., Qiu, L., Ye, G., Shen, Y., Zhang, M., Chen, P., Zhao, S., Lin, S., Jiang, D., Yin, D., Gao, P., Li, K., Li, H., Sun, X.: A challenger to gpt-4v? early explorations of gemini in visual expertise. arXiv preprint arXiv:2312.12436 (2023) 5
18. Goel, P., Wang, K.C., Liu, C.K., Fatahalian, K.: Iterative motion editing with natural language. arXiv preprint arXiv:2312.11538 (2023) 4, 5
19. Guerrero, P., Hašan, M., Sunkavalli, K., Měch, R., Boubekeur, T., Mitra, N.J.: Matformer: A generative model for procedural materials. arXiv preprint arXiv:2207.01044 (2022) 4
20. Henzler, P., Deschaintre, V., Mitra, N.J., Ritschel, T.: Generative modelling of brdf textures from flash images. arXiv preprint arXiv:2102.11861 (2021) 4

21. Hu, Y., Xie, Q., Jain, V., Francis, J., Patrikar, J., Keetha, N., Kim, S., Xie, Y., Zhang, T., Zhao, Z., et al.: Toward general-purpose robots via foundation models: A survey and meta-analysis. arXiv preprint arXiv:2312.08782 (2023) 4

22. Hu, Y., Guerrero, P., Hasan, M., Rushmeier, H., Deschaintre, V.: Node graph optimization using differentiable proxies. In: ACM SIGGRAPH 2022 conference proceedings. pp. 1–9 (2022) 4

23. Hu, Y., He, C., Deschaintre, V., Dorsey, J., Rushmeier, H.: An inverse procedural modeling pipeline for svbrdf maps. ACM Transactions on Graphics (TOG) **41**(2), 1–17 (2022) 4

24. Huang, I., Krishna, V., Atekha, O., Guibas, L.: Aladdin: Zero-shot hallucination of stylized 3d assets from abstract scene descriptions. arXiv preprint arXiv:2306.06212 (2023) 4, 8, 29

25. Jiang, A.Q., Sablayrolles, A., Mensch, A., Bamford, C., Chaplot, D.S., Casas, D.d.l., Bressand, F., Lengyel, G., Lample, G., Saulnier, L., et al.: Mistral 7b. arXiv preprint arXiv:2310.06825 (2023) 2, 4

26. Li, C., Wong, C., Zhang, S., Usuyama, N., Liu, H., Yang, J., Naumann, T., Poon, H., Gao, J.: Llava-med: Training a large language-and-vision assistant for biomedicine in one day. Advances in Neural Information Processing Systems **36** (2024) 2, 5

27. Liang, J., Huang, W., Xia, F., Xu, P., Hausman, K., Ichter, B., Florence, P., Zeng, A.: Code as policies: Language model programs for embodied control. In: 2023 IEEE International Conference on Robotics and Automation (ICRA). pp. 9493–9500. IEEE (2023) 4

28. Liu, J., Gan, Y., Dong, J., Qi, L., Sun, X., Jian, M., Wang, L., Yu, H.: Perception-driven procedural texture generation from examples. Neurocomputing **291**, 21–34 (2018) 4

29. Olausson, T.X., Inala, J.P., Wang, C., Gao, J., Solar-Lezama, A.: Is self-repair a silver bullet for code generation? In: The Twelfth International Conference on Learning Representations (2023) 4

30. OpenAI: Gpt-4 system card. OpenAI (2023), https://cdn.openai.com/papers/gpt-4-system-card.pdf 2, 4

31. OpenAI: Gpt-4v(ision) system card. OpenAI (2023), https://api.semanticscholar.org/CorpusID:263218031 2, 4, 5

32. Park, J.S., O'Brien, J., Cai, C.J., Morris, M.R., Liang, P., Bernstein, M.S.: Generative agents: Interactive simulacra of human behavior. In: Proceedings of the 36th Annual ACM Symposium on User Interface Software and Technology. pp. 1–22 (2023) 4

33. Patil, S.G., Zhang, T., Wang, X., Gonzalez, J.E.: Gorilla: Large language model connected with massive apis. arXiv preprint arXiv:2305.15334 (2023) 4

34. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: International conference on machine learning. pp. 8748–8763. PMLR (2021) 10

35. Raistrick, A., Lipson, L., Ma, Z., Mei, L., Wang, M., Zuo, Y., Kayan, K., Wen, H., Han, B., Wang, Y., et al.: Infinite photorealistic worlds using procedural generation. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12630–12641 (2023) 4, 6, 9, 12, 30, 31

36. Richardson, E., Metzer, G., Alaluf, Y., Giryes, R., Cohen-Or, D.: Texture: Text-guided texturing of 3d shapes. arXiv preprint arXiv:2302.01721 (2023) 4

37. Ritchie, D., Guerrero, P., Jones, R.K., Mitra, N.J., Schulz, A., Willis, K.D., Wu, J.: Neurosymbolic models for computer graphics. In: Computer Graphics Forum. vol. 42, pp. 545–568. Wiley Online Library (2023) 4, 9

38. Romera-Paredes, B., Barekatain, M., Novikov, A., Balog, M., Kumar, M.P., Dupont, E., Ruiz, F.J., Ellenberg, J.S., Wang, P., Fawzi, O., et al.: Mathematical discoveries from program search with large language models. Nature **625**(7995), 468–475 (2024) 4

39. Schick, T., Dwivedi-Yu, J., Dessì, R., Raileanu, R., Lomeli, M., Hambro, E., Zettlemoyer, L., Cancedda, N., Scialom, T.: Toolformer: Language models can teach themselves to use tools. Advances in Neural Information Processing Systems **36** (2024) 2, 4, 15

40. Sharma, P., Jampani, V., Li, Y., Jia, X., Lagun, D., Durand, F., Freeman, W.T., Matthews, M.: Alchemist: Parametric control of material properties with diffusion models. arXiv preprint arXiv:2312.02970 (2023) 4

41. Shi, L., Li, B., Hašan, M., Sunkavalli, K., Boubekeur, T., Mech, R., Matusik, W.: Match: Differentiable material graphs for procedural material capture. ACM Transactions on Graphics (TOG) **39**(6), 1–15 (2020) 4

42. Shimizu, E., Fisher, M., Paris, S., McCann, J., Fatahalian, K.: Design adjectives: A framework for interactive model-guided exploration of parameterized design spaces. In: Proceedings of the 33rd Annual ACM Symposium on User Interface Software and Technology. pp. 261–278 (2020) 4

