

Arbor: Explicit Geometric Conditioning for Controllable 3D Asset Generation

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Text and image conditioned 3D models now generate convincing assets, but they
2 still offer little direct control over the space an object should occupy or avoid. In
3 authoring, this spatial intent is often known before generation starts. A chair should
4 fit a seating envelope, a prop should leave clearance for motion, or a part should
5 expose a contact surface. Prompts and image views are poor carriers for such
6 constraints, requiring the need for an explicit control interface.

7 We present Arbor¹, a trainable attachment for text conditioned latent 3D generation.
8 Arbor introduces constraint meshes as a native 3D control interface. The interface
9 uses hull regions where geometry should exist, avoidance regions that should
10 remain empty, and touch regions the object should contact. Unlike completion
11 or whole object scaffold control, these meshes are not target evidence. They are
12 local typed requirements and can include regions where no surface should appear.
13 Arbor keeps this signal as geometry by converting constraint meshes into tokens
14 and learning a routed attachment inside a frozen denoiser. Each latent region can
15 therefore receive the part of the constraint that matters for its spatial location.

16 We evaluate Arbor on automatic and artist curated control benchmarks with hull,
17 avoidance, and touch constraints, and compare the metric trends to a user preference
18 study. Even without dedicated compliance losses, Arbor improves constraint
19 obedience while preserving object quality and variation under fixed constraints.

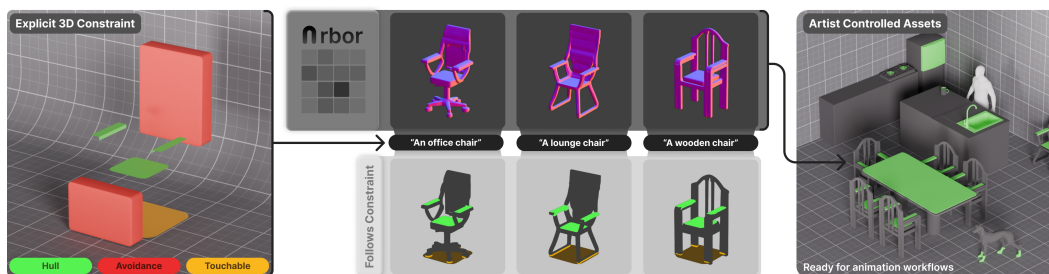


Figure 1: **Arbor overview.** Arbor turns simple 3D control objects into an explicit constraint signal for text-conditioned 3D generation. Hull regions mark where generated geometry should exist, touch regions mark contact patches, and avoidance regions mark free space that should remain empty. This enables artist to co-author the generation process, making asset generation more reliable and therefore more likely to be used in production.

¹Named after an arched support trellis to guide plant growth.

20 1 Introduction

21 3D asset creation begins with an idea, but must often align with precise spatial requirements for the
22 object to be production ready. A chair should fit a seating template. A handle should stay within
23 reach. A kit bash part should expose a clean attachment region. A prop should leave room for motion,
24 gameplay, or neighboring objects. For an artist, these requirements are easier to express with explicit
25 geometry than with words, but current 3D generators do not yet offer a strong interface for this. The
26 result is slot machine behavior: the artist tries many prompts, never lands on the exact constraint, and
27 ends up doing manual cleanup.

28 3D generation from text prompts and from images has advanced quickly, from early optimization
29 pipelines to modern feed forward and latent 3D priors [27, 13, 11, 37, 2, 36, 43, 33, 7], and now
30 produces convincing geometry from text, images, or partial observations. What these methods do not
31 provide is explicit author control over occupied and empty space when free volume starts to matter.
32 Prompts describe semantics, not space. Image conditions are tied to viewpoint and appearance,
33 making spatially absolute constraints inconvenient to encode. Arbor adds this missing interface,
34 providing explicit volumetric control over 3D asset generation without interfering with the generator’s
35 variability in shape creation.

36 We introduce a novel control mechanism which tackles this attachment and interface problem. This is
37 essential for game production as animations are often reused for multiple assets. Hence they require
38 absolute accurate interaction surfaces or the animation will not interact with the object properly
39 (clip through the surface, attach to the wrong place, etc.). While 3D control methods exist, they
40 tackle different problems on editing [6, 14], generation under global structural priors [8, 31, 29],
41 completion of partial geometry [3, 34, 25], or steering inference at sampling time [9]. For us the
42 goal is to generate a full object from a text prompt while respecting a dense local geometric hull
43 that marks semantically meaningful volumes for the object (handles, seats, wheels). Recent shape
44 guidance methods focus mainly on regions where geometry should appear. Arbor also specifies
45 regions that must stay empty (holes, gaps, clearance). We further show that the condition space can
46 define custom constraints such as touch, which marks surfaces where the asset should make contact
47 without extending to much past them (for example where legs meet the ground plane). Our constraint
48 pipeline combines hull, avoidance, and touch signals as meshes that an artist prepares in a familiar
49 modeling workflow before controlled generation, see Fig. 1.

50 Enforcing these constraints inside an existing 3D generator is not straightforward. Modern generators
51 operate on a compressed 3D latent [36, 33, 43] while the constraint meshes are dense. Two sub-
52 problems follow. We need to turn dense meshes into compact tokens that preserve both local detail
53 and typed signals, and we need to inject those tokens into the latent regions where they matter. For
54 encoding, Arbor reuses frozen geometric encoders and repurposes their material channels to carry the
55 typed signals for hull, touch, and avoidance. For injection, a geometry router groups latent queries
56 by region, retrieves the local constraint evidence relevant to each group, summarizes the full control
57 object into a small set of global tokens, and injects both through lightweight residual branches into
58 the frozen generator. Training this adapter alongside the frozen backbone on the original generation
59 objective ties the constraint latents to the generator’s own latent space to obtain semantically and
60 geometrically correct objects.

61 In summary, Arbor lets artists steer generation by defining explicit constraint meshes the generator
62 must obey, acting as an adapter that maps those constraints into the generator’s latent space. We make
63 three contributions.

- 64 • **Explicit geometric control interface.** We introduce a unified set of typed regions (hull, touch,
65 avoidance) that users specify directly as 3D meshes.
- 66 • **Encoded geometry as condition.** We show that frozen geometric encoders can be repurposed to
67 turn constraint meshes and typed signals into compact latent tokens that preserve local structure
68 and serve as a model input.
- 69 • **Geometry router and adapter.** We introduce a router that assigns local constraint evidence to
70 the latent regions where it matters, and a residual branch that injects this evidence into a frozen
71 3D generator.

72 We evaluate Arbor against the backbone without geometry and against baselines that steer during
73 sampling. Because this setting has to assess both constraint adherence and generation quality, we
74 introduce Ctrl Score and compare its trends to a user preference study.

75 2 Related work

76 **3D generation** has advanced along three practical lines. Optimization based methods lift text or
77 image supervision into explicit 3D forms such as NeRFs or Gaussian splats [16, 24, 21, 27, 28].
78 They are flexible, but often slow and sensitive to the optimization objective. Reconstruction based
79 methods infer 3D from one or a few images and are often used for generation by first producing an
80 image and then recovering 3D [13, 30, 37, 2, 15, 7]. Native latent 3D models such as Direct3D [33],
81 TRELIS [36], and Hunyuan3D 2.0 [43] instead denoise compact 3D states directly. A related family
82 of generators produces structured outputs such as parts, layouts, assemblies, or CAD primitives [12,
83 18, 39, 19, 42, 38, 20], where authoring happens by editing or completing the structured output
84 after generation. Arbor instead brings the artist into the generation step itself, giving native latent
85 generators the direct interface they currently lack.

86 ControlNet and Adapter methods showed that a pretrained text to image diffusion model can absorb
87 new spatial conditions through a control branch trained alongside the frozen backbone [41, 22, 40].
88 The 3D version of the problem is harder, because a condition must remain meaningful across
89 viewpoints, compressed latent states, and the final output representation. Arbor inherits the attached
90 control intuition, but rebuilds it for native 3D generation where the condition is itself geometry rather
91 than a 2D map.

92 **Editing** methods receive a complete asset or scene up front and modify it in a chosen region.
93 SIGNeRF edits NeRF scenes with depth conditioned reference sheets [6]. Instant3dit and ObjFiller-
94 3D inpaint 3D objects through multiview image edits before mapping the result back to 3D [1, 10],
95 while Easy3E performs feed forward asset editing directly in a voxel flow [14]. Each method requires
96 the user to bring the asset they want to modify. Arbor instead conditions a generator before any asset
97 exists, with only typed local constraint regions on top of a text prompt.

98 **Reconstruction** methods condition on an input that already describes the asset’s global structure.
99 Several methods rely on a 3D structural prior. SK-Adapter encodes a skeleton into tokens injected
100 into a frozen TRELIS backbone [31]. Coin3D conditions on a coarse primitive proxy through a 3D
101 adapter [8]. Others retrieve a 3D reference shape as a similarity prior [32]. Other methods rely on
102 an image. SPAR3D conditions on an image plus an editable point cloud [15]. Hunyuan3D-Omni
103 combines an image with several geometric modalities including points, boxes, voxels, and skeletons.
104 It accepts partial observations but still requires the image to anchor the asset [29]. Each input fixes
105 the asset’s global structure ahead of generation, whether as a skeleton, a proxy, a reference, or an
106 image. Arbor specifies only the parts that matter and leaves the rest to the prompt and the generator,
107 allowing different assets to satisfy the same constraint set.

108 **Inpainting and completion** methods start from sparse or partial geometry and grow it into a finished
109 asset. DiffComplete and Points-to-3D condition a 3D diffusion model on partial geometry and
110 complete the rest around the input surface [3, 34], while Spice-E uses cross entity attention between
111 a coarse guidance shape and the noisy sample [25]. SpaceControl steers a pretrained 3D denoiser
112 during sampling without extra training, exposing a global tradeoff between constraint fidelity and
113 generative variation [9]. None of these distinguish typed roles for the input geometry, treating it as a
114 single signal to complete, match, or steer toward. Arbor uses geometry as a typed specification. Hull
115 regions indicate where the asset should exist, touch regions mark contact surfaces, and avoidance
116 regions are defined precisely by not becoming surface. In contrast to methods that prefill the latent
117 or attend to a coarse proxy, Arbor guides denoising with explicit structure while the artist specifies
118 volumes and signals up front.

119 3 Method

120 The main goal of Arbor is the creation of explicit geometry control methods for 3D generation. Given
121 a text prompt y and a constraint object C , the goal is to sample a 3D asset $x \sim p(x | y, C)$.

122 3.1 Overview

123 Arbor builds on existing 3D generative models from the TRELIS family. We use Trellis 1 (T1) [36]
124 as this strong text-conditioned backbone for 3D geometry generation. To impose surface constraints,
125 we leverage TRELIS.2 (T2) [35], whose direct surface encoding mechanism provides a flexible

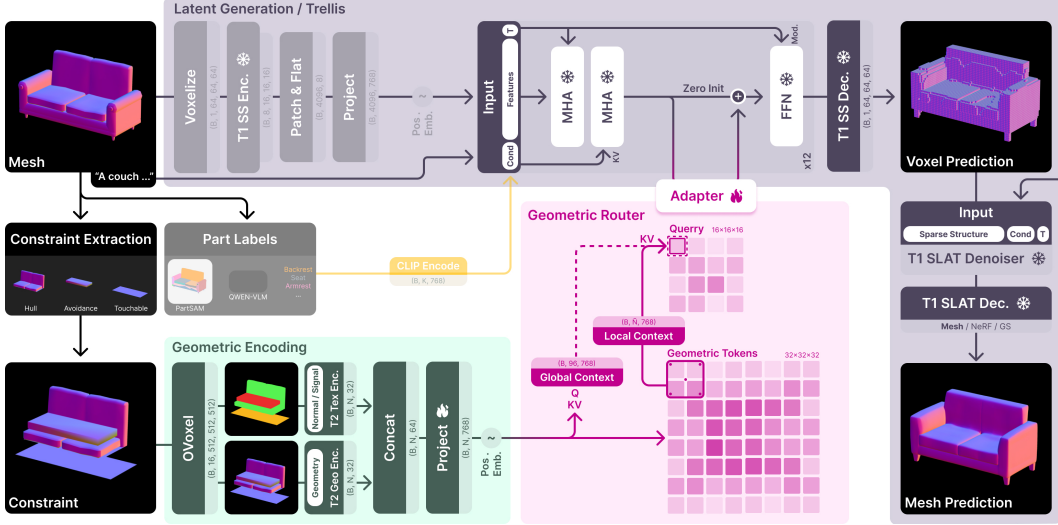


Figure 2: **Constraint conditioning pipeline.** Arbor converts a typed constraint object into TRELIS.2 OVoxels, encodes geometry and signal attributes with frozen encoders (Sec. 3.2), aligns the resulting latents into geometry tokens, and routes those tokens into the TRELIS sparse structure denoiser (Sec. 3.3). Local routing gives each query group the nearby constraint evidence it needs, while learned global summaries provide object scale context. (Part Labels and yellow path are for Arbor Semantics A.3)

126 interface for conditioning generation on explicit geometric signals. Both models follow a common
 127 compression generation design: raw 3D representations are first mapped to compact latent states, and
 128 then generation is performed in the compressed domain.

129 T2 also supports auxiliary per-surface channels, originally intended for material parameters used
 130 in rendering. We repurpose these channels to encode constraint information as binary flags, where
 131 each flag indicates the presence or absence of a specific surface constraint. For example, channel (c)
 132 0-2 encode normals, c3 hull, c4 avoid, and c5 touch. We then inject them through a learned routed
 133 residual branch as shown in Fig. 2.

