

# SceneCritic: A Symbolic Evaluator for 3D Indoor Scene Synthesis

Anonymous CVPR submission

Paper ID 20

## Abstract

001 *Scaling spatial intelligence in embodied agents demands*  
002 *environments that capture rich compositional structure and*  
003 *precise spatial relationships and can support such diverse*  
004 *task requirements. Large Language Models (LLMs) and*  
005 *Vision-Language Models (VLMs) increasingly generate in-*  
006 *door scenes via intermediate structures like layouts and*  
007 *scene graphs, yet evaluation still relies on LLM or VLM*  
008 *judges whose scores are sensitive to viewpoint, prompt*  
009 *phrasing, and hallucination. This instability makes it hard*  
010 *to disentangle spatial plausibility from evaluation artifacts.*  
011 *We introduce **SceneCritic**, a symbolic evaluator for floor-*  
012 *plan-level layouts grounded in **SceneOnto**, a structured*  
013 *spatial ontology aggregated from 3D-FRONT, ScanNet, and*  
014 *Visual Genome. SceneCritic jointly verifies semantic, ori-*  
015 *entation, and geometric coherence across object relation-*  
016 *ships, providing object- and relationship-level assessments*  
017 *that pinpoint specific violations. We further propose an it-*  
018 *erative refinement testbed that probes how models revise*  
019 *spatial structure under three critic modalities: a rule-based*  
020 *collision critic, an LLM critic operating on layout text, and*  
021 *a VLM critic operating on rendered observations. Extensive*  
022 *experiments show that (a) SceneCritic aligns substantially*  
023 *better with human judgments than VLM-based evaluators,*  
024 *(b) text-only LLMs can outperform VLMs on semantic lay-*  
025 *out quality, and (c) image-based VLM refinement is most*  
026 *effective for semantic and orientation correction.*

## 027 1. Introduction

028 The generation of 3D indoor environments has become  
029 central to a range of applications, from training embodied  
030 agents that must navigate and manipulate objects in realistic  
031 spaces [11, 20, 30, 37], to virtual reality and robotics simu-  
032 lation. Recent work has demonstrated that Large Language  
033 Models (LLMs) and Vision-Language Models (VLMs) can  
034 serve as powerful priors for this task, leveraging their world  
035 knowledge to produce object layouts that are both diverse  
036 and semantically meaningful [4, 14, 32, 42]. These ap-  
037 proaches generate scenes through explicit structured repre-

038 sentations (e.g., floor plans, scene graphs, or sets of spatial  
039 constraints), often coupled with refinement strategies and  
040 constraint-based optimization to improve physical and se-  
041 mantic plausibility before asset placement and rendering.  
042 Implicitly, this pipeline frames scene generation as a test of  
043 whether a model can construct and maintain a coherent spa-  
044 tial representation that remains consistent across successive  
045 placement and refinement steps.

046 Despite this progress, the evaluation of generated scenes  
047 remains surprisingly underexplored. The most common  
048 protocol is to ask a VLM to judge a small number of ren-  
049 dered views, but such evaluators are unstable: they miss  
050 object relations due to viewpoint and occlusion, are prone  
051 to hallucination, and are highly sensitive to prompt design,  
052 where minor rephrasing can substantially alter scores [23].  
053 As Figure 1 illustrates, a VLM evaluator assigns scores of  
054 75% and 55% to the same scene depending on the rendered  
055 view, highlighting the need for view-independent evalua-  
056 tion. Several works acknowledge these limitations and re-  
057 sort to human user studies [24, 28, 29, 46], which are more  
058 reliable but expensive to scale and still produce only single  
059 scores or coarse rankings without identifying which spatial  
060 constraints were satisfied or violated.

061 As the field of 3D scene generation continues to grow  
062 rapidly, especially in embodied and interactive settings [7,  
063 24, 36, 40], there is a pressing need for a trustworthy and  
064 comprehensive evaluation framework. However, design-  
065 ing such an evaluator is far from straightforward, as spatial  
066 reasoning involves multiple interrelated concepts includ-  
067 ing semantic compatibility, geometric property constraints,  
068 and inter-object relationships [6]. Existing programmatic  
069 metrics have addressed specific dimensions of this prob-  
070 lem but remain narrow in scope. Distributional statistics  
071 such as FID, KL divergence over furniture categories, and  
072 out-of-bound rates [14] capture global similarity but can-  
073 not diagnose individual placements. Physics-oriented con-  
074 straints such as collision avoidance, room-layout IoU, and  
075 reachability [40] measure geometric feasibility while leav-  
076 ing semantic coherence unaddressed. Other systems com-  
077 bine collision-free and in-boundary scores with LLM rat-  
078 ings [8, 44], but still outsource semantic judgment to a