43. Shinn, N., Cassano, F., Labash, B., Gopinath, A., Narasimhan, K., Yao, S.: Reflexion: Language agents with verbal reinforcement learning.(2023). arXiv preprint cs.AI/2303.11366 (2023) 4

44. Singh, I., Blukis, V., Mousavian, A., Goyal, A., Xu, D., Tremblay, J., Fox, D., Thomason, J., Garg, A.: Progprompt: Generating situated robot task plans using large language models. In: 2023 IEEE International Conference on Robotics and Automation (ICRA). pp. 11523–11530. IEEE (2023) 4

45. Sun, C., Han, J., Deng, W., Wang, X., Qin, Z., Gould, S.: 3d-gpt: Procedural 3d modeling with large language models. arXiv preprint arXiv:2310.12945 (2023) 5, 9, 20, 21

46. Tchapmi, L.P., Ray, T., Tchapmi, M., Shen, B., Martin-Martin, R., Savarese, S.: Generating procedural 3d materials from images using neural networks. In: 2022 4th International Conference on Image, Video and Signal Processing. pp. 32–40 (2022) 4

47. Team, G., Anil, R., Borgeaud, S., Wu, Y., Alayrac, J.B., Yu, J., Soricut, R., Schalkwyk, J., Dai, A.M., Hauth, A., et al.: Gemini: a family of highly capable multimodal models. arXiv preprint arXiv:2312.11805 (2023) 2, 5

48. Touvron, H., Lavril, T., Izacard, G., Martinet, X., Lachaux, M.A., Lacroix, T., Rozière, B., Goyal, N., Hambro, E., Azhar, F., et al.: Llama: Open and efficient foundation language models. arXiv preprint arXiv:2302.13971 (2023) 2, 4

49. Vecchio, G., Martin, R., Roullier, A., Kaiser, A., Rouffet, R., Deschaintre, V., Boubekeur, T.: Controlmat: A controlled generative approach to material capture. arXiv preprint arXiv:2309.01700 (2023) 4

50. Vecchio, G., Sortino, R., Palazzo, S., Spampinato, C.: Matfuse: Controllable material generation with diffusion models. arXiv preprint arXiv:2308.11408 (2023) 4

51. Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., Anandkumar, A.: Voyager: An open-ended embodied agent with large language models. arXiv preprint arXiv:2305.16291 (2023) 4, 15, 21

52. Wei, J., Wang, X., Schuurmans, D., Bosma, M., Xia, F., Chi, E., Le, Q.V., Zhou, D., et al.: Chain-of-thought prompting elicits reasoning in large language models. Advances in Neural Information Processing Systems **35**, 24824–24837 (2022) 2, 4, 8

53. Wen, Z., Liu, Z., Sridhar, S., Fu, R.: Anyhome: Open-vocabulary generation of structured and textured 3d homes. arXiv preprint arXiv:2312.06644 (2023) 4, 8

54. Wu, T., Yang, G., Li, Z., Zhang, K., Liu, Z., Guibas, L., Lin, D., Wetzstein, G.: Gpt-4v (ision) is a human-aligned evaluator for text-to-3d generation. arXiv preprint arXiv:2401.04092 (2024) 5, 7

55. Xiao, X., Liu, J., Wang, Z., Zhou, Y., Qi, Y., Cheng, Q., He, B., Jiang, S.: Robot learning in the era of foundation models: A survey. arXiv preprint arXiv:2311.14379 (2023) 4

56. Yamada, Y., Chandu, K., Lin, Y., Hessel, J., Yildirim, I., Choi, Y.: L3go: Language agents with chain-of-3d-thoughts for generating unconventional objects. arXiv preprint arXiv:2402.09052 (2024) 4, 5, 9, 15, 20, 21

57. Yang, H., Chen, Y., Pan, Y., Yao, T., Chen, Z., Mei, T.: 3dstyle-diffusion: Pursuing fine-grained text-driven 3d stylization with 2d diffusion models. In: Proceedings of the 31st ACM International Conference on Multimedia. pp. 6860–6868 (2023) 4

58. Yang, Y., Sun, F.Y., Weihs, L., VanderBilt, E., Herrasti, A., Han, W., Wu, J., Haber, N., Krishna, R., Liu, L., et al.: Holodeck: Language guided generation of 3d embodied ai environments. In: The IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2024). vol. 30, pp. 20–25. IEEE/CVF (2024) 4, 8

59. Yang, Z., Li, L., Lin, K., Wang, J., Lin, C.C., Liu, Z., Wang, L.: The dawn of lmms: Preliminary explorations with gpt-4v (ision). arXiv preprint arXiv:2309.17421 **9**(1), 1 (2023) 2, 5

60. Yao, S., Yu, D., Zhao, J., Shafran, I., Griffiths, T., Cao, Y., Narasimhan, K.: Tree of thoughts: Deliberate problem solving with large language models. Advances in Neural Information Processing Systems **36** (2024) 2, 4, 7

61. Yin, S., Fu, C., Zhao, S., Li, K., Sun, X., Xu, T., Chen, E.: A survey on multimodal large language models. arXiv preprint arXiv:2306.13549 (2023) 2, 5

62. Yin, S., Fu, C., Zhao, S., Xu, T., Wang, H., Sui, D., Shen, Y., Li, K., Sun, X., Chen, E.: Woodpecker: Hallucination correction for multimodal large language models. arXiv preprint arXiv:2310.16045 (2023) 2, 5

63. Zeng, X.: Paint3d: Paint anything 3d with lighting-less texture diffusion models. arXiv preprint arXiv:2312.13913 (2023) 4

64. Zhou, H., Yao, X., Meng, Y., Sun, S., BIng, Z., Huang, K., Knoll, A.: Language-conditioned learning for robotic manipulation: A survey. arXiv preprint arXiv:2312.10807 (2023) 4

65. Zsolnai-Fehér, K., Wonka, P., Wimmer, M.: Gaussian material synthesis. arXiv preprint arXiv:1804.08369 (2018) 2, 4

We discuss the societal impact and limitations of our work in Sections A and B. We then outline the prompts we use for the state evaluator and edit generator in Section C. We investigate the qualitative effects of ablating parts of our system in Section D. In Section E, we provide insights on the code edits that BlenderAlchemy produces, and the resultant changes to the material node graph. Finally, we showcase some renderings of scenes that feature BlenderAlchemy materials in Section F.