134 **TRELIS 1 - Text conditioned 3D generation** T1 is a text-conditioned 3D generator that separates
 135 generation into a sparse-structure stage followed by latent refinement. Its sparse-structure stage
 136 operates on a 16^3 latent grid with 8 channels per lattice cell and is trained with flow matching. We
 137 use this stage as the prompt-conditioned backbone in Arbor, since it determines the coarse object
 138 structure and exposes a text-conditioned denoising process.

139 At each transformer block, the latent state is updated by self-attention, text cross-attention, and a
 140 feed-forward network. Omitting residual connections, normalization, and adaptive time modulation,
 141 one block can be written as

$$z_{\text{text}}^{(\ell)} = \text{CrossAttn}_{\text{text}}^{(\ell)}\left(\text{SelfAttn}^{(\ell)}(z^{(\ell-1)}), y\right), \quad z^{(\ell)} = \text{FFN}^{(\ell)}(z_{\text{text}}^{(\ell)}). \quad (1)$$

142 Here, ℓ indexes the transformer layer, $z^{(\ell-1)}$ is the incoming sparse-structure state, $z_{\text{text}}^{(\ell)}$ is the state
 143 after cross-attention to the text prompt, and y denotes the encoded text context. The diffusion or flow
 144 time t enters through adaptive modulation, which is omitted from Eq. (1) for clarity.

145 During inference, T1 denoises an initial noise latent under the prompt conditioning y and time t .
 146 The resulting sparse structure is decoded into a 64^3 occupancy grid and passed to a second latent
 147 refinement stage. This refinement stage operates on active voxels and produces the final representation,
 148 which can be decoded as a mesh, radiance field, or 3D Gaussian representation.

149 **TRELIS 2 - Constraint Encoding** T2 is designed around a direct sparse voxel representation,
 150 OVoxels, which jointly represents geometry and aligned surface attributes. While T2 is image-
 151 conditioned as a generator, Arbor uses only its frozen encoder stack. This makes T2 suitable as a
 152 constraint encoder without changing Arbor’s prompt-driven generation setting.

153 Given a mesh M with aligned surface attributes A , T2 first converts the input into aligned shape and
 154 material fields on a 512^3 voxel grid:

$$\text{OVoxelize}(M, A) = (\widetilde{M}_{\text{shape}}, \widetilde{M}_{\text{mat}}). \quad (2)$$

155 The frozen T2 shape and material encoders then map these fields to compact sparse latents,

$$z_{\text{shape}} = E_{\text{T2,shape}}(\widetilde{M}_{\text{shape}}), \quad z_{\text{mat}} = E_{\text{T2,mat}}(\widetilde{M}_{\text{mat}}). \quad (3)$$

156 Each encoder reduces spatial resolution by a factor of 16 and produces sparse 32-dimensional tokens.

157 Arbor repurposes this T2 encoding path to represent geometric constraints. In particular, we encode
 158 constraint signals through the auxiliary attribute channels normally used for material parameters.
 159 These channels provide a direct, surface-aligned interface for injecting binary constraint flags, while
 160 keeping the T1 text-conditioned generation backbone unchanged.

161 3.2 Constraints

162 Constraints in Arbor are 3D meshes that specify desired spatial behavior. We use hull, avoidance,
 163 and touch constraints, and encode them into compact latent tokens for the text conditioned generator.

164 **Types** Hull constraints define regions where generated geometry should exist, e.g., the seat of a
 165 chair. Avoidance constraints define regions where generated geometry should not exist, e.g., empty
 166 space above the seat. Touch constraints define surfaces the object should contact, e.g., chair legs
 167 touching the ground. A.1 describes constraint creation. Here, we focus on how constraints are
 168 introduced into the model.

169 **Encoding** The natural first option would be to encode constraints with the native T1 encoder, since
 170 T1 is the generator we want to influence. This representation, however, is poorly suited to model
 171 interactive regional control. T1 encoding was built for training the latent space, not for interactive
 172 conditioning. A constraint would need to be voxelized, rendered from 150 views, and projected
 173 with DINOv2 [23] features. The process is slow, expects complete objects, and has no direct place
 174 for typed control signals. We therefore encode constraints with the frozen T2 stack instead. Arbor
 175 fuses the separate constraint meshes into a single mesh C_{mesh} with surface normals C_{normals} and
 176 signal channels C_{signal} that encode the constraint type. Applying Eq. 2 to this constraint object yields
 177 the OVoxel shape field $\widetilde{C}_{\text{shape}}$ and aligned fields $\widetilde{C}_{\text{normals}}$ and $\widetilde{C}_{\text{signal}}$. The shape encoder receives
 178 $\widetilde{C}_{\text{shape}}$, while the attribute encoder is repurposed. Instead of the original 6 material channels, it
 179 receives 3 voxel aligned normal channels $\widetilde{C}_{\text{normals}}$ and 3 binary control channels for typed signals
 180 $\widetilde{C}_{\text{signal}}$. Eq. 4 is the Arbor version of the T2 encoder call in Eq. 3.

$$c_{\text{shape}} = E_{\text{T2,shape}}(\widetilde{C}_{\text{shape}}), \quad c_{\text{signal}} = E_{\text{T2,mat}}(\widetilde{C}_{\text{normals}} ; \widetilde{C}_{\text{signal}}). \quad (4)$$

181 The semicolon denotes channel-wise concatenation. We concatenate both latent streams into the
 182 geometry memory $c_{\text{geo}} = (c_{\text{shape}} \ c_{\text{signal}})$. Each geometry token keeps its OVoxel position p_i . While
 183 this encoding is simple and effective, two issues remain. The tokens are not native to the T1 latent
 184 space, and they live on a finer 32^3 grid than the T1 16^3 state. We map every geometry token to the T1
 185 model width with a learned projection and add a learned 3D positional embedding from p_i . We keep
 186 the notation c_{geo} for these prepared tokens below. The router handles the resolution mismatch by
 187 choosing which geometry tokens each local group of T1 queries receives.

188 3.3 Control

189 Generated geometry should obey the constraint, but still follow the text prompt. These goals can
 190 conflict. A model can improve hull overlap by overfilling the guide while losing prompt fit, structure,
 191 or variation. This is why sampling time steering [9], including our Gradient baseline (Appendix A.3),
 192 are an important but incomplete alternative. Arbor instead learns the attachment inside the generator.
 193 An adapter injects constraint tokens, and a router decides which tokens each latent region receives
 194 (Fig. 2, pink path). Each SS block carries one hidden query token for every cell of the 16^3 lattice.
 195 Arbor partitions this query lattice into 64 groups indexed by g , each covering a $4 \times 4 \times 4$ block
 196 of neighboring queries. For each group, the router builds a geometry context from the projected
 197 constraint tokens. The adapter uses the current T1 hidden states as queries, attends to this context,
 198 and writes the result back as a residual update.

199 **Adapter** Arbor injects geometry with a separate grounding branch inside each T1 block, after
 200 frozen text cross attention and before the feed forward (FFN) update. Eq. 5 modifies the T1 block
 201 from Eq. 1 by inserting a geometry residual before the original FFN. For the queries in group g at
 202 block ℓ , with router context $c_{\text{ctx}}^{(g)}$ defined below, we write

$$\Delta \mathbf{z}_{\text{geo}}^{(\ell,g)} = \text{CrossAttn}_{\text{geo}}^{(\ell)}(\mathbf{z}_{\text{text}}^{(\ell,g)}, c_{\text{ctx}}^{(g)}), \quad \mathbf{z}^{(\ell,g)} = \text{FFN}^{(\ell)}(\mathbf{z}_{\text{text}}^{(\ell,g)} + W_{\text{geo}}^{(\ell)} \Delta \mathbf{z}_{\text{geo}}^{(\ell,g)}). \quad (5)$$

203 Here $\mathbf{z}_{\text{text}}^{(\ell,g)}$ denotes the subset of text conditioned hidden states in query group g , $\Delta \mathbf{z}_{\text{geo}}^{(\ell,g)}$ is the
 204 geometry update, and $W_{\text{geo}}^{(\ell)}$ is zero initialized. $\text{CrossAttn}_{\text{geo}}^{(\ell)}$ is one packed attention pass over
 205 local tokens and global summary tokens. This placement keeps the pretrained text path intact while
 206 allowing geometry to affect the state before the FFN. Empirically, later insertion, routing geometry
 207 through the text conditioning path, or unfreezing larger parts of the pretrained model led to weaker
 208 results. The remaining challenge is scale, since Arbor conditions on dense encoded geometry rather
 209 than the small global token set used by skeleton adapters such as SK-Adapter [31]. Attending from
 210 all T1 queries to all constraint tokens would be expensive and would require truncation.

211 **Router** The router makes dense geometry memory usable by the SS denoiser. Routing is separate
 212 from the position embeddings introduced above. The embedding marks where each geometry token
 213 lies, while routing chooses which tokens are read by each query group. Before denoising, Arbor
 214 prepares the projected geometry memory c_{geo} , token positions p_i , and the 64 fixed query groups,
 215 which depend only on the constraint object and the T1 lattice. At block ℓ , the original T1 stream
 216 supplies the attention queries, namely the text conditioned hidden states $\mathbf{z}_{\text{text}}^{(\ell,g)}$ in Eq. 5. The router
 217 supplies the keys and values by retrieving local geometry for each query group and appending a
 218 compact global context. All queries in one group share the context,

$$c_{\text{ctx}}^{(g)} = [c_{\text{local}}^{(g)}; c_{\text{global}}]. \quad (6)$$

219 For local routing, each group uses its center and eight corners as routing anchors R_g . We score every
 220 geometry token by its distance to the nearest anchor and keep the $K = 2048$ nearest tokens,

$$c_{\text{local}}^{(g)} = \text{TopK}_{2048} \left(\{ (c_{\text{geo}}^{(i)}, p_i) \}_{i=1}^N, \min_{r \in R_g} \|p_i - r\|_2 \right). \quad (7)$$

221 Here, p_i is the 3D position of geometry token $c_{\text{geo}}^{(i)}$, N is the number of geometry tokens, and r ranges
 222 over the anchors in R_g . This TopK step is not attention. It selects a bounded local memory for Eq. 5,
 223 keeping cost fixed while tying the context to the relevant constraint region. To retain object level
 224 information that local routing can miss, Arbor also summarizes the full geometry memory with 96
 225 learned summary tokens,

$$c_{\text{global}} = \text{MLP} \left(\text{CrossAttn}(S_{\text{global}}, \{c_{\text{geo}}^{(i)}\}_{i=1}^N) \right), \quad (8)$$

226 where S_{global} is the learned summary set. These tokens capture broader signals such as extent, touch
 227 direction, and the relation between occupied and forbidden regions.

228 **Training.** We train only the geometry facing modules. These are the geometry projection and
 229 position embedding, the global summary modules, the semantic part token modules, and the grounding
 230 adapters. The T1 self attention, text cross attention, and feed forward weights remain frozen.
 231 Optimization uses the standard SS flow matching objective of the T1 backbone. No explicit constraint
 232 compliance loss is used in this run family. Text and geometry conditions are dropped independently
 233 for classifier free guidance. Training examples are built by sampling an object, constructing or
 234 selecting a typed constraint set, encoding it as above, and using the original T1 sparse structure latent
 235 as the flow matching target. Details, schedules, and implementation settings are deferred to A.1.

236 4 Evaluation

237 We evaluate Arbor as a control interface for text-conditioned 3D generation. The goal is not to copy
 238 a guide into the output, but to integrate constraint meshes while still producing a plausible object
 239 that follows the prompt. We test control over the same backbone without geometry, comparison to
 240 geometry guided baselines, and variation under fixed constraints.

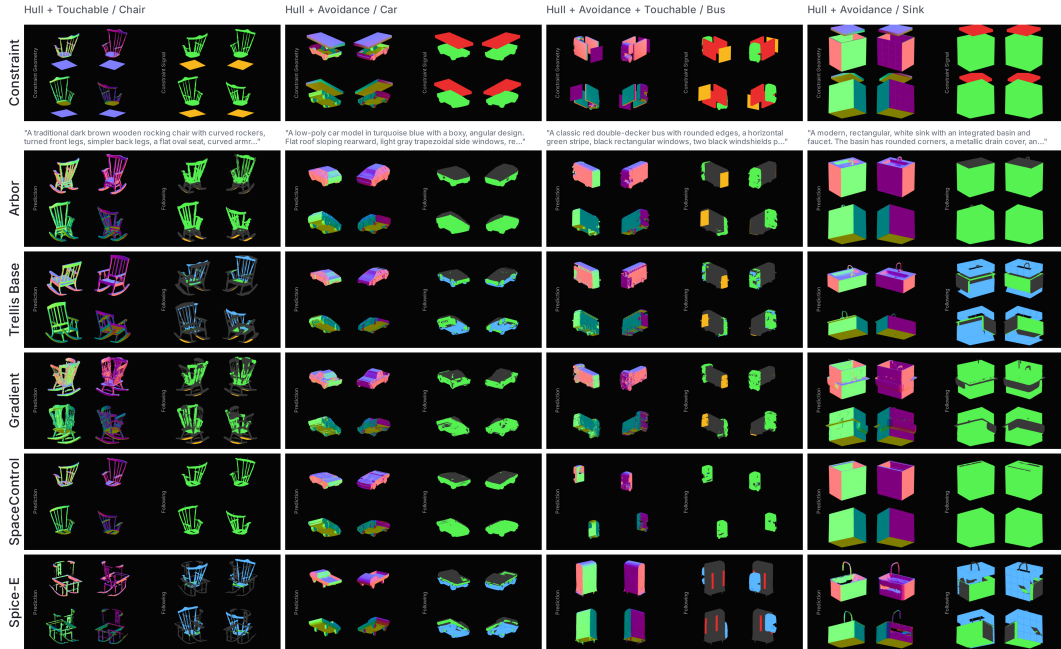


Figure 3: **Controlled generation comparison.** Each column shows one prompt and constraint object. The constraint is rendered as normal shaded geometry with signal regions colored **hull**, **touch**, and **avoidance**. Rows compare predictions and their constraint following. Here, **green** indicates a hull match and **blue** indicates missing hull. Arbor keeps readable objects while following local roles.