## A    Societal Impact

***Impact on the creative community.*** Blender is the most popular free and open sourced 3D creation suite for both professional and hobbyist 3D designers, and has developed a huge community of artists (**r/blender** has, as of March 2024, close to 1.2 million members [6]) and developers (the Blender github has 4.1K followers, and 134K commits to their official Blender repo [1]). Community efforts like Blender open movies ( https://studio.blender.org/films/) not only demonstrate the versatility and expressive power of the tool itself, but also the collective drive of human imagination. Lowering the barrier to entry for the regular user to participate in such community art projects is the main motivation of BlenderAlchemy. BlenderAlchemy *does not* (and *does not intend to*) replace the intention of the human creator, and is the reason why we've made the starting point of BlenderAlchemy *an existing project file*, and *a specification of intent*. This is in contrast to all the other works like [2, 14, 45, 56], which attempt to generate everything from scratch. The goal of this project is not about the creator *doing* less as much as it is about enabling the creator *to create* more.

***Systematic Biases.*** LLMs and VLMs have inherent biases that are inherited by BlenderAlchemy. For now, safeguarding against such biases from making it into the final 3D creation requires the careful eye of the human user in the process – the fact that edits are made directly in Blender means that an established UI already exist to correct / remove / reject problematic edits. The multi-hypothesis nature of the edit generator also allows multiple possibilities to be generated, increasing the chances that there is one that is both usable and unproblematic.

## B    Limitations

***Cost and speed of inference*** Our system uses state of the art vision-language models. We've demonstrated this with GPT-4V, which as of March of 2024, remains very expensive and high-latency. For each of the material examples shown in the paper synthesized by a $4 \times 8$ dimensional system, the cost stands at just under \$3 per material, the bulk of which is spent on the edit generator $G$. In practice, it's likely that BlenderAlchemy will need to synthesize many candidate materials for one to be usable for the end product, further increasing the average costs per *usable* material. Beyond optimizing our system further (e.g.

instead of pairwise comparisons for the visual state evaluator, ask it to choose the best candidate in batches), we expect that the cost and speed of VLMs will substantially improve in the near future. Efforts that develop applications that can run large pretrained models locally [5] also hold promise that we can further lower the latency/costs by running open-source pretrained models locally, concurrently with 3D design processes.

**Library of skills** We do not develop a library of skills that our procedure can use, which have shown to be important in [14, 45, 51]. Such libraries are likely to be extremely domain-specific (library tools used by a material editor would be very different than animation), and will be the subject of our future work.

**Edit-only** We've made *editing* 3D graphics our main objective, and though the method can, in theory, be trivially extended to iteratively edit an initial empty scene into the full generated scene, we do not demonstrate this. We believe that stronger VLMs than what exists today will be necessary to do end-to-end generation without humans in the loop, unless such generation is done in simple settings like [56] and [45].

**VLM not fine-tuned to Blender scripting** GPT-4V still hallucinates when producing Blender python code with imports of libraries that don't exist or assigning values of the wrong dimensions to fields. In such cases, our system does rejection sampling of edit generations based on whether errors are returned when run in Blender, but such a method can incur penalties in runtime and cost. Future work can look to fine-tune on real and synthetic datasets of Blender scripts and resultant visual outputs. We believe in such cases, the iterative refinement procedure of BlenderAlchemy may likely still be useful due to the inherent multi-step, trial-and-error nature of human design process in pursuit of a vague intent specification.

## C    Prompting the State Evaluator and the Edit Generator

For the material-editing task setting, the prompts used for our edit generator are shown in Figures 9 (for leap edits) and 10 (for tweak edits). The prompt used for our visual state evaluator is shown in Figure 11. For the ablations of visual perception and visual imagination in our main paper, we adapt each of these prompts. For ablation of visual imagination, we remove the target image from the $G$ prompt and the $V$ prompt, and use the target description instead. An example of this for $V$'s prompt can be seen in Figure 13. For ablation of visual perception of either (or both) $G$ and/or $V$, we modify the prompts for the state evaluator to judge between two candidates based on their code only (Figure 12) and do the equivalent with $G$, where we provide only the target description and the source code, and ask it to generate code directly. For the images of materials, we render a $512 \times 512$ image from the same camera in the Blender design space, on one side of a sphere onto which a material is applied.

The prompts that we use for editing lighting configurations are essentially the same as the ones shown for materials, with minor changes to swap mentions of materials to mentions of lighting setup.

```
The following Blender code was used to produce a material:
'''python
[INITIAL MATERIAL CODE]
'''

The final material is assigned to the object 'material_obj', a sphere
    , and produces the rendering on the left below:

[IMAGES OF CURRENT MATERIAL AND TARGET MATERIAL]

The desired material is shown in the image on the right. Please
    describe the difference between the two materials, and edit the
    code above to reflect this desired change. Pay special attention
     to the base color of the materials.
MAKE SURE YOUR CODE IS RUNNABLE. MAKE SURE TO ASSIGN THE FINAL
    MATERIAL TO 'material_obj' (through 'apply(material_obj)') AS
    THE LAST LINE OF YOUR CODE.
DO NOT BE BRIEF IN YOUR CODE. DO NOT ABBREVIATE YOUR CODE WITH "..."
    -- TYPE OUT EVERYTHING.
```

**Fig. 9:** Prompts used to generate **leap** edits on an input material.

## D    Explaining System Design Decisions Through Qualitative Examples

So far, we've investigated the quantitative effects of ablating various parts of our system on metrics measuring the alignment of a material with a user's intention. We now discuss the qualitative effects of these ablations on the output on the material editing task, with reference to a sample shown in Figure 14.

***Tweak and leap edits.*** The columns "- Tweak" and "- Leap" correspond to leap-only and tweak-only versions of BlenderAlchemy. When leap edits are disabled ("- leap"), we can see that the edits fail to change the structure of the swirls, but instead produce darker stripes in an attempt to make the output look more marble-like, a change that can be associated with continuous parameters of certain graph nodes. On the other hand, when tweak edits are disabled,

```
The following Blender code was used to produce a material:
```python
[INITIAL MATERIAL CODE]
```


The final material is assigned to the object 'material_obj', a sphere
    , and produces the rendering on the left below:

[IMAGES OF CURRENT MATERIAL AND TARGET MATERIAL]


The desired material is shown in the image on the right.
Answer the following questions:
1) What is the SINGLE most visually obvious difference between the
    two materials in the image above?
2) Look at the code. Which fields/variables which are set to
    numerical values are most likely responsible for the obvious
    visual difference in your answer to question 1?
3) Copy the code above (COPY ALL OF IT) and replace the assignments
    of such fields/variables accordingly!
MAKE SURE YOUR CODE IS RUNNABLE. MAKE SURE TO ASSIGN THE FINAL
    MATERIAL TO 'material_obj' (through 'apply(material_obj)') AS
    THE LAST LINE OF YOUR CODE.
DO NOT BE BRIEF IN YOUR CODE. DO NOT ABBREVIATE YOUR CODE WITH "..."
    -- TYPE OUT EVERYTHING.
```