241 **Models.** Arbor is the model from Sec. 3. Both the TRELIS [36] text generator and the TREL-
 242 LIS.2 [35] OVOxel encoders stay frozen and only the Arbor geometry modules are trained. We also
 243 report Arbor Semantics, which adds stronger per query semantic text cues, and Arbor Compliance,
 244 which finetunes Arbor with explicit hull, avoidance, and touch losses (Appendix A.3). We com-
 245 pare against TRELIS, our Gradient baseline at sampling time (Appendix A.3), SpaceControl [9],
 246 and Spice-E [25] for controlled generation. The variation track adds Point-E, SPAR3D [15], and
 247 Hunyuan3D-Omni [29], which receive image cues that the others do not.

248 **Datasets.** Arbor is trained on roughly 50k objects sampled from ABO [4], HSSD [17], and a
 249 Sketchfab subset of Objaverse-XL [5], which is about 10% of the training volume reported by
 250 TRELIS [36]. Evaluation uses Toys4K [26], the dataset on which TRELIS reports its official
 251 numbers. We construct two control benchmarks on Toys4K with hull, avoidance, and touch signals.
 252 The automatic split contains 128 procedurally generated constraints; the manual split contains 32
 253 hand authored constraints that are not sampled by the training program (Appendix A.1).

254 **Metrics.** We report Hull Hit, Avoidance Violation (Avoid Viol.), Touch Hit, Volume Match (Vol.
 255 Match), multiview CLIP (MV-CLIP), and Control Score (Ctrl. Scr.), our new addition (Appendix A.5).
 256 All geometry terms use a shared 64^3 voxel grid, matching the sparse structure resolution used by
 257 the backbone. Ctrl. Scr. is a per sample harmonic mean over the terms where higher values are
 258 better (Hull Hit, Touch Hit, $1 - \text{Avoid Viol.}$, Vol. Match, MV-CLIP), so a method has to do well on
 259 all of them at once. Vol. Match is a coarse size guard that prevents overfilled outputs from looking
 260 artificially complete.

261 4.1 Controlled generation

262 Fig. 3 and Tab. 1 show the main result. TRELIS keeps the object prior strong, but has no mechanism
 263 for the typed geometry and only matches the guide by chance. Gradient and SpaceControl move
 264 mass toward the control object, yet this often comes at the cost of noisy geometry, missing structure,
 265 or a collapsed shape. Spice-E can preserve recognizable form, but treats the guide as a shape signal
 266 and does not reliably separate hull, avoidance, and touch roles. Arbor keeps both requirements
 267 visible. The outputs remain readable assets, and the following views show local roles respected in the
 268 intended regions.

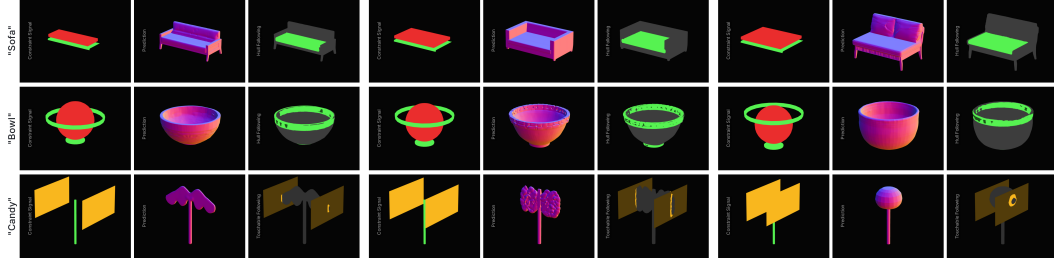


Figure 4: **Constraint sweeps.** The prompt is fixed and the constraint region is continuously moved, scaled, or rotated. Arbor follows the deformation without snapping to a small set of canonical layouts.

Table 1: **Controlled generation benchmark.** Manual ($n=32$) and automatic ($n=128$) splits. Pref. is the pairwise user study win rate over 404 trials from 27 participants, with the parenthesized value merging the three Arbor variants. Geometry metrics use 64^3 voxels and Ctrl. Scr. is defined in Appendix A.5.

Method	User Study		Manual (n=32)					Auto (n=128)					
	Pref. (%) \uparrow	Ctrl. Scr. \uparrow	Hull Hit \uparrow	Avoid Viol. \downarrow	Touch Hit \uparrow	Vol. Match \uparrow	MV-CLIP \uparrow	Ctrl. Scr. \uparrow	Hull Hit \uparrow	Avoid Viol. \downarrow	Touch Hit \uparrow	Vol. Match \uparrow	MV-CLIP \uparrow
Arbor	45.9 (59.2)	0.402	0.714	0.025	0.857	0.469	0.229	0.472	0.786	0.006	0.984	0.487	0.240
Arbor Semantics	27.7	0.355	0.600	0.019	0.714	0.481	0.231	0.467	0.775	0.005	0.984	0.485	0.238
Arbor Compliance	29.9	0.401	0.772	0.015	0.857	0.441	0.233	0.466	0.802	0.003	0.968	0.481	0.240
Trellis	33.3	0.253	0.283	0.026	0.714	0.423	0.235	0.267	0.269	0.019	0.667	0.340	0.245
Gradient	18.2	0.347	0.690	0.046	0.571	0.471	0.230	0.436	0.775	0.019	0.833	0.474	0.236
SpaceControl	9.6	0.151	0.823	0.000	0.000	0.188	0.220	0.214	0.729	0.001	0.492	0.166	0.234
Spice-E	5.2	0.151	0.108	0.085	0.286	0.469	0.228	0.227	0.196	0.031	0.508	0.465	0.239

269 The table mirrors this qualitative behavior. Arbor and Compliance are essentially tied on the manual
 270 split, and all Arbor variants separate clearly from the non-Arbor baselines. The Arbor family wins
 271 59.2% of pairwise user choices in our 27 participant, 404 trial study, the single Arbor row is preferred
 272 most often before merging variants. TRELLIS is the next strongest human baseline because it
 273 produces clean objects, even though it misses many constraints. For this reason, Ctrl. Scr. combines
 274 control adherence, volume agreement, and MV-CLIP instead of reporting geometry overlap alone.

275 4.2 Variation under fixed constraints

276 The prompt steerable baselines above either miss the control object or degrade the sample. We
 277 therefore also compare to recent reconstruction models with explicit 3D conditioning. These methods
 278 solve a different task because they receive an image input, but test whether stronger visual evidence
 279 is enough to combine control and variation. We fix one control object and vary the seed, while
 280 keeping the image input fixed for Point-E, SPAR3D, and Hunyuan3D-Omni. In Fig. 5, Arbor changes
 281 the truck back and tires while respecting hull and avoidance, and it builds new sofa cushions and
 282 surrounding frame geometry around the controlled seating area. Point-E reaches similar raw variation,
 283 but the following views show that much of this variation is unstable and leaves the hull. SPAR3D
 284 and Hunyuan3D-Omni stay closer to the image and therefore vary less. Tab. 2 confirms the tradeoff.
 285 Arbor reaches the highest variation and Ctrl. Scr. even without an image anchor, while image bound
 286 models either lose control or lose generative range.

287 4.3 Ablations

288 **Routing.** We probe the conditioning path at a fixed checkpoint to identify which signal Arbor actually
 289 uses. Tab. 3 reports constraint IoU retained relative to the full model. Local routing alone keeps
 290 84.3% of the constraint IoU, showing that routed local evidence is the main mechanism. The gap
 291 to the full model shows that global summaries still add useful object level context, while global
 292 summaries alone are not sufficient. Removing the signal stream is most damaging because the model
 293 still sees geometry but no longer knows which regions should be filled, touched, or avoided. The
 294 shape and normal rows show that Arbor uses more than token location. Encoded shape carries
 295 substantial information, and normals provide a smaller but still visible gain.

296 **Sweeps.** Fig. 4 tests whether this control stays smooth outside the discrete benchmark. The prompt is
 297 fixed while a single hull region moves through position, scale, and orientation. The generated asset
 298 follows the moving constraint without snapping to a small set of layouts or losing object identity. The

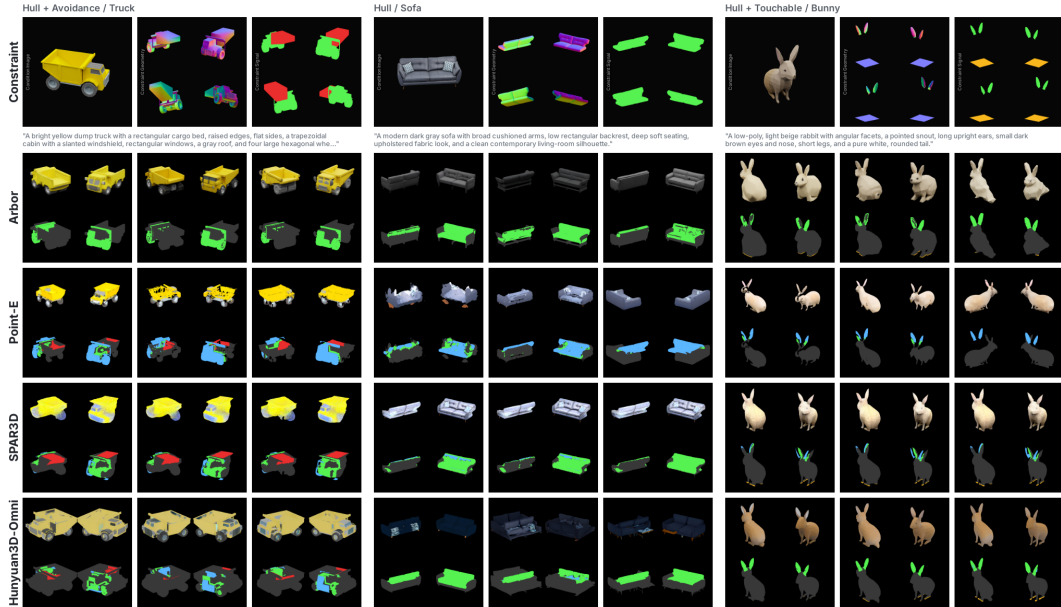


Figure 5: **Variation under a fixed constraint.** Each block keeps the hull fixed and varies the seed; image conditioned baselines also fix the input image. Arbor changes details and proportions across seeds while still satisfying the constraint, where the image anchored methods stay close to their input.

Method	Var.↑	Ctrl. Scr.↑	Hull Hit↑	Avoid Viol.↓	Vol. Match↑	MV-CLIP↑
Arbor	0.740	0.361 ± 0.010	0.707 ± 0.044	0.016 ± 0.010	0.460 ± 0.024	0.229 ± 0.002
Point-E	0.731	0.105 ± 0.009	0.100 ± 0.014	0.089 ± 0.015	0.364 ± 0.023	0.224 ± 0.002
SPAR3D	0.141	0.162 ± 0.004	0.182 ± 0.001	0.198 ± 0.003	0.565 ± 0.002	0.224 ± 0.000
Hunyuan3D-Omni	0.526	0.318 ± 0.034	0.533 ± 0.028	0.040 ± 0.027	0.534 ± 0.023	0.230 ± 0.001

Table 2: **Variation under a fixed constraint.** Var. is mean pairwise 1-IoU across three seeds; other columns report mean \pm standard deviation. Arbor keeps high variation while preserving control.

Ablation type	IoU retained↑
Only local routing	84.3%
Only global summary	27.0%
No shape stream	54.9%
No signal stream	39.2%
No normals	83.8%

Table 3: **Conditioning ablation.** Constraint IoU retained relative to full Arbor.

299 sweep grid also isolates hull, touch, and avoidance controls, showing that Arbor reacts to each signal
 300 rather than only to a single combined constraint. We include additional sheets and video sweeps in
 301 the supplemental material.

302 5 Conclusion

303 Arbor introduces an explicit geometric conditioning interface for a frozen text conditioned 3D
 304 generator. Constraint meshes are turned into compact tokens by frozen 3D encoders and injected
 305 through a routed residual branch into the denoising blocks, so that hull, touch, and avoidance regions
 306 become part of the generator’s input rather than only a sampling time correction. On our Toys4K
 307 benchmarks, Arbor improves constraint adherence over the backbone without geometry, sampling
 308 time baselines, and trained baselines, while preserving the variation of the underlying prior under
 309 fixed constraints.

310 The current system still exposes two limits. First, constraint regions carry geometry and typed signals,
 311 but not full semantic function. A seat region gives the generator a volume and an orientation, but
 312 not a guarantee that the volume will be used as a seat when the prompt conflicts with the constraint.
 313 Our semantic variant did not yet outperform the routed geometry path itself. Second, Arbor acts only
 314 at the sparse structure stage and does not directly control later refinement, where surface detail and
 315 material attributes enter. These limits point to the next step: richer part labels and the same routed
 316 attachment extended to later stages, so that structure, detail, and material can follow one explicit
 317 geometric specification. Further limitations are recorded in A.6.

318 **References**

- 319 [1] Amir Barda, Matheus Gadelha, Vladimir G. Kim, Noam Aigerman, Amit H. Bermano, and
320 Thibault Groueix. Instant3dit: Multiview inpainting for fast editing of 3D objects. In *Conference*
321 *on Computer Vision and Pattern Recognition (CVPR)*, pages 16273–16282, 2025. URL <https://arxiv.org/abs/2412.00518>.
322
- 323 [2] Mark Boss, Zixuan Huang, Aaryaman Vasishta, and Varun Jampani. SF3D: Stable fast 3D
324 mesh reconstruction with UV-unwrapping and illumination disentanglement. In *Conference on*
325 *Computer Vision and Pattern Recognition (CVPR)*, pages 16240–16250, 2025. URL <https://arxiv.org/abs/2408.00653>.
326
- 327 [3] Ruihang Chu, Enze Xie, Shentong Mo, Zhenguo Li, Matthias Nießner, Chi-Wing Fu, and
328 Jiaya Jia. DiffComplete: Diffusion-based generative 3D shape completion. In *Advances in*
329 *Neural Information Processing Systems (NeurIPS)*, 2023. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2306.16329)
330 [2306.16329](https://arxiv.org/abs/2306.16329).
- 331 [4] Jasmine Collins, Shubham Goel, Kenan Deng, Achleshwar Luthra, Leon Xu, Erhan Gundogdu,
332 Xi Zhang, Tomas F. Yago Vicente, Thomas Dideriksen, Himanshu Arora, Matthieu Guillaumin,
333 and Jitendra Malik. ABO: Dataset and benchmarks for real-world 3D object understanding.
334 In *Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2110.06199)
335 [2110.06199](https://arxiv.org/abs/2110.06199).
- 336 [5] Matt Deitke, Ruoshi Liu, Matthew Wallingford, Huong Ngo, Oscar Michel, Aditya Kusupati,
337 Alan Fan, Christian Laforte, Vikram Voleti, Samir Yitzhak Gadre, Eli VanderBilt, Aniruddha
338 Kembhavi, Carl Vondrick, Georgia Gkioxari, Kiana Ehsani, Ludwig Schmidt, and Ali Farhadi.
339 Objaverse-XL: A universe of 10M+ 3D objects. In *Advances in Neural Information Processing*
340 *Systems (NeurIPS)*, 2023. URL <https://arxiv.org/abs/2307.05663>.
- 341 [6] Jan-Niklas Dihlmann, Andreas Engelhardt, and Hendrik Lensch. SIGNeRF: Scene integrated
342 generation for neural radiance fields. In *Conference on Computer Vision and Pattern Recognition*
343 *(CVPR)*, 2024. doi: 10.1109/CVPR52733.2024.00638. URL [https://doi.org/10.1109/](https://doi.org/10.1109/CVPR52733.2024.00638)
344 [CVPR52733.2024.00638](https://doi.org/10.1109/CVPR52733.2024.00638).
- 345 [7] Jan-Niklas Dihlmann, Mark Boss, Simon Donne, Andreas Engelhardt, Hendrik P. A. Lensch,
346 and Varun Jampani. ReLi3D: Relightable multi-view 3D reconstruction with disentangled
347 illumination. In *International Conference on Learning Representations (ICLR)*, 2026. URL
348 <https://openreview.net/forum?id=B1SKgQb3Vd>.
- 349 [8] Wenqi Dong, Bangbang Yang, Lin Ma, Xiao Liu, Liyuan Cui, Hujun Bao, Yuewen Ma,
350 and Zhaopeng Cui. Coin3D: Controllable and interactive 3D assets generation with proxy-
351 guided conditioning. In *ACM SIGGRAPH*, 2024. doi: 10.1145/3641519.3657425. URL
352 <https://doi.org/10.1145/3641519.3657425>.
- 353 [9] Elisabetta Fedele, Francis Engelmann, Ian Huang, Or Litany, Marc Pollefeys, and Leonidas
354 Guibas. SpaceControl: Introducing test-time spatial control to 3D generative modeling.
355 In *International Conference on Learning Representations (ICLR)*, 2026. URL <https://openreview.net/forum?id=mEqsCVI5sN>.
356
- 357 [10] Haitang Feng, Jie Liu, Jie Tang, Gangshan Wu, Beiqi Chen, Jianhuang Lai, and Guangcong
358 Wang. ObjFiller-3D: Consistent multi-view 3D inpainting via video diffusion models. *arXiv*
359 *preprint*, 2025. URL <https://arxiv.org/abs/2508.18271>.
- 360 [11] Zexin He and Tengfei Wang. OpenLRM: Open-source large reconstruction models. [https://](https://github.com/3DTopia/OpenLRM)
361 github.com/3DTopia/OpenLRM, 2023. URL <https://github.com/3DTopia/OpenLRM>.
362 GitHub repository; open-source implementation of LRM, not a primary paper.
- 363 [12] Amir Hertz, Or Perel, Raja Giryes, Olga Sorkine-Hornung, and Daniel Cohen-Or. SPAGHETTI.
364 *ACM Transactions on Graphics (TOG)*, 2022. doi: 10.1145/3528223.3530084. URL <https://doi.org/10.1145/3528223.3530084>.
365
- 366 [13] Yicong Hong, Kai Zhang, Jiuxiang Gu, Sai Bi, Yang Zhou, Difan Liu, Feng Liu, Kalyan
367 Sunkavalli, Trung Bui, and Hao Tan. LRM: Large reconstruction model for single image to
368 3D. In *International Conference on Learning Representations (ICLR)*, 2024. URL <https://openreview.net/forum?id=s11U8vvsFF>.
369
- 370 [14] Shimin Hu, Yuanyi Wei, Fei Zha, Yudong Guo, and Juyong Zhang. Easy3E: Feed-forward 3D
371 asset editing via rectified voxel flow. *arXiv preprint*, 2026. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2602.21499)
372 [2602.21499](https://arxiv.org/abs/2602.21499). CVPR 2026.

- 373 [15] Zixuan Huang, Mark Boss, Aaryaman Vasishta, James Matthew Rehg, and Varun Jampani.
374 SPAR3D: Stable point-aware reconstruction of 3D objects from single images. In *Conference*
375 *on Computer Vision and Pattern Recognition (CVPR)*, pages 16860–16870, 2025. URL <https://arxiv.org/abs/2501.04689>.
376
- 377 [16] Ajay Jain, Ben Mildenhall, Jonathan T. Barron, Pieter Abbeel, and Ben Poole. Zero-shot
378 text-guided object generation with dream fields. In *Conference on Computer Vision and Pattern*
379 *Recognition (CVPR)*, pages 857–866, 2022. doi: 10.1109/CVPR52688.2022.00094. URL
380 <https://doi.org/10.1109/CVPR52688.2022.00094>.
- 381 [17] Mukul Khanna, Yongsen Mao, Hanxiao Jiang, Sanjay Haresh, Brennan Schacklett, Dhruv Batra,
382 Alexander Clegg, Eric Undersander, Angel X. Chang, and Manolis Savva. Habitat synthetic
383 scenes dataset (HSSD-200): An analysis of 3D scene scale and realism tradeoffs for ObjectGoal
384 navigation. In *IEEE International Conference on Robotics and Automation (ICRA)*, 2024. URL
385 <https://arxiv.org/abs/2306.11290>.
- 386 [18] Juil Koo, Seungwoo Yoo, Minh Hieu Nguyen, and Minhyuk Sung. SALAD: Part-level latent
387 diffusion for 3D shape generation and manipulation. In *International Conference on Computer*
388 *Vision (ICCV)*, 2023. URL <https://arxiv.org/abs/2303.12236>.
- 389 [19] Juil Koo, Wei-Tung Lin, Chanho Park, Chanhyeok Park, and Minhyuk Sung. BoxSplitGen:
390 A generative model for 3D part bounding boxes in varying granularity. In *Winter Conference*
391 *on Applications of Computer Vision (WACV)*, pages 1777–1787, 2026. URL <https://arxiv.org/abs/2602.20666>.
392
- 393 [20] Mingi Lee, Dongsu Zhang, Clément Jambon, and Young Min Kim. BrepDiff: Single-stage
394 b-rep diffusion model. In *Proceedings of the Special Interest Group on Computer Graphics*
395 *and Interactive Techniques Conference Conference Papers, SIGGRAPH Conference Papers '25*,
396 New York, NY, USA, 2025. Association for Computing Machinery. ISBN 9798400715402. doi:
397 10.1145/3721238.3730698. URL <https://doi.org/10.1145/3721238.3730698>.
- 398 [21] Chen-Hsuan Lin, Jun Gao, Luming Tang, Towaki Takikawa, Xiaohui Zeng, Xun Huang, Karsten
399 Kreis, Sanja Fidler, Ming-Yu Liu, and Tsung-Yi Lin. Magic3D: High-resolution text-to-3D
400 content creation. In *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
401 300–309, 2023. doi: 10.1109/CVPR52729.2023.00037. URL <https://doi.org/10.1109/CVPR52729.2023.00037>.
402
- 403 [22] Chong Mou, Xintao Wang, Liangbin Xie, Yanze Wu, Jian Zhang, Zhongang Qi, and Ying Shan.
404 T2I-adapter: Learning adapters to dig out more controllable ability for text-to-image diffusion
405 models. In *AAAI Conference on Artificial Intelligence (AAAI)*, pages 4296–4304, 2024. doi:
406 10.1609/AAAI.V38I5.28226. URL <https://doi.org/10.1609/AAAI.V38I5.28226>.
- 407 [23] Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy V. Vo, Marc Szafraniec, Vasil
408 Khalidov, Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, Mahmoud
409 Assran, Nicolas Ballas, Wojciech Galuba, Russell Howes, Po-Yao Huang, Shang-Wen Li, Ishan
410 Misra, Michael Rabbat, Vasu Sharma, Gabriel Synnaeve, Hu Xu, Hervé Jégou, Julien Mairal,
411 Patrick Labatut, Armand Joulin, and Piotr Bojanowski. DINOv2: Learning robust visual
412 features without supervision. *Transactions on Machine Learning Research (TMLR)*, 2024. URL
413 <https://openreview.net/forum?id=a68SUt6zFt>.
- 414 [24] Ben Poole, Ajay Jain, Jonathan T. Barron, and Ben Mildenhall. DreamFusion: Text-to-3D using
415 2D diffusion. In *International Conference on Learning Representations (ICLR)*, 2023. URL
416 <https://openreview.net/forum?id=FjNys5c7VyY>.
- 417 [25] Etai Sella, Gal Fiebelman, Noam Atia, and Hadar Averbuch-Elor. Spice-e: Structural priors
418 in 3D diffusion using cross-entity attention. In *ACM SIGGRAPH*, pages 1–11, 2024. URL
419 <https://arxiv.org/abs/2311.17834>.
- 420 [26] Stefan Stojanov, Anh Thai, and James M. Rehg. Using shape to categorize: Low-shot learning
421 with an iterative categorization-discrimination loop. In *Conference on Computer Vision and*
422 *Pattern Recognition (CVPR)*, 2021. URL <https://arxiv.org/abs/2104.07371>.
- 423 [27] Jingxiang Sun, Bo Zhang, Ruizhi Shao, Lizhen Wang, Wen Liu, Zhenda Xie, and Yebin Liu.
424 DreamCraft3D: Hierarchical 3D generation with bootstrapped diffusion prior. In *International*
425 *Conference on Learning Representations (ICLR)*, 2024. URL [https://openreview.net/](https://openreview.net/forum?id=DDX1u29Gqr)
426 [forum?id=DDX1u29Gqr](https://openreview.net/forum?id=DDX1u29Gqr).