**Fig. 10:** Prompts used to generate **tweak** edits on an input material.

```
Here is the target material rendering:
[IMAGE OF TARGET MATERIAL]

Below, I show two different materials. Which one is visually more
    similar to the target material rendering? The one on the left or
     right?

[CONCATENATED IMAGES OF 2 CANDIDATES]
```
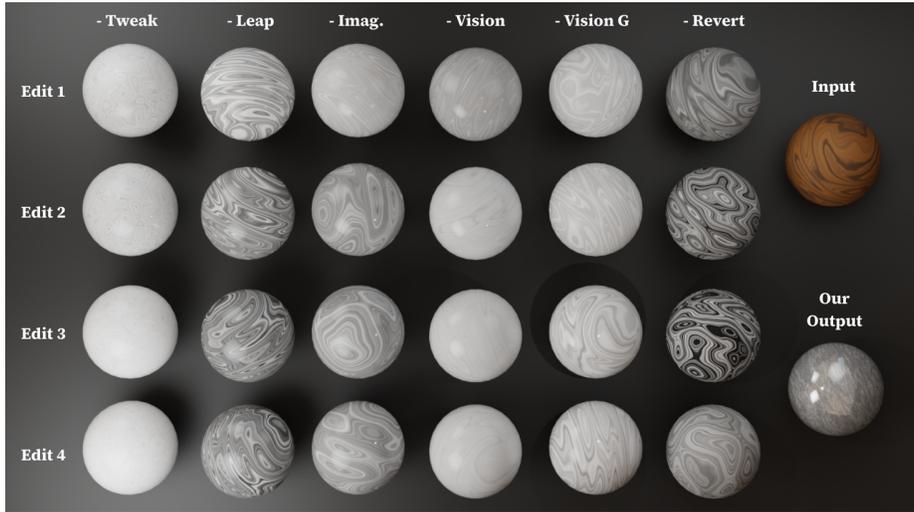
**Fig. 11:** Prompts used to evaluate visual state of the edited material.

```
Our desired target material can be described by:
    [TARGET DESCRIPTION] .
Imagine I'm showing you two Blender python scripts for materials, and
    they're side by side. Which one has the highest chance of
    producing the desired target material in Blender? The one on the
    left or right?
Code on the LEFT:
'''python
 [CODE OF CANDIDATE 1]
'''
Code on the RIGHT:
'''python
 [CODE OF CANDIDATE 2]
'''
Make sure that your final answer indicates which one has the highest
    chance of producing the desired material -- left or right.
    Answer by putting left or right in '''s.
```

**Fig. 12:** Prompts used to evaluate visual state of the edited material, **when used without vision**.

```
Our desired target material can be described by:
    [TARGET DESCRIPTION] .
Below, I show two different materials. Which one is visually more
    similar to the desired material described? The one on the left
    or right?
 [CONCATENATED IMAGES OF 2 CANDIDATES]
```

**Fig. 13:** Prompts used to evaluate visual state of the edited material, **when used without visual imagination for the target material**.

**Fig. 14: Qualitative samples of outputs at every editing step of a $4\times8$ version of BlenderAlchemy, across different ablations.** For reference, we show the input and the output of our unablated system on the right.

BlenderAlchemy produces drastic changes to the swirling patterns of the wood in Edit 1, but leads to a plateauing of progress, as all subsequent edits are drastic enough that reversion reverts back to edit 1, making no progress beyond the pale white material in edit 1.

***Visual imagination.*** The column "- Imag." corresponds to ablations of visual imagination for the text-based material editing task. Note how though edits are being made in every edit iteration, the verisimilitude of the material plateaus very quickly. Without a visual target to compare against, the edit generator has a difficult time knowing how to adjust the parameters of the shader node graph.

***Visual perception*** Columns "-Vision" and "- Vision G" correspond to ablating (1) the vision of the edit generator *and* the visual state evaluator and (2) ablating the vision of the edit generator *only*. In both cases, we see that it's mostly adjusting the color of the input wood material, making it light grey to match the prompt. However, the end result does not look like marbled granite.

***Edit hypothesis reversion.*** Column "- Revert" shows what happens when edit hypothesis reversion is disabled. As can be seen in Edit 3, the best candidate among the edit hypotheses is chosen to be a material that is *less* similar to granite marble than Edit 1. Edit 4 recuperates a little, but the instability has costed BlenderAlchemy 2 edit cycles, all just to eventually end up with a material that fits the prompt as much as Edit 1. This shows the importance of providing BlenderAlchemy the ability to revert to earlier edit hypotheses.
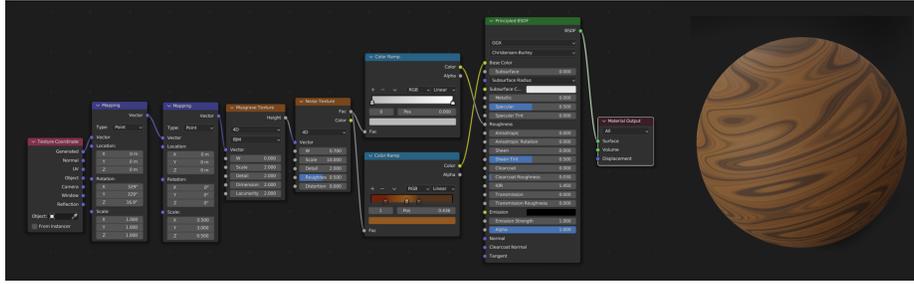
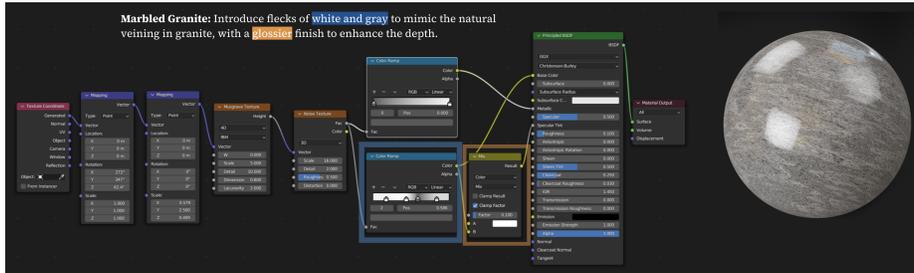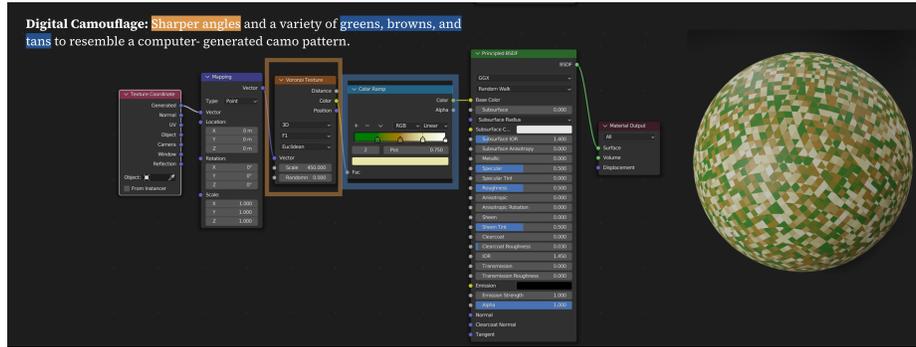**Fig. 15: Blender material graph of the starter wood material.**



**Fig. 16: Blender material graph of the "marbled granite" material.** Note the correspondence between "white and gray" with the colors chosen for the color ramp, and "glossy finish" with the input into the specularity port of the principled BSDF node.
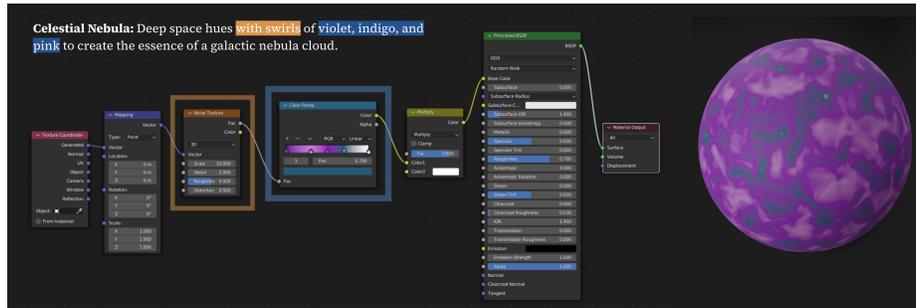
## E  Analysis of Program Edits

What kinds of qualitative changes to programs and material ndoe graphs does BlenderAlchemy affect? Figure 15 shows the wood material shader node graph corresponding to the input material to the text-based material editing task setting. As a few samples, Figures 16, 18, 17 show the material shader node graph of 3 materials edited from the input material.

Starting with the wood shader graph, one can observe changes of graph node types (e.g. insertion of the Voronoi node in camo material) or connectivity (e.g. changing the ports of the Principled BSDF node that receive inputs from the color ramps in the marbled granite material), continuous values (e.g. color values of the color ramps of the celestial nebula material) or categorical values in nodes (e.g. changing 4D to 3D noise textures in the marbled granite material).

We can observe this in the code as well. An example an be seen in the code edits for the digital camouflage material, in Figure 19. Figure 19a shows the input Blender script of the wooden material (whose graph and rendering are shown in Figure 15). BlenderAlchemy eventually edits the material by outputting the code in 19b. We see that not only does the code replace many lines, eliminating numerous nodes in the wooden material, but also instantiates new nodes like the Voronoi Texture node with appropriate parameters to match the input text description (e.g. see the line setting `voronoi_texture_1.inputs['Randomness']`,

**Fig. 17: Blender material graph of the "digital camo" material.** To achieve the "sharp angles", our system chose to use a Voronoi texture node, and chooses the right colors in the color ramp to match the "greens, browns and tans" mentioned.



**Fig. 18: Blender material graph of the "celestial nebula" material.** Note the correspondence between "swirls" and the noise texture node, as well as the colors "violet, indigo and pink" being reflected in the color ramp.

with the comment, "# Eliminate randomness to create sharper edges"). See top left of Figure 17 for the original text description.

Figure 20 shows the size of the code edits at each iteration of a $4 \times 8$ version of BlenderAlchemy, measured in terms of *the total number of characters added/deleted* (Figure 20a), and *the number of lines added/deleted* (20b). The size of the edit of the best of iteration $i$ is measured with respect to the best of iteration $i-1$, starting at the input script of the wooden material.

We see that even though the restrictions of edits to tweak and leap edits are not strictly enforced in our procedure (and only softly done through VLM prompting), the distribution of edit size measured in both ways suggest that (1) the size of leap edits are substantially larger than those of tweak edits, and that (2) our method of oscillating between tweak and leap edits (tweaking for iterations 1 and 3, leaping for iterations 2 and 4) allows the distribution of edit size to spend every other iteration looking more similar to either tweak or leap edits. Interestingly, consistently across all the graphs shown, even when our method is producing leap edits, the average size is still lower than those of leap-

```
def shader_wood(nw: NodeWrangler, rand=False, **input_kwargs):
    # Code generated using version 2.4.3 of the node_transpiler
    texture_coordinate_1 = nw.new_node(Nodes.TextureCoord)
    mapping_2 = nw.new_node(Nodes.Mapping,
        input_kwargs={'Vector': texture_coordinate_1.outputs["Generated"], 'Rotation':
uniform(0,ma.pi*2, 3)})
    mapping_1 = nw.new_node(Nodes.Mapping,
        input_kwargs={'Vector': mapping_2, 'Scale': (0.5, sample_range(2, 4) if rand else 3,
0.5)})

    musgrave_texture_2 = nw.new_node(Nodes.MusgraveTexture,
        input_kwargs={'Vector': mapping_1, 'Scale': 2.0},
        attrs={'musgrave_dimensions': '4D'})
    if rand:
        musgrave_texture_2.inputs['W'].default_value = sample_range(0, 5)
        musgrave_texture_2.inputs['Scale'].default_value = sample_ratio(2.0, 3/4, 4/3)
    noise_texture_1 = nw.new_node(Nodes.NoiseTexture,
        input_kwargs={'Vector': musgrave_texture_2, 'W': 0.7, 'Scale': 10.0},
        attrs={'noise_dimensions': '4D'})
    if rand:
        noise_texture_1.inputs['W'].default_value = sample_range(0, 5)
        noise_texture_1.inputs['Scale'].default_value = sample_ratio(5, 0.5, 2)
    colorramp_2 = nw.new_node(Nodes.ColorRamp,
        input_kwargs={'Fac': noise_texture_1.outputs["Fac"]})
    colorramp_2.color_ramp.elements.new(0)
    colorramp_2.color_ramp.elements[0].position = 0.1727
    colorramp_2.color_ramp.elements[0].color = (0.1567, 0.0162, 0.0017, 1.0)
    colorramp_2.color_ramp.elements[1].position = 0.4364
    colorramp_2.color_ramp.elements[1].color = (0.2908, 0.1007, 0.0148, 1.0)
    colorramp_2.color_ramp.elements[2].position = 0.5864
    colorramp_2.color_ramp.elements[2].color = (0.0814, 0.0344, 0.0125, 1.0)
    if rand:
        colorramp_2.color_ramp.elements[0].position += sample_range(-0.05, 0.05)
        colorramp_2.color_ramp.elements[1].position += sample_range(-0.1, 0.1)
        colorramp_2.color_ramp.elements[2].position += sample_range(-0.05, 0.05)
        for e in colorramp_2.color_ramp.elements:
            sample_color(e.color, offset=0.03)
    colorramp_4 = nw.new_node(Nodes.ColorRamp, input_kwargs={'Fac':
noise_texture_1.outputs["Fac"]})
    colorramp_4.color_ramp.elements[0].position = 0.0
    colorramp_4.color_ramp.elements[0].color = (0.4855, 0.4855, 0.4855, 1.0)
    colorramp_4.color_ramp.elements[1].position = 1.0
    colorramp_4.color_ramp.elements[1].color = (1.0, 1.0, 1.0, 1.0)
    principled_bsdf_1 = nw.new_node(Nodes.PrincipledBSDF,
        input_kwargs={'Base Color': colorramp_2.outputs["Color"], 'Roughness':
colorramp_4.outputs["Color"]},
        attrs={'subsurface_method': 'BURLEY'})
    material_output = nw.new_node(Nodes.MaterialOutput,
        input_kwargs={'Surface': principled_bsdf_1})

def apply(obj, geo_kwargs=None, shader_kwargs=None, **kwargs):
    surface.add_material(obj, shader_wood, reuse=False, input_kwargs=shader_kwargs)

# material_obj is what we want to apply this onto.
apply(material_obj)
```