- 427 [28] Jiaxiang Tang, Jiawei Ren, Hang Zhou, Ziwei Liu, and Gang Zeng. DreamGaussian: Generative
428 gaussian splatting for efficient 3D content creation. In *International Conference on Learning*
429 *Representations (ICLR)*, 2024. URL <https://openreview.net/forum?id=UyNXMqnN3c>.
- 430 [29] Team Hunyuan3D. Hunyuan3D-omni: A unified framework for controllable generation of 3D
431 assets. *arXiv preprint*, 2025. URL <https://arxiv.org/abs/2509.21245>.
- 432 [30] Dmitry Tochilkin, David Pankratz, Zexiang Liu, Zixuan Huang, Adam Letts, Yangguang Li,
433 Ding Liang, Christian Laforte, Varun Jampani, and Yan-Pei Cao. TripoSR: Fast 3D object
434 reconstruction from a single image. *arXiv preprint*, 2024. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2403.02151)
435 [2403.02151](https://arxiv.org/abs/2403.02151).
- 436 [31] Anbang Wang, Yuzhuo Ao, Shangzhe Wu, and Chi-Keung Tang. SK-adapter: Skeleton-based
437 structural control for native 3D generation. *arXiv preprint*, 2026. URL [https://arxiv.org/](https://arxiv.org/abs/2603.14152)
438 [abs/2603.14152](https://arxiv.org/abs/2603.14152).
- 439 [32] Zhenwei Wang, Tengfei Wang, Zexin He, Gerhard Hancke, Ziwei Liu, and Rynson W. H. Lau.
440 Phidias: A generative model for creating 3D content from text, image, and 3D conditions
441 with reference-augmented diffusion. In *International Conference on Learning Representations*
442 *(ICLR)*, 2025. URL [https://proceedings.iclr.cc/paper_files/paper/2025/hash/](https://proceedings.iclr.cc/paper_files/paper/2025/hash/50ca96a1a9ebe0b5e5688a504feb6107-Abstract-Conference.html)
443 [50ca96a1a9ebe0b5e5688a504feb6107-Abstract-Conference.html](https://proceedings.iclr.cc/paper_files/paper/2025/hash/50ca96a1a9ebe0b5e5688a504feb6107-Abstract-Conference.html).
- 444 [33] Shuang Wu, Youtian Lin, Feihu Zhang, Yifei Zeng, Jingxi Xu, Philip Torr, Xun
445 Cao, and Yao Yao. Direct3D: Scalable image-to-3D generation via 3D latent
446 diffusion transformer. In *Advances in Neural Information Processing Systems*
447 *(NeurIPS)*, volume 37, pages 121859–121881, 2024. doi: 10.52202/079017-3873.
448 URL [https://proceedings.neurips.cc/paper_files/paper/2024/hash/](https://proceedings.neurips.cc/paper_files/paper/2024/hash/dc970c91c0a82c6e4cb3c4af7bff5388-Abstract-Conference.html)
449 [dc970c91c0a82c6e4cb3c4af7bff5388-Abstract-Conference.html](https://proceedings.neurips.cc/paper_files/paper/2024/hash/dc970c91c0a82c6e4cb3c4af7bff5388-Abstract-Conference.html).
- 450 [34] Jiatong Xia, Zicheng Duan, Anton van den Hengel, and Lingqiao Liu. Points-to-3D: Structure-
451 aware 3D generation with point cloud priors. In *Conference on Computer Vision and Pattern*
452 *Recognition (CVPR)*, 2026. URL <https://jiatongxia.github.io/points2-3D/>.
- 453 [35] Jianfeng Xiang, Xiaoxue Chen, Sicheng Xu, Ruicheng Wang, Zelong Lv, Yu Deng, Hongyuan
454 Zhu, Yue Dong, Hao Zhao, Nicholas Jing Yuan, and Jiaolong Yang. Native and compact
455 structured latents for 3D generation. *arXiv preprint*, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2512.14692)
456 [2512.14692](https://arxiv.org/abs/2512.14692).
- 457 [36] Jianfeng Xiang, Zelong Lv, Sicheng Xu, Yu Deng, Ruicheng Wang, Bowen Zhang, Dong Chen,
458 Xin Tong, and Jiaolong Yang. Structured 3D latents for scalable and versatile 3D generation. In
459 *Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 21469–21480, 2025.
460 URL <https://arxiv.org/abs/2412.01506>.
- 461 [37] Jiale Xu, Weihao Cheng, Yiming Gao, Xintao Wang, Shenghua Gao, and Ying Shan. In-
462 stantMesh: Efficient 3D mesh generation from a single image with sparse-view large recon-
463 struction models. *arXiv preprint*, 2024. URL <https://arxiv.org/abs/2404.07191>.
- 464 [38] Xiang Xu, Joseph G. Lambourne, Pradeep Kumar Jayaraman, Zhengqing Wang, Karl D.D.
465 Willis, and Yasutaka Furukawa. BrepGen: A b-rep generative diffusion model with structured
466 latent geometry. *ACM Transactions on Graphics (TOG)*, 43(4):1–14, 2024. doi: 10.1145/
467 [3658129](https://brepgen.github.io/). URL <https://brepgen.github.io/>.
- 468 [39] Yunhan Yang, Yufan Zhou, Yuan-Chen Guo, Zi-Xin Zou, Yukun Huang, Ying-Tian Liu, Hao Xu,
469 Ding Liang, Yan-Pei Cao, and Xihui Liu. OmniPart: Part-aware 3D generation with semantic
470 decoupling and structural cohesion. *arXiv preprint*, 2025. URL [https://arxiv.org/abs/](https://arxiv.org/abs/2507.06165)
471 [2507.06165](https://arxiv.org/abs/2507.06165). SIGGRAPH Asia 2025.
- 472 [40] Hu Ye, Jun Zhang, Sibio Liu, Xiao Han, and Wei Yang. IP-adapter: Text compatible image
473 prompt adapter for text-to-image diffusion models. *arXiv preprint*, 2023. URL [https://](https://arxiv.org/abs/2308.06721)
474 arxiv.org/abs/2308.06721.
- 475 [41] Lvmin Zhang, Anyi Rao, and Maneesh Agrawala. Adding conditional control to text-to-image
476 diffusion models. In *International Conference on Computer Vision (ICCV)*, 2023. URL
477 <https://arxiv.org/abs/2302.05543>.
- 478 [42] Wang Zhao, Yan-Pei Cao, Jiale Xu, Yuejiang Dong, and Ying Shan. Assembler: Scalable 3D
479 part assembly via anchor point diffusion. *arXiv preprint*, 2025. URL [https://arxiv.org/](https://arxiv.org/abs/2506.17074)
480 [abs/2506.17074](https://arxiv.org/abs/2506.17074). SIGGRAPH Asia 2025.

- 481 [43] Zibo Zhao, Zeqiang Lai, Qingxiang Lin, Yunfei Zhao, Haolin Liu, Shuhui Yang, Yifei Feng,
482 Mingxin Yang, Sheng Zhang, Xianghui Yang, Huiwen Shi, Sicong Liu, Junta Wu, Yihang Lian,
483 Fan Yang, Ruining Tang, Zebin He, Xinzhou Wang, Jian Liu, Xuhui Zuo, Zhuo Chen, Biwen
484 Lei, Haohan Weng, Jing Xu, Yiling Zhu, Xinhai Liu, Lixin Xu, Changrong Hu, Shaoxiong
485 Yang, Song Zhang, Yang Liu, Tianyu Huang, Lifu Wang, Jihong Zhang, Meng Chen, Liang
486 Dong, Yiwen Jia, Yulin Cai, Jiaao Yu, Yixuan Tang, Hao Zhang, Zheng Ye, Peng He, Runzhou
487 Wu, Chao Zhang, Yonghao Tan, Jie Xiao, Yangyu Tao, Jianchen Zhu, Jinbao Xue, Kai Liu,
488 Chongqing Zhao, Xinming Wu, Zhichao Hu, Lei Qin, Jianbing Peng, Zhan Li, Minghui Chen,
489 Xipeng Zhang, Lin Niu, Paige Wang, Yingkai Wang, Haozhao Kuang, Zhongyi Fan, Xu Zheng,
490 Weihao Zhuang, YingPing He, Tian Liu, Yong Yang, Di Wang, Yuhong Liu, Jie Jiang, Jingwei
491 Huang, and Chunchao Guo. Hunyuan3D 2.0: Scaling diffusion models for high resolution
492 textured 3D assets generation. *arXiv preprint*, 2025. URL <https://arxiv.org/abs/2501.12202>.
493
- 494 [44] Zhe Zhu, Le Wan, Rui Xu, Yiheng Zhang, Honghua Chen, Zhiyang Dou, Cheng Lin, Yuan
495 Liu, and Mingqiang Wei. PartSAM: A scalable promptable part segmentation model trained on
496 native 3D data. *arXiv preprint*, 2026. URL <https://arxiv.org/abs/2509.21965>. ICLR
497 2026.

498 A Supplementary Material

499 This supplement adds the details that support the experimental claims but would interrupt the flow of
500 the main paper. Section A.1 covers data creation, typed constraint synthesis, PartSAM preprocessing,
501 and semantic annotation generation. Section A.2 adds qualitative coverage beyond the main figures.
502 Section A.3 summarizes the secondary Arbor variants and the internal design choices that did not
503 become the final method. Section A.4 records the paper configuration. Section A.5 defines the
504 evaluator protocol and the metrics used in the review.

505 A.1 Data Creation

506 As mentioned in Sec. 4, Arbor is trained on about 50k objects from ABO [4], HSSD [17], and
507 Objaverse XL [5]. This is a large corpus, but still much smaller than the pretraining scale of
508 TRELIS [36]. During development we also saw similar behavior when training only on ABO
509 and HSSD, which reduces the corpus to roughly 10k objects. In this section we describe constraint
510 creation, evaluation data, part segmentation, and semantic extraction. Table 4 gives a compact
511 overview of the datasets and their role in training.

Table 4: **Dataset Overview:** Distribution of data used and their corresponding role in training.

	ABO [4]	HSSD [17]	Objaverse XL [5]	Toys4K [26]
512 Role	train	train	train	benchmark / eval
Segmented count	4,484	6,660	39,846	3,225

513 **Constraint creation.** While the main paper describes constraints as artist authored meshes, there is
514 a practical gap: there is no dataset of hull, avoidance, and touch meshes that can be used directly for
515 training. We therefore spent a large part of the project on an automatic constraint creation system. We
516 call this system the orchestra of constraints. It is used online during data loading, so every training
517 batch can sample fresh control geometry. For evaluation, we freeze the exported constraints and store
518 them in benchmark manifests.

519 The system is organized into three families shown in Fig. 6: hull in green, avoidance in red, and touch
520 in yellow. Each family contains several fast geometry samplers. For hull constraints, one important
521 option is Part Union Hull. Here we use the offline PartSAM [44] segmentation, pick one or more
522 semantic parts, and merge them into a positive region. This gives meaningful hulls such as chair seats,
523 back rests, or lamp heads. Another option is a random surface patch, where we sample connected
524 triangles directly on the mesh and convert them into a positive control region. This produces less
525 semantic but more diverse support geometry. We also use section based hulls and simple proxy
526 shapes when we want broad geometric coverage.

527 Avoidance constraints are sampled differently because they describe free space. One common case is
528 a layout blocker, where we place a forbidden region in a location that should stay empty, for example
529 above a seat or inside the opening of a shelf. Another case is a surface clearance region, where we
530 offset a region away from the mesh and ask the model not to generate into that space. These negative
531 constraints are important because they force Arbor to reason not only about where geometry should
532 exist, but also where it should not.

533 Touch constraints couple a small support region with a forbidden half space behind it. The support
534 region marks where contact should happen, while the forbidden side prevents the model from simply
535 filling both sides of the plane. This is useful for object placement such as feet on the ground or
536 attachment points on walls and support surfaces.

537 The most useful training rows are often combinations. For example, we may place an avoidance
538 region above a hull region, or sample a touch region together with a nearby hull patch. This is
539 what turns the constraint object into a more realistic authoring signal instead of a single isolated
540 mask. Figure 6 shows the full family set. The exact sampling balance of the paper model is given in
541 Appendix A.4.

542 **Evaluation data.** As mentioned in Sec. 4, we evaluate on Toys4K [26], the same object set used
543 by TRELIS [36]. This gives a fair base corpus with broad category coverage and existing prompts.



Figure 6: **Automatic constraint families used by Arbor.** The figure shows the concrete generators that make up Arbor’s typed constraint program. Green columns are positive hull families, yellow columns are touch/contact families, and red columns are avoidance families. Each column lists the family intent at the top and example outputs on several objects below. These families are sampled online during training, while benchmark manifests freeze their exported meshes and metadata for evaluation.

544 However, Toys4K does not provide hull, avoidance, or touch constraints. We therefore had to build
 545 our own benchmark.

546 We created two benchmark splits. The first is an automatic split, which samples 128 rows from the
 547 same constraint families that we also use during training. We call this split auto. The second is a
 548 manual split with 32 rows. Here we tried to work with artist intent instead of the training program.
 549 We loaded reference meshes in Blender, removed parts, added support regions, or blocked space in
 550 the way a human author would reason about the object. This is more time consuming, but it is also
 551 more important, because it gives us a cleaner out of distribution benchmark.

552 For the qualitative figures in the main paper we only use the manual split. It is the strongest test of
 553 whether Arbor can follow human intent rather than only replaying the program it saw during training.
 554 In the appendix we also include examples from the automatic split to show broader coverage. We
 555 keep the split sizes fixed because several baselines require more than ten minutes per object, so larger
 556 suites would make the comparison depend more on compute budget than on method behavior. We
 557 plan to release the benchmark manifests and typed control meshes where redistribution is allowed, so
 558 that later work can compare against the same control setup.

559 **Part segmentation.** As mentioned above, we use PartSAM [44] to extract semantically meaningful
 560 object parts. This helps us sample better hull regions and later connect geometry to semantic labels.
 561 Part segmentation is a difficult problem on its own, so we treat it as offline preprocessing. The
 562 practical benefit is that the expensive step only happens once per object.

563 PartSAM is a feed forward method. In our pipeline the mesh is encoded once, then many part masks
 564 are predicted from cached features instead of re encoding the mesh for every part query. Optional

565 graph cut refinement improves difficult boundaries, and the final labels are written back to the original
566 mesh. In practice this takes about 12 seconds per object and is therefore feasible at our scale. We use
567 these segments both for constraint selection and for semantic extraction.

568 **Semantic extraction.** As already mentioned in the limitations and future work of the main paper,
569 we believe that semantic labels attached to hull parts are one of the most promising next steps for
570 constrained control. This is why we also report Arbor Semantics in Appendix A.3. Since there is no
571 public 3D dataset with both part segments and usable part labels, we built this data ourselves.

572 The final production path stays fully offline and uses a local Qwen3-VL image prompting workflow.
573 For each object we keep the raw PartSAM RGB regions, load four whole-object context renders,
574 and select a subset of visible regions that covers most of the segmented area. The first Qwen3-VL
575 prompt receives the captions together with the context views and proposes a compact object-specific
576 vocabulary of candidate part labels, usually about 10–24 short noun phrases with a few aliases.