**(a)** Code for the input synthetic wood material.

```
def shader_digital_camo(nw: NodeWrangler, rand=False, **input_kwargs):
    # Create a texture coordinate node to provide coordinates for texture mapping
    texture_coordinate_1 = nw.new_node(Nodes.TextureCoord)

    # Create a mapping node to control the placement of procedural textures
    mapping_1 = nw.new_node(Nodes.Mapping,
        input_kwargs={'Vector': texture_coordinate_1.outputs["Generated"]})

    # Create a Voronoi texture that will generate the base pattern for the digital camouflage
    voronoi_texture_1 = nw.new_node(Nodes.VoronoiTexture,
        input_kwargs={'Vector': mapping_1.outputs['Vector']},
        attrs={'voronoi_dimensions': '3D'})
    voronoi_texture_1.inputs['Scale'].default_value = 150.0  # Increase the scale for a tighter
pattern
    voronoi_texture_1.inputs['Randomness'].default_value = 0.0  # Eliminate randomness to
create sharper edges
    # Set up a color ramp node to create the digital camouflage color palette
    colorramp_1 = nw.new_node(Nodes.ColorRamp,
        input_kwargs={'Fac': voronoi_texture_1.outputs['Color']})
    # Define color stop positions and colors for the digital camouflage pattern
    colorramp_1.color_ramp.elements[0].position = 0.25
    colorramp_1.color_ramp.elements[0].color = (0.0, 0.2, 0.0, 1)  # Dark green shade
    colorramp_1.color_ramp.elements.new(0.5)
    colorramp_1.color_ramp.elements[1].position = 0.50
    colorramp_1.color_ramp.elements[1].color = (0.4, 0.2, 0.0, 1)  # Brown shade
    colorramp_1.color_ramp.elements.new(0.75)
    colorramp_1.color_ramp.elements[2].position = 0.75
    colorramp_1.color_ramp.elements[2].color = (0.8, 0.8, 0.4, 1)  # Tan shade

    # Use a Principled BSDF shader to create a realistic material with the defined color
pattern
    principled_bsdf_1 = nw.new_node(Nodes.PrincipledBSDF,
        input_kwargs={'Base Color': colorramp_1.outputs["Color"]})

    # Create a material output node to connect the shader to the material's surface output
    material_output = nw.new_node(Nodes.MaterialOutput,
        input_kwargs={'Surface': principled_bsdf_1})

def apply(obj, geo_kwargs=None, shader_kwargs=None, **kwargs):
    # Apply the newly created digital camouflage shader to the provided object
    surface.add_material(obj, shader_digital_camo, reuse=False, input_kwargs=shader_kwargs)

# Assume we have a Blender object named 'material_obj' to which we want to apply the material
# Apply the new digital camouflage material to 'material_obj'
apply(material_obj)
```
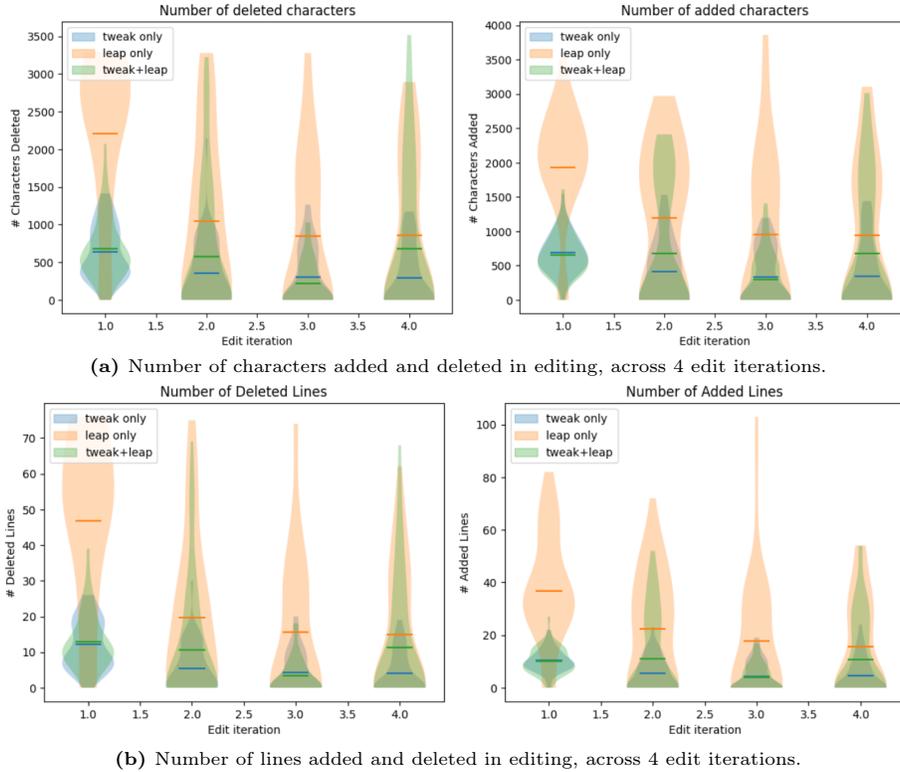
**(b)** Final output from BlenderAlchemy given the "digital camo" prompt.