577 We then issue a second prompt over lettered region cards. Each card shows one raw PartSAM region
578 highlighted in up to two views, while the context renders stay visible on the same page. The prompt
579 forces the VLM to choose exactly one label from the candidate vocabulary or return null, rather than
580 inventing a free-form phrase. This constrained vocabulary step matters in practice: it keeps the label
581 set stable across objects and reduces the generic responses that appear when the model is asked to
582 name parts without a shared candidate list. If a page fails to parse or the confidence stays low, we
583 re-query that region with a smaller single-region prompt instead of accepting the ambiguous answer.

584 The stored annotation for each confident region includes the selected label, short variants, confidence,
585 a brief visual-evidence note, supporting views, and the original region geometry statistics. Weak or
586 ambiguous regions remain unlabeled. We write the result as an offline sidecar JSON together with
587 resumable CSV and diagnostics files, so the semantic branch stays out of the training-time dataloader.

588 These labels are the basis of Arbor Semantics, where labeled queries receive an additional semantic
589 text signal.

590 A.2 Extended Results

591 We provide more qualitative results for Arbor in Fig. 7. While the main paper highlights the manual
592 benchmark, this figure adds more rows from both the manual and automatic Toys4K splits. One
593 can see that Arbor generalizes across different object categories and across different mixes of hull,
594 avoidance, and touch constraints.

595 We also provide more sweep states in the supplementary media. These are easier to judge in motion
596 than in still frames, because one can directly see whether the object moves smoothly with the
597 constraint or starts to collapse.

598 A.3 Arbor Variants

599 In the main paper we chose Arbor as the final method. The variants in this section explain what we
600 tried next, what each change was supposed to solve, and why Arbor remained the strongest overall
601 choice. They are useful because they show which parts of the control problem are still open.

602 **Compliance.** Arbor Compliance starts from a trained Arbor checkpoint and adds explicit decoded
603 constraint losses during finetuning. The idea was simple: if Arbor already knows how to use the
604 constraint, perhaps a small amount of direct pressure on the decoded occupancy could improve
605 difficult cases such as weak hulls or small touch regions.

606 We decode the geo conditioned sparse structure prediction \hat{z}_{geo} and evaluate it against the hull,
607 avoidance, and touch targets. Using the notation from Sec. 3, the added objective is

$$\mathcal{L}_{\text{comp}} = \mathcal{L}_{\text{flow}} + \lambda_{\text{hull}}\mathcal{L}_{\text{hull}} + \lambda_{\text{avoid}}\mathcal{L}_{\text{avoid}} + \lambda_{\text{touch}}\mathcal{L}_{\text{touch}} + \lambda_{\text{margin}}\mathcal{L}_{\text{margin}}.$$

608 The reported run uses weight 0.10 for hull, 0.06 for avoidance, 0.06 for touch support, 0.06 for touch
609 forbidden space, and 0.04 for the margin term against the no geometry counterfactual. The margin
610 term is important because it asks the geo conditioned branch to be better than its own no geometry
611 prediction, rather than only to minimize an absolute occupancy error.

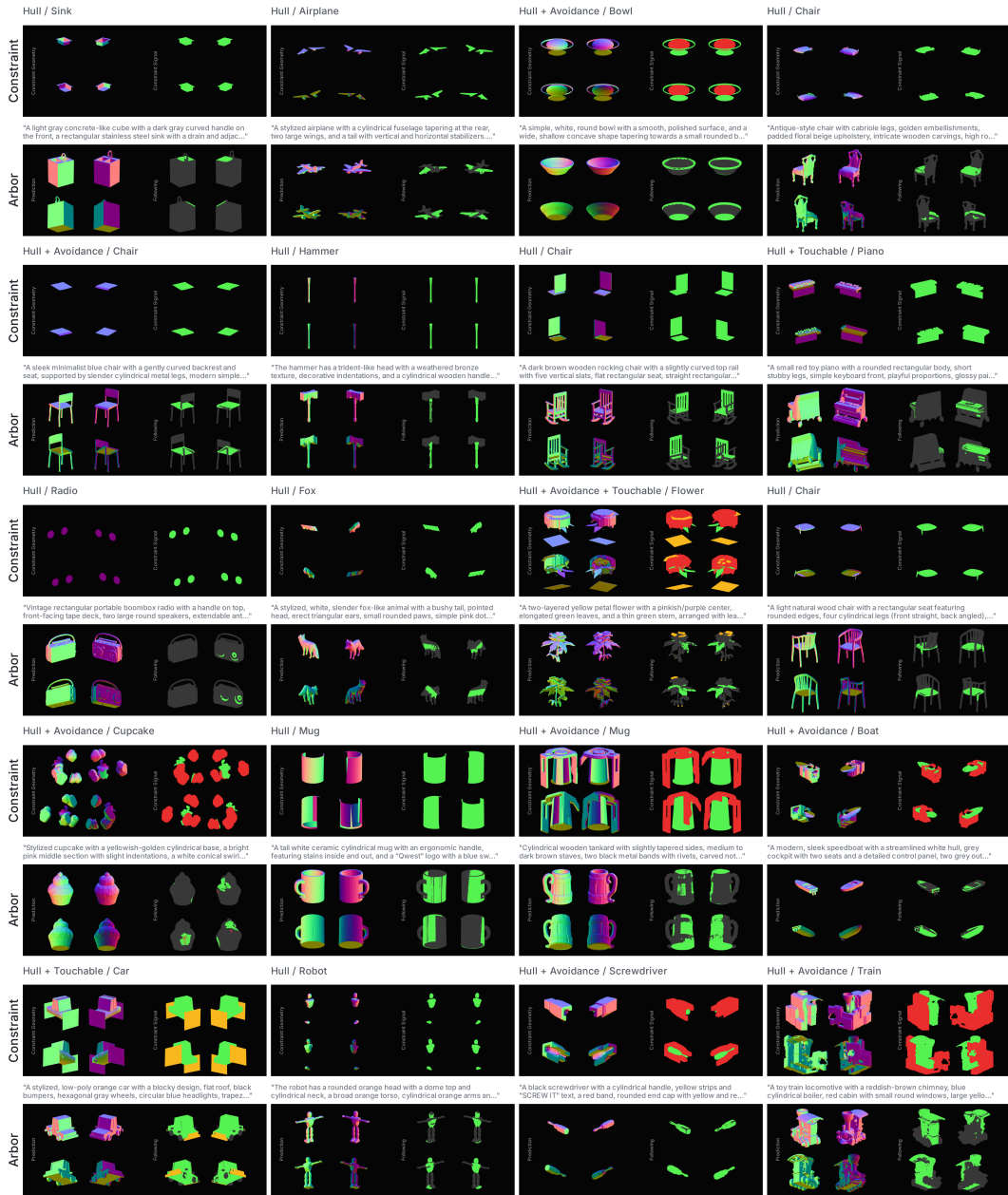


Figure 7: Additional Arbor results on selected Toys4K constraints. Showing manual and automatic benchmark cases.

612 In practice, this variant was not the best solution. It does improve direct constraint pressure, but it
 613 also moves the model toward a failure mode where following the constraint becomes easier than
 614 generating a plausible object from the prompt. This is close to the behavior seen in SpaceControl,
 615 where the guide itself can become the dominant object. From both the experiment section and the
 616 user study we therefore conclude that Compliance is useful as a probe, but too brittle to be the default
 617 method.

618 **Semantics.** Arbor Semantics keeps the geometry path of Arbor and changes only the text path. The
 619 motivation was that semantics should probably not be treated as just another geometry token. Instead,
 620 they should help the model understand what a local region means. For example, it is more useful
 621 to tell the network that a region is the seat of a chair than to append a weak semantic vector to the
 622 routed geometry memory.

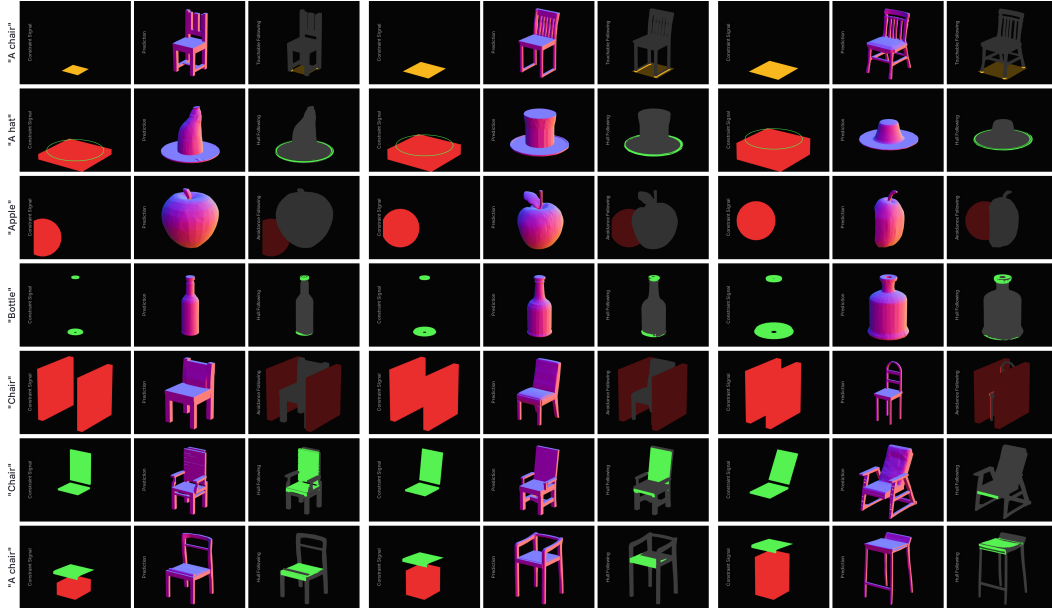


Figure 8: **Extended constraint sweeps.** Additional sweep selections using the same rendering language as Fig. 4, collected in the appendix to show further controlled variations beyond the main paper examples.

623 For a labeled query q with label set S_q and prompt y , we render a semantic first prompt

$$\tilde{y}_q = \text{render}(S_q, y), \quad \tilde{c}_q = E_{\text{text}}(\tilde{y}_q),$$

624 and replace the normal text context $E_{\text{text}}(y)$ with \tilde{c}_q only for those queries. The routed geometry
 625 update from Eq. 5 stays unchanged. Only the text cross attention and its input normalization are
 626 reopened. Self attention, FFN, and time modulation remain frozen.

627 This is a more natural semantic injection site, because it touches the part of the model that already
 628 turns language into structure. We also tested semantic labels inside the routed geometry path, but those
 629 signals were mostly ignored. The text route is more promising. Still, the current Arbor Semantics
 630 model is not yet preferred over Arbor in the user study, and it remains less validated than the main
 631 model. We therefore present it as the clearest future direction, not as the final method.

632 **Gradient baseline.** The Gradient method is a separate idea that never becomes part of Arbor
 633 itself. It asks whether control can be injected only at inference time by modifying the denoising
 634 trajectory, instead of being learned during training. We keep TRELIS frozen, decode the current
 635 sparse structure state, evaluate a hull loss on the occupied target region H , and backpropagate that
 636 loss into the latent:

$$g_k = \nabla_{x_k} \mathcal{L}_{\text{hull}}(x_k), \quad x_k \leftarrow x_k - \eta_k \frac{g_k}{\|g_k\|_2}.$$

637 In the reported setup this guidance is applied only during the last denoising steps.

638 This method does steer the sample toward the constraint, and it often does so better than SpaceControl.
 639 However, it also introduces new artifacts. The model can walk in different directions at once, partially
 640 reconstruct the guide, or copy in the constraint while damaging the rest of the object. This is why the
 641 Gradient row is useful as a test time baseline, but not as the main answer to the control problem.

642 **Variants comparison.** As mentioned in the method section, we also experimented with the place-
 643 ment of the adapter and with deeper interference into the TRELIS block. Table 5 compares the final
 644 Arbor model against the two strongest alternatives that were fully evaluated. Late Adapter moves the
 645 geometry branch later in the block. Fused Attention+FFN merges text and geometry into one main
 646 conditioning site and also reopens the following FFN. Both variants can work, but Arbor remains the
 647 best balanced solution.

648 We also experimented with different encoding structures. One idea was to use O voxels directly and
 649 reduce them with a learned convolutional encoder. This failed badly, because the model first had to

Table 5: **Arbor and its two strongest internal alternatives under the full Toys4K protocol.** Manual ($n=32$) and automatic ($n=128$) splits. Late Adapter moves the geometry branch later in the sparse structure block. Fused Attention+FFN replaces Arbor’s separate geometry branch with one retrained joint conditioning site followed by an unfrozen feed forward block. Values are reported under the same evaluator and control score definition as Tab. 1.

Method	Manual ($n = 32$)						Auto ($n = 128$)					
	Ctrl. Scr.↑	Hull Hit↑	Avoid Viol.↓	Touch Hit↑	Vol. Match↑	MV-CLIP↑	Ctrl. Scr.↑	Hull Hit↑	Avoid Viol.↓	Touch Hit↑	Vol. Match↑	MV-CLIP↑
Arbor	0.4021	0.7143	0.0245	0.8571	0.4685	0.22947	0.4718	0.7860	0.0055	0.9841	0.4866	0.24020
Late Adapter	0.3828	0.6299	0.0257	0.8571	0.4818	0.22933	0.4707	0.7827	0.0049	0.9762	0.4873	0.24019
Fused Attention+FFN	0.3372	0.6099	0.0161	0.5714	0.4578	0.22675	0.4574	0.7497	0.0102	0.9603	0.4695	0.24017

650 learn its own control representation before it could even start to use the signal. We also tried a coarse
651 16^3 control field that matches the TRELIS lattice. This was easy for the network to detect, but too
652 coarse to provide meaningful control at the final 64^3 output scale. In another version the network
653 learned to copy the coarse signal instead of integrating it into the object.