**Fig. 19:** The input and output code for the "digital camouflage" example. Red parts correspond to parts that are deleted, and green corresponds to parts added by Blender-Alchemy. The darker red/green correspond to the editing of the bulk of the procedural node setup for both materials. Import statements omitted.

**(a)** Number of characters added and deleted in editing, across 4 edit iterations.



**(b)** Number of lines added and deleted in editing, across 4 edit iterations.

**Fig. 20:** Analysis of the size of code edits at every step of the editing discovery process of a $4 \times 8$ version of our system. We shown distributions for every edit iteration (1-4) of the number of characters/lines deleted (left) and added (right), for three variants of the system. **tweak-only** indicates the version where every edit iteration is requested to be tweak edits. **leap-only** is the equivalent for leap edits. **tweak+leap** alternates between the two kinds of edits – edits 1 and 3 are tweak edits, and the 2 and 4 are leap edits. Lines in each distribution indicates the mean.

only edits. We suspect that this is because tweaking in iterations (3 and 4) gets the material closer to the desired outcome, and the need for radical changes is lowered in the 2nd and 4th iterations, compared to the potentially destabilizing effects of leap-only edits.

## F    BlenderAlchemy Materials In Scenes

In this section we show the results of applying the material outputs of Blender-Alchemy onto meshes that we download from the internet. We start with a concise language description, like "demascus steel", and then expand the description into a more detailed one by prompting GPT4, using the prompt in Figure 21, following the semantic-upsampling idea in [24].

```
Come up with an in-detail caption of a material, describing the
    details of its appearance, including colors, textures, surface
    characteristics. For instance, for "Marbled Granite", output
    something like "Marbled Granite: Introduce flecks of white and
    gray to mimic the natural veining in granite, with a glossier
    finish to enhance the depth." Now do this for
    [INPUT DESCRIPTION] . Write no more than 1 sentence.
```

**Fig. 21:** Prompt used to derive more detailed description of the appearance of a material, given an abstract short description.

As the outcome of this, here's the list of (expanded) language descriptions that we use to synthesize the materials for this section:

1. Damascus Steel: Swirls of contrasting steely grays and blacks, interwoven to create a mesmerizing, wavelike pattern on the metal's surface, exhibiting a unique blend of toughness and flexibility with a semi-matte finish.
2. Brushed Aluminum: Present a sleek, matte finish with fine, unidirectional satin lines, exuding an industrial elegance in shades of silver that softly diffuse light.
3. Surface of the Sun: Envision a vibrant palette of fiery oranges, deep reds, and brilliant yellows, swirling and blending in a tumultuous dance, with occasional brilliant white flares erupting across a textured, almost liquid-like surface that seems to pulse with light and heat.
4. Ice Slabs: Crystal-clear with subtle blue undertones, showcasing intricate patterns of frozen bubbles and fractures that glisten as they catch the light, embodying the serene, raw beauty of nature's artistry.
5. Tron Legacy Material: A sleek, electric blue and black surface with a high gloss finish, featuring circuit-like patterns that glow vibrantly against the dark backdrop, evoking the futuristic aesthetic of the Tron digital world.
6. Paint Splash Material: A vibrant array of multicolored droplets scattered randomly across a stark, matte surface, creating a playful yet chaotic texture that evokes a sense of spontaneity and artistic expression.
7. Rusted Metal: A textured blend of deep oranges and browns, with irregular patches and streaks that convey the material's weathered and corroded surface, giving it a rough, tactile feel.

Since BlenderAlchemy requires a starter material script as input, we choose this for each of the above among (1) available procedural materials in Infinigen [35] and (2) materials BlenderAlchemy has synthesized thus far. For the starter materials for each, the following were decided manually:

1. Damascus Steel: We start from the **metallic swirl** material (previously synthesized by BlenderAlchemy from an input wood material – see main paper).

2. Brushed Aluminum: We start from the **metallic swirl** material (previously synthesized by BlenderAlchemy from an input wood material – see main paper).
3. Surface of the sun: we start from the **acid trip** material (previously synthesized by BlenderAlchemy from an input wood material – see main paper).
4. Ice Slabs: we start from the **chunky rock** material from the Infinigen [35] procedural material library.
5. Tron Legacy Material: we start from the **acid trip** material (previously synthesized by BlenderAlchemy from an input wood material – see main paper).
6. Paint Splash Material: we start from the **celestial nebula** material (previously synthesized by BlenderAlchemy from an input wood material – see main paper).
7. Rusted Metal : We start from the **stone** material from the Infinigen [35] procedural material library.

We use a 4x8 version of Blender Alchemy with visual imagination enabled. Look at Figures 22, 23, 24 for the application of the brushed aluminum, paint splash and Tron Legacy materials applied to a McLaren. Figures 25 and 26 show the application of the ice slabs and surface-of-the-sun materials onto a product visualization of a pair of Nikes. Finally, Figure 27 shows the application of Damascus steel and rusted metal onto two different katanas within an intriguing scene. All scenes are created using assets found on BlenderKit [?].

**(a)** Front view



**(b)** Back view



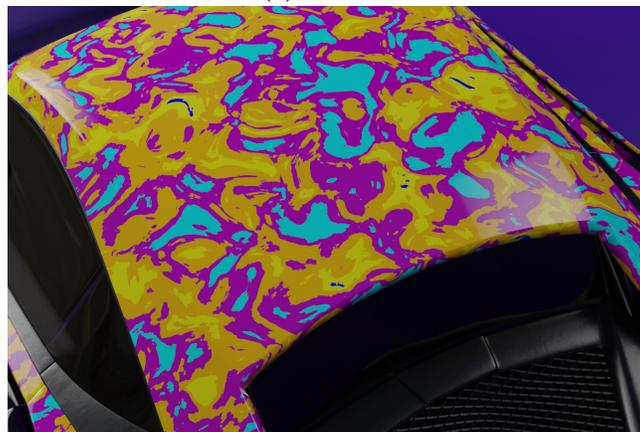**(c)** Close up view of brushed aluminum material.