654 Taken together, these experiments explain why the TRELIS.2 encoding [35] is the right choice even
655 if it looks unusual at first. It gives Arbor a strong frozen geometry representation with explicit signal
656 channels. This lets the control path learn how to use geometry, instead of first having to invent the
657 geometry representation itself.

658 A.4 Implementation Details

659 **Paper model.** Arbor is built on TRELIS [36] sparse structure generation from text. The denoiser
660 has 12 blocks, width 768, and operates on a 16^3 latent lattice with 8 channels per site. Constraint
661 meshes are voxelized at 512^3 , encoded by the frozen TRELIS.2 [35] shape and attribute encoders
662 into sparse 32^3 tokens, projected into model width, and capped at 2048 local geometry tokens per
663 query group. The router partitions the sparse-structure lattice into 64 query groups of size $4 \times 4 \times 4$.
664 Each group also receives 96 learned global summary tokens.

665 **Trainable scope.** Only the modules that face geometry are trained in the final Arbor run: the
666 geometry projection, geometry position embedding, routed grounding adapters, global summary
667 modules, and the small semantic part layers that attach confident labels to selected geometry regions.
668 The pretrained TRELIS self attention, text cross attention, and feed forward weights remain frozen.
669 This keeps the base text prior intact and makes Arbor a real conditioning attachment rather than a full
670 backbone finetune.

671 **Training recipe.** The paper model is trained on 8 GPUs with batch size 4 per GPU, AdamW at
672 learning rate 10^{-4} , EMA rate 0.9999, fp16 training, and adaptive gradient clipping. The frozen
673 paper checkpoint was reached after roughly 80 hours of training. Classifier free dropout is applied
674 independently to text and geometry with probabilities 0.1 and 0.1. No explicit compliance loss is used
675 in the final Arbor run. Progression boards use a fixed evaluation subset of 72 samples balanced across
676 ABO, HSSD, and ObjaverseXL Sketchfab, with fixed prompts, fixed constraints, and fixed noise.
677 This recipe was selected because it remained stable across long multi GPU runs without reopening
678 the frozen TRELIS backbone.

679 **Constraint family balance.** Every training sample contains one hull family. Avoidance and
680 touchable families are activated independently with probability 0.5 each. Within the hull sampler,
681 part based hulls dominate the mass with weight 0.6, while random patch hulls, section crops, planar
682 patches, support patches, symmetry anchors, and primitive proxies each carry weight 0.0667. Within
683 avoidance, layout blockers and inverted carvings each carry weight 0.35, surface clearance regions
684 carry 0.2, and coupled keep out patterns carry 0.1. Touchable uses one family but varies the support
685 side, patch shape, and forbidden half space.

686 **Constraint source geometry.** The paper run uses reduced segmented meshes for part based
687 constraint creation, but preserves exact PartSAM identities for semantic joins. This choice matters
688 for both speed and data quality: Arbor keeps the geometry source compact enough for stable training
689 while still letting the semantic annotations refer to the same part identities.

690 **A.5 Evaluation Details**

691 **Evaluator protocol.** All paper benchmarks are frozen manifests. Each row fixes the prompt, typed
 692 constraint meshes, and the benchmark canonical frame. All compared methods are first aligned
 693 to one common evaluation interface and always return a mesh, regardless of their native internal
 694 representation. The evaluator then voxelizes or rerenders those meshes in one shared protocol. This
 695 is what makes the 64^3 control metrics comparable across latent generators, training-free guidance
 696 baselines, and modality-mismatched reconstruction models.

697 **Metric definitions.** All control metrics use a shared 64^3 voxel grid. This matches the sparse
 698 structure stage and gives one common grid for all methods, but very thin contact errors still require
 699 qualitative inspection. Hull Hit is hull support recall, i.e. the fraction of required hull support surface
 700 reached by the prediction. Avoid Viol. is the occupied fraction inside forbidden avoidance volume.
 701 Touch Hit is the hit rate over touch constraints: a touchable region counts as satisfied if the prediction
 702 reaches it anywhere. Volume Match is a bounded occupancy count agreement term,

$$\text{VolMatch} = \frac{\min(|V_{\text{pred}}|, |V_{\text{gt}}|)}{\max(|V_{\text{pred}}|, |V_{\text{gt}}|)},$$

703 where V_{gt} is the source Toys4K asset occupancy used to construct the benchmark row. This is not a
 704 reconstruction metric; it is only a coarse size guard that prevents methods that simply overfill the
 705 guide from looking artificially complete. MV-CLIP is computed from canonical semantic render
 706 views and acts as a coarse prompt semantic check rather than as a realism metric.

707 **Ctrl. Scr.** Ctrl. Scr. is computed per sample as a harmonic mean over the positive terms that apply
 708 to that sample. Let

$$\mathcal{S}_i = \{\text{HullHit}_i, \text{VolMatch}_i, 1 - \text{AvoidViol}_i, \text{MVCLIP}_i\},$$

709 and when touch constraints are present, append TouchHit_i . The sample score is

$$\text{CtrlScr}_i = \frac{|\mathcal{S}_i|}{\sum_{s \in \mathcal{S}_i} \frac{1}{\max(s, \epsilon)}},$$

710 and the table reports the mean over the split. The harmonic mean is deliberate: a method cannot hide
 711 a severe control failure behind one strong metric. It is a compact summary, not a replacement for the
 712 component columns reported in Tab. 1.

713 **Fair comparison policy.** The main controlled generation table is restricted to methods that can
 714 reasonably be evaluated as text plus geometry generation under the same benchmark protocol. Image
 715 conditioned or point conditioned methods are moved to the fixed hull variation track rather than
 716 being mixed into the main control table with extra cues. In that track, the required image or point
 717 set condition is derived from the same benchmark hull rather than giving those methods additional
 718 authored supervision. Internally we also keep track of whether a comparison row is official, surrogate,
 719 modality mismatched, or an internal baseline, so those distinctions stay explicit during evaluation.
 720 We also keep diagnostic GT metrics such as Chamfer and ICP aligned scores in the evaluator, but
 721 we do not surface them in the main paper table because the paper claim is controlled generation, not
 722 reconstruction accuracy.

723 **User study.** We compare the table trends with a small preference study of 404 unlabeled pairwise
 724 choices from 27 participants. Each trial shows the prompt, the constraint render, and candidate outputs
 725 without method names. Participants are asked to choose the preferred result under the combined
 726 criterion of control following and object plausibility. The reported percentages are pairwise win rates
 727 over the trials in which a method appears, not a single multinomial distribution over methods. For
 728 this study only, Arbor, Arbor Semantics, and Arbor Compliance are merged into one Arbor family,
 729 because the question is whether the Arbor conditioning approach is preferred overall, not which small
 730 internal variant wins within that family. That is why the preference percentage in Tab. 1 is shown
 731 both for the base Arbor row and for the merged Arbor family number in parentheses.

732 **Responsible research notes.** The user study is a low risk visual preference task over rendered 3D
733 objects. The interface shows the prompt, the constraint image, and anonymized method outputs, and
734 asks participants to select the result that best balances constraint following with object plausibility.
735 The instruction shown to participants is: choose the output that best follows the shown constraint
736 while remaining a plausible object for the prompt. No paid crowd work dataset was collected for
737 this study. No personal or sensitive participant data is used in the model or benchmark. All training
738 and evaluation assets are drawn from cited 3D datasets and models, and their original sources are
739 credited in the paper. Our new assets are the Arbor method, the typed constraint generation pipeline,
740 the Toys4K control benchmark manifests, and the evaluation code. The code and data package is
741 not part of the initial submission. The planned public release will include source citations, license
742 notes, configuration files, benchmark manifests, and evaluation scripts after review and packaging
743 checks. The main risks are the same as for other 3D asset generators: generated assets may be low
744 quality, misleading, or incompatible with source asset licenses if used without review. We therefore
745 treat Arbor as an authoring aid whose outputs should remain subject to human inspection and license
746 checks before use.

747 **A.6 Extended Limitations**

748 The main paper already states the most important limitations, but three additional points matter for
749 interpretation. First, Arbor still operates only at the sparse structure stage, so later SLAT refinement
750 and decoding can soften very local geometric details even when the coarse support is correct. Second,
751 constraint meshes can conflict sharply with the text prior: very small hulls, weak semantic cues, or
752 aggressive keep out regions can force the model into compromises where either prompt following
753 or local obedience becomes visibly worse. Third, the benchmark remains heterogeneous because
754 nearby methods come from editing, completion, sampling time steering, and image conditioned
755 reconstruction rather than from one clean text plus geometry benchmark family. We therefore keep
756 modality specific methods in the track where their inputs are honest, and we do not claim that Arbor
757 solves every form of 3D control. Semantic labels remain the clearest next step, but the current
758 semantic variants still do not match routed geometry as the main control carrier.

759 **NeurIPS Paper Checklist**

760 **1. Claims**

761 Question: Do the main claims made in the abstract and introduction accurately reflect the
762 paper’s contributions and scope?

763 Answer: [Yes].

764 Justification: The abstract and Sec. 1 state the interface, encoding, routing, and adapter
765 claims. The evaluation in Sec. 4, the user study described in A, and the limitations folded
766 into Sec. 5 together with A.6 state the tested scope and current boundaries.

767 Guidelines:

- 768 • The answer [N/A] means that the abstract and introduction do not include the claims
769 made in the paper.
- 770 • The abstract and/or introduction should clearly state the claims made, including the
771 contributions made in the paper and important assumptions and limitations. A [No] or
772 [N/A] answer to this question will not be perceived well by the reviewers.
- 773 • The claims made should match theoretical and experimental results, and reflect how
774 much the results can be expected to generalize to other settings.
- 775 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
776 are not attained by the paper.

777 **2. Limitations**

778 Question: Does the paper discuss the limitations of the work performed by the authors?

779 Answer: [Yes].

780 Justification: Sec. 5 discusses Arbor’s geometry without roles limitation and its sparse
781 structure stage focus. A.6 extends this with prompt and constraint conflict regimes, the weak
782 semantic extension, and benchmark heterogeneity. A.3 reports the gradient, compliance,
783 and semantic Arbor variants.

784 Guidelines:

- 785 • The answer [N/A] means that the paper has no limitation while the answer [No] means
786 that the paper has limitations, but those are not discussed in the paper.
- 787 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 788 • The paper should point out any strong assumptions and how robust the results are to
789 violations of these assumptions (e.g., independence assumptions, noiseless settings,
790 model well-specification, asymptotic approximations only holding locally). The authors
791 should reflect on how these assumptions might be violated in practice and what the
792 implications would be.
- 793 • The authors should reflect on the scope of the claims made, e.g., if the approach was
794 only tested on a few datasets or with a few runs. In general, empirical results often
795 depend on implicit assumptions, which should be articulated.
- 796 • The authors should reflect on the factors that influence the performance of the approach.
797 For example, a facial recognition algorithm may perform poorly when image resolution
798 is low or images are taken in low lighting. Or a speech-to-text system might not be
799 used reliably to provide closed captions for online lectures because it fails to handle
800 technical jargon.
- 801 • The authors should discuss the computational efficiency of the proposed algorithms
802 and how they scale with dataset size.
- 803 • If applicable, the authors should discuss possible limitations of their approach to
804 address problems of privacy and fairness.
- 805 • While the authors might fear that complete honesty about limitations might be used by
806 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
807 limitations that aren’t acknowledged in the paper. The authors should use their best
808 judgment and recognize that individual actions in favor of transparency play an impor-
809 tant role in developing norms that preserve the integrity of the community. Reviewers
810 will be specifically instructed to not penalize honesty concerning limitations.

811 **3. Theory assumptions and proofs**

812 Question: For each theoretical result, does the paper provide the full set of assumptions and
813 a complete (and correct) proof?

814 Answer: [N/A].

815 Justification: The paper does not present theoretical results or proofs. The equations in
816 Sec. 3 define the architecture and training objective.

817 Guidelines:

- 818 • The answer [N/A] means that the paper does not include theoretical results.
- 819 • All the theorems, formulas, and proofs in the paper should be numbered and cross-
820 referenced.
- 821 • All assumptions should be clearly stated or referenced in the statement of any theorems.
- 822 • The proofs can either appear in the main paper or the supplemental material, but if
823 they appear in the supplemental material, the authors are encouraged to provide a short
824 proof sketch to provide intuition.
- 825 • Inversely, any informal proof provided in the core of the paper should be complemented
826 by formal proofs provided in appendix or supplemental material.
- 827 • Theorems and Lemmas that the proof relies upon should be properly referenced.

828 4. Experimental result reproducibility

829 Question: Does the paper fully disclose all the information needed to reproduce the main ex-
830 perimental results of the paper to the extent that it affects the main claims and/or conclusions
831 of the paper (regardless of whether the code and data are provided or not)?

832 Answer: [Yes].

833 Justification: Sec. 3, Sec. 4, A.1, and A.5 describe the architecture, trainable modules, loss,
834 benchmark splits, baselines, and metrics used for the main claims. Some release details such
835 as public code packaging are handled separately in the code access question below.

836 Guidelines:

- 837 • The answer [N/A] means that the paper does not include experiments.
- 838 • If the paper includes experiments, a [No] answer to this question will not be perceived
839 well by the reviewers: Making the paper reproducible is important, regardless of
840 whether the code and data are provided or not.
- 841 • If the contribution is a dataset and/or model, the authors should describe the steps taken
842 to make their results reproducible or verifiable.
- 843 • Depending on the contribution, reproducibility can be accomplished in various ways.
844 For example, if the contribution is a novel architecture, describing the architecture fully
845 might suffice, or if the contribution is a specific model and empirical evaluation, it may
846 be necessary to either make it possible for others to replicate the model with the same
847 dataset, or provide access to the model. In general, releasing code and data is often
848 one good way to accomplish this, but reproducibility can also be provided via detailed
849 instructions for how to replicate the results, access to a hosted model (e.g., in the case
850 of a large language model), releasing of a model checkpoint, or other means that are
851 appropriate to the research performed.
- 852 • While NeurIPS does not require releasing code, the conference does require all submis-
853 sions to provide some reasonable avenue for reproducibility, which may depend on the
854 nature of the contribution. For example
 - 855 (a) If the contribution is primarily a new algorithm, the paper should make it clear how
856 to reproduce that algorithm.
 - 857 (b) If the contribution is primarily a new model architecture, the paper should describe
858 the architecture clearly and fully.
 - 859 (c) If the contribution is a new model (e.g., a large language model), then there should
860 either be a way to access this model for reproducing the results or a way to reproduce
861 the model (e.g., with an open-source dataset or instructions for how to construct
862 the dataset).
 - 863 (d) We recognize that reproducibility may be tricky in some cases, in which case
864 authors are welcome to describe the particular way they provide for reproducibility.
865 In the case of closed-source models, it may be that access to the model is limited in

866 some way (e.g., to registered users), but it should be possible for other researchers
867 to have some path to reproducing or verifying the results.

868 5. Open access to data and code

869 Question: Does the paper provide open access to the data and code, with sufficient instruc-
870 tions to faithfully reproduce the main experimental results, as described in supplemental
871 material?

872 Answer: [No].

873 Justification: The paper discloses the method, data construction, benchmarks, metrics, and
874 implementation details in Sec. 3, Sec. 4, Appendix A.1, and Appendix A.5, but the code and
875 data package is not part of the initial submission. We plan a public release after review and
876 packaging checks.

877 Guidelines:

- 878 • The answer [N/A] means that paper does not include experiments requiring code.
- 879 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
880 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 881 • While we encourage the release of code and data, we understand that this might not
882 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
883 including code, unless this is central to the contribution (e.g., for a new open-source
884 benchmark).
- 885 • The instructions should contain the exact command and environment needed to run to
886 reproduce the results. See the NeurIPS code and data submission guidelines ([https:
887 //neurips.cc/public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 888 • The authors should provide instructions on data access and preparation, including how
889 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 890 • The authors should provide scripts to reproduce all experimental results for the new
891 proposed method and baselines. If only a subset of experiments are reproducible, they
892 should state which ones are omitted from the script and why.
- 893 • At submission time, to preserve anonymity, the authors should release anonymized
894 versions (if applicable).
- 895 • Providing as much information as possible in supplemental material (appended to the
896 paper) is recommended, but including URLs to data and code is permitted.

897 6. Experimental setting/details

898 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
899 rameters, how they were chosen, type of optimizer) necessary to understand the results?

900 Answer: [Yes].

901 Justification: Sec. 4 specifies the training mixture, benchmark sizes, baselines, and metric
902 protocol. A.1, A.3, and A.5 provide additional data, baseline, and metric details.

903 Guidelines:

- 904 • The answer [N/A] means that the paper does not include experiments.
- 905 • The experimental setting should be presented in the core of the paper to a level of detail
906 that is necessary to appreciate the results and make sense of them.
- 907 • The full details can be provided either with the code, in appendix, or as supplemental
908 material.

909 7. Experiment statistical significance

910 Question: Does the paper report error bars suitably and correctly defined or other appropriate
911 information about the statistical significance of the experiments?

912 Answer: [No].

913 Justification: Tab. 2 reports mean and standard deviation across seeds for the fixed constraint
914 variation track, but the main controlled generation table reports deterministic benchmark
915 means over the frozen manual and automatic splits without confidence intervals as computa-
916 tion is expensive for these methods such as Spice-E and to not clutter the table even more,
917 with the broad analysis we already make in the paper.

918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The authors should answer [Yes] if the results are accompanied by error bars, confidence intervals, or statistical significance tests, at least for the experiments that support the main claims of the paper.
- The factors of variability that the error bars are capturing should be clearly stated (for example, train/test split, initialization, random drawing of some parameter, or overall run with given experimental conditions).
- The method for calculating the error bars should be explained (closed form formula, call to a library function, bootstrap, etc.)
- The assumptions made should be given (e.g., Normally distributed errors).
- It should be clear whether the error bar is the standard deviation or the standard error of the mean.
- It is OK to report 1-sigma error bars, but one should state it. The authors should preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis of Normality of errors is not verified.
- For asymmetric distributions, the authors should be careful not to show in tables or figures symmetric error bars that would yield results that are out of range (e.g., negative error rates).
- If error bars are reported in tables or plots, the authors should explain in the text how they were calculated and reference the corresponding figures or tables in the text.

8. Experiments compute resources

Question: For each experiment, does the paper provide sufficient information on the computer resources (type of compute workers, memory, time of execution) needed to reproduce the experiments?

Answer: [Yes].

Justification: Appendix A.4 reports the main Arbor training setup, including the eight GPU run, batch size, optimizer, precision, and approximate training time. Baseline and ablation compute is smaller in scope and follows the evaluation protocol described in Appendix A.5.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- The paper should indicate the type of compute workers CPU or GPU, internal cluster, or cloud provider, including relevant memory and storage.
- The paper should provide the amount of compute required for each of the individual experimental runs as well as estimate the total compute.
- The paper should disclose whether the full research project required more compute than the experiments reported in the paper (e.g., preliminary or failed experiments that didn't make it into the paper).

9. Code of ethics

Question: Does the research conducted in the paper conform, in every respect, with the NeurIPS Code of Ethics [https://neurips.cc/public/EthicsGuidelines?](https://neurips.cc/public/EthicsGuidelines)

Answer: [Yes].

Justification: The work focuses on controllable generation of object assets, uses cited 3D datasets and models, and reports a low risk preference study over rendered outputs. Appendix A.5 notes that no personal or sensitive participant data is used in the model or benchmark.

Guidelines:

- The answer [N/A] means that the authors have not reviewed the NeurIPS Code of Ethics.
- If the authors answer [No], they should explain the special circumstances that require a deviation from the Code of Ethics.
- The authors should make sure to preserve anonymity (e.g., if there is a special consideration due to laws or regulations in their jurisdiction).

971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
1001
1002
1003
1004
1005
1006
1007
1008
1009
1010
1011
1012
1013
1014
1015
1016
1017
1018
1019
1020
1021
1022
1023
1024

10. Broader impacts

Question: Does the paper discuss both potential positive societal impacts and negative societal impacts of the work performed?

Answer: [Yes].

Justification: Appendix A.5 discusses the positive impact of a more controllable 3D authoring interface and the negative risks of low quality, misleading, or rights infringing 3D assets. It also states that generated outputs should remain subject to human review and source license checks.

Guidelines:

- The answer [N/A] means that there is no societal impact of the work performed.
- If the authors answer [N/A] or [No], they should explain why their work has no societal impact or why the paper does not address societal impact.
- Examples of negative societal impacts include potential malicious or unintended uses (e.g., disinformation, generating fake profiles, surveillance), fairness considerations (e.g., deployment of technologies that could make decisions that unfairly impact specific groups), privacy considerations, and security considerations.
- The conference expects that many papers will be foundational research and not tied to particular applications, let alone deployments. However, if there is a direct path to any negative applications, the authors should point it out. For example, it is legitimate to point out that an improvement in the quality of generative models could be used to generate Deepfakes for disinformation. On the other hand, it is not needed to point out that a generic algorithm for optimizing neural networks could enable people to train models that generate Deepfakes faster.
- The authors should consider possible harms that could arise when the technology is being used as intended and functioning correctly, harms that could arise when the technology is being used as intended but gives incorrect results, and harms following from (intentional or unintentional) misuse of the technology.
- If there are negative societal impacts, the authors could also discuss possible mitigation strategies (e.g., gated release of models, providing defenses in addition to attacks, mechanisms for monitoring misuse, mechanisms to monitor how a system learns from feedback over time, improving the efficiency and accessibility of ML).

11. Safeguards

Question: Does the paper describe safeguards that have been put in place for responsible release of data or models that have a high risk for misuse (e.g., pre-trained language models, image generators, or scraped datasets)?

Answer: [Yes].

Justification: Appendix A.5 states that released assets should include source citations, license notes, documentation, and human review before use. The submission does not release an unchecked scraped dataset dump.

Guidelines:

- The answer [N/A] means that the paper poses no such risks.
- Released models that have a high risk for misuse or dual-use should be released with necessary safeguards to allow for controlled use of the model, for example by requiring that users adhere to usage guidelines or restrictions to access the model or implementing safety filters.
- Datasets that have been scraped from the Internet could pose safety risks. The authors should describe how they avoided releasing unsafe images.
- We recognize that providing effective safeguards is challenging, and many papers do not require this, but we encourage authors to take this into account and make a best faith effort.

12. Licenses for existing assets

Question: Are the creators or original owners of assets (e.g., code, data, models), used in the paper, properly credited and are the license and terms of use explicitly mentioned and properly respected?

1025
1026
1027
1028
1029
1030
1031
1032
1033
1034
1035
1036
1037
1038
1039
1040
1041
1042
1043
1044
1045
1046
1047
1048
1049
1050
1051
1052
1053
1054
1055
1056
1057
1058
1059
1060
1061
1062
1063
1064
1065
1066
1067
1068
1069
1070
1071
1072
1073
1074
1075

Answer: [Yes].

Justification: Existing models, baselines, and datasets are credited through citations in the main paper and supplement. Appendix A.5 states that released assets should include source citations and license notes before use.

Guidelines:

- The answer [N/A] means that the paper does not use existing assets.
- The authors should cite the original paper that produced the code package or dataset.
- The authors should state which version of the asset is used and, if possible, include a URL.
- The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- For scraped data from a particular source (e.g., website), the copyright and terms of service of that source should be provided.
- If assets are released, the license, copyright information, and terms of use in the package should be provided. For popular datasets, `paperswithcode.com/datasets` has curated licenses for some datasets. Their licensing guide can help determine the license of a dataset.
- For existing datasets that are re-packaged, both the original license and the license of the derived asset (if it has changed) should be provided.
- If this information is not available online, the authors are encouraged to reach out to the asset's creators.

13. New assets

Question: Are new assets introduced in the paper well documented and is the documentation provided alongside the assets?

Answer: [No].

Justification: The paper introduces Arbor and two Toys4K control benchmarks, and documents their construction, constraint semantics, evaluation protocol, and metrics in Sec. 3, Sec. 4, Appendix A.1, and Appendix A.5. The assets themselves are not included with the initial submission; documentation will accompany the later public release.

Guidelines:

- The answer [N/A] means that the paper does not release new assets.
- Researchers should communicate the details of the dataset/code/model as part of their submissions via structured templates. This includes details about training, license, limitations, etc.
- The paper should discuss whether and how consent was obtained from people whose asset is used.
- At submission time, remember to anonymize your assets (if applicable). You can either create an anonymized URL or include an anonymized zip file.

14. Crowdsourcing and research with human subjects

Question: For crowdsourcing experiments and research with human subjects, does the paper include the full text of instructions given to participants and screenshots, if applicable, as well as details about compensation (if any)?

Answer: [Yes].

Justification: Appendix A.5 describes the participant task, trial content, and preference criterion. The study is a small preference study over rendered examples rather than a paid crowd work dataset collection.

Guidelines:

- The answer [N/A] means that the paper does not involve crowdsourcing nor research with human subjects.
- Including this information in the supplemental material is fine, but if the main contribution of the paper involves human subjects, then as much detail as possible should be included in the main paper.

1076 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
1077 or other labor should be paid at least the minimum wage in the country of the data
1078 collector.

1079 **15. Institutional review board (IRB) approvals or equivalent for research with human**
1080 **subjects**

1081 Question: Does the paper describe potential risks incurred by study participants, whether
1082 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)
1083 approvals (or an equivalent approval/review based on the requirements of your country or
1084 institution) were obtained?

1085 Answer: [Yes].

1086 Justification: Appendix A.5 describes the study as a low risk visual preference task over
1087 rendered 3D outputs and states that no personal or sensitive participant data is used in the
1088 model or benchmark. The study followed the applicable internal process for such low risk
1089 preference judgments.

1090 Guidelines:

- 1091 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
1092 with human subjects.
- 1093 • Depending on the country in which research is conducted, IRB approval (or equivalent)
1094 may be required for any human subjects research. If you obtained IRB approval, you
1095 should clearly state this in the paper.
- 1096 • We recognize that the procedures for this may vary significantly between institutions
1097 and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the
1098 guidelines for their institution.
- 1099 • For initial submissions, do not include any information that would break anonymity (if
1100 applicable), such as the institution conducting the review.

1101 **16. Declaration of LLM usage**

1102 Question: Does the paper describe the usage of LLMs if it is an important, original, or
1103 non-standard component of the core methods in this research? Note that if the LLM is used
1104 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
1105 scientific rigor, or originality of the research, declaration is not required.

1106 Answer: [Yes].

1107 Justification: Appendix A.1 describes the offline semantic extraction pipeline, where a local
1108 Qwen model proposes candidate part phrases for optional semantic annotations. These
1109 labels support secondary semantic variants and are not the central mechanism behind the
1110 main Arbor result.

1111 Guidelines:

- 1112 • The answer [N/A] means that the core method development in this research does not
1113 involve LLMs as any important, original, or non-standard components.
- 1114 • Please refer to our LLM policy in the NeurIPS handbook for what should or should not
1115 be described.