**Fig. 22:** Brushed aluminum material synthesized by BlenderAlchemy, applied onto the body of a car.

**(a)** Front view



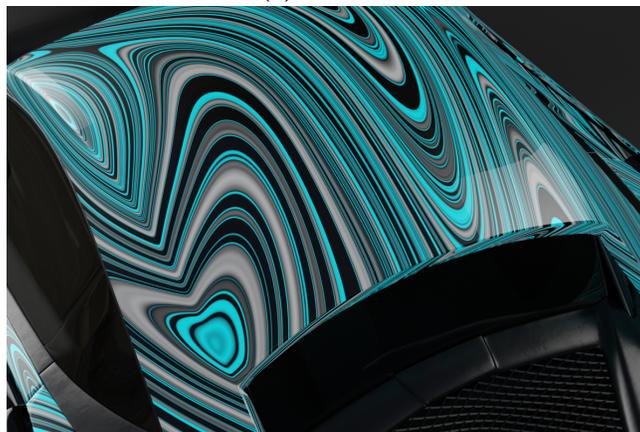**(b)** Back view



**(c)** Close up view of the "paint splash" material.

**Fig. 23:** "Paint splash" material synthesized by BlenderAlchemy, applied onto the body of a car.
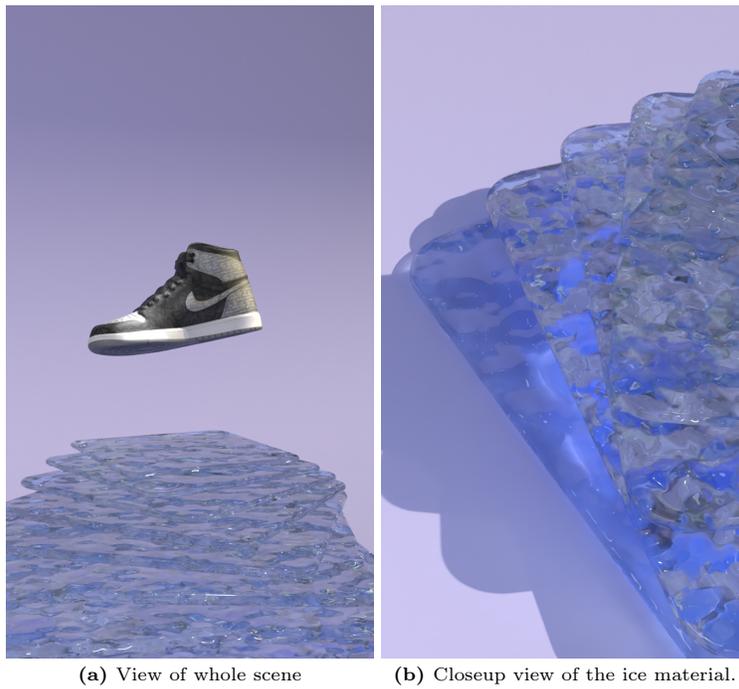
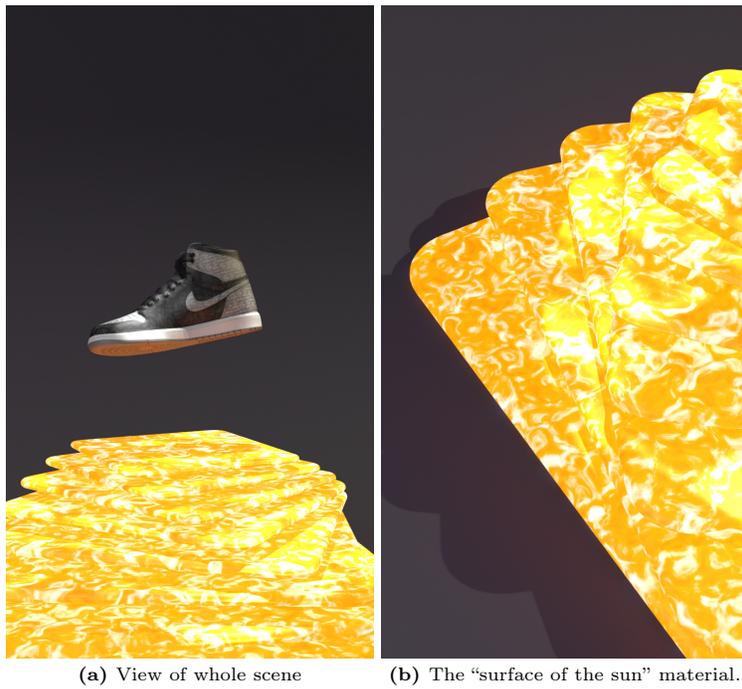**(a)** Front view



**(b)** Back view



**(c)** Close up view of the "Tron Legacy" material.

**Fig. 24:** "Tron Legacy" material synthesized by BlenderAlchemy, applied onto the body of a car.

(a) View of whole scene          (b) Closeup view of the ice material.

**Fig. 25:** "Ice" material synthesized by BlenderAlchemy, applied onto the slabs supporting the shoe.
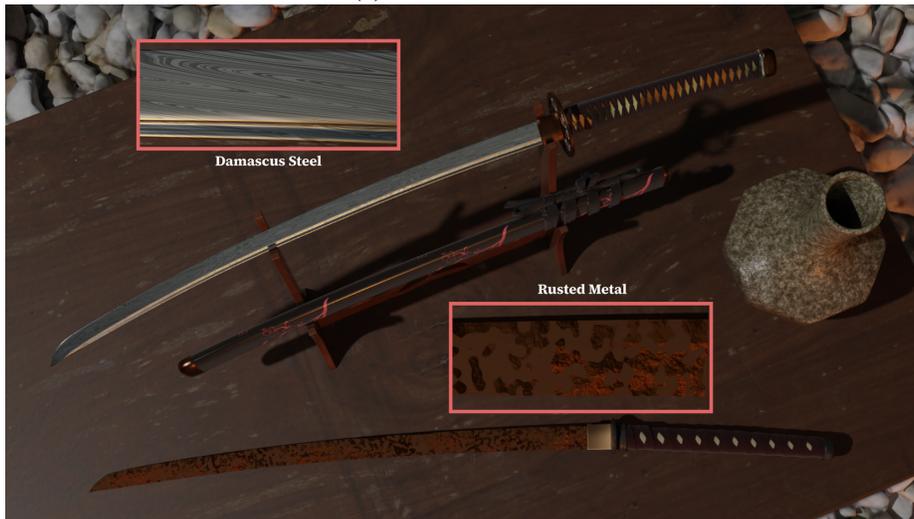
(a) View of whole scene          (b) The "surface of the sun" material.

**Fig. 26:** "Surface of the sun" material synthesized by BlenderAlchemy, applied onto the slabs supporting the shoe.

(a) View of whole scene



Damascus Steel

Rusted Metal

(b) Close up of materials

**Fig. 27:** Damascus steel and rusted metal synthesized by BlenderAlchemy, applied onto the blades of two katanas in the scene.