		Spatial Semantics		Object Orientation		Object Overlap		Top View
A piano room with a table, four chairs, a coffee table, two plants, and a piano		VLM	SceneCritic: 93%	VLM	SceneCritic: 75%	VLM	SceneCritic: 100%	
View 1	75%	<ul style="list-style-type: none"> <li>Scale: 77.78%</li> <li>Co-occurrence: 100%</li> <li>Completeness: 100%</li> </ul>	60%	<ul style="list-style-type: none"> <li>Violations: Chair 0 should face table. Chair 2 should face table.</li> <li>Successes: Chair 1 orientation OK. Chair 2 orientation OK.</li> </ul>	100%	<ul style="list-style-type: none"> <li>Successes: Chair 0,1,2,3: No Overlap. Table: No Overlap. Coffee Table: No Overlap. Piano: No Overlap. Plant 1, 2: No Overlap.</li> </ul>		
View 2	55%	<ul style="list-style-type: none"> <li>Violations: Table Scale Issue: height=1.387m (expected 0.308-1.191m)</li> </ul>	75%		100%			
Human Agreement		🤔?	😊	🤔?	😊	😊	😊	VLM: 77% SceneCritic: 89% 🟢

Figure 1. **VLM Evaluation vs. SceneCritic.** Given a piano room with nine objects, the VLM assigns different scores to the same scene depending on which view is rendered, while *SceneCritic* evaluates the layout directly by traversing relational constraints from *SceneOnto* (e.g., incorrect table scales, chairs not facing the table they co-occur with). This produces stable, object-level assessments that identify specific violations alongside successful placements. *SceneCritic* scores also align more closely with human preferences.

079 model-based judge. Among existing frameworks, SceneE-  
 080 val [33] goes furthest by combining text-scene fidelity met-  
 081 rics with physics-based plausibility checks, but does not  
 082 evaluate semantic layout coherence such as whether object  
 083 co-occurrences, scales, and orientations reflect plausible in-  
 084 door configurations. Across this landscape, semantic lay-  
 085 out plausibility is consistently either delegated to LLM and  
 086 VLM judges, or completely unaddressed.

087 In this paper, we introduce *SceneCritic*, a symbolic eval-  
 088 uator for indoor floor-plan layouts. *SceneCritic*'s con-  
 089 straints are grounded in *SceneOnto*, a structured spatial  
 090 ontology dataset we construct by aggregating priors from  
 091 3D-FRONT, ScanNet, and Visual Genome, encoding ob-  
 092 ject co-occurrence statistics, expected scales, and canon-  
 093 ical orientations. While prior work has encoded indoor spa-  
 094 tial knowledge implicitly through learned embeddings [47],  
 095 per-scene graph construction [16], or hand-coded proce-  
 096 dural rules [11], these representations are designed for  
 097 generation or perception, not evaluation. *SceneOnto* in-  
 098 stead constructs explicit relational graphs of conditional  
 099 co-occurrence probabilities, scale distributions, and canon-  
 100 ical orientations per room type, aggregated across multiple  
 101 datasets. This structure is what enables *SceneCritic* to tra-  
 102 verse object relationships during evaluation and localize vi-  
 103 olations to specific object pairs. From these priors, *Scene-*  
 104 *Critic* derives interpretable constraints that score layouts  
 105 along three axes: semantic coherence, orientation correct-  
 106 ness, and overlap. Because all scene generation pipelines  
 107 produce a layout before rendering, *SceneCritic* targets this  
 108 shared intermediate representation directly. By evaluating  
 109 at the object and relationship level, it explicitly identifies  
 110 which objects succeed or violate which constraints.

111 A stable symbolic evaluator also enables a deeper ques-  
 112 tion: *how do models with different post-training objectives*  
 113 *build and revise spatial structure over successive refinement*  
 114 *steps?* We pair *SceneCritic* with an iterative refinement

testbed in which a model repeatedly improves an initial lay-  
 out under one of several critic modalities: (a) a rule-based  
 critic using collision constraints as feedback, (b) an LLM  
 critic operating on the layout as text, or (c) a VLM critic  
 operating on rendered observations. Because *SceneCritic*  
 scores every intermediate layout, the test bed produces re-  
 finement trajectories, making it possible to study which  
 errors different models can correct, where they plateau,  
 and whether fixing one violation introduces new ones else-  
 where.

Our key contributions are: (i) we construct *SceneOnto*,  
 a dataset-grounded ontology for indoor layout evaluation,  
 including object dimensions, co-occurrence statistics, sup-  
 port surfaces, and orientation preferences; (ii) we develop  
*SceneCritic*, a symbolic evaluator for indoor floor-plan lay-  
 outs with interpretable constraints over semantic coherence,  
 orientation correctness, and overlap, and validate it against  
 human judgment, showing it aligns substantially better than  
 VLM-based evaluators; (iii) we introduce an iterative re-  
 finement testbed using heuristic, LLM, and VLM critics to  
 probe how models trained under different post-training ob-  
 jectives build, maintain, and revise spatial structure over  
 multiple correction steps; (iv) we use *SceneCritic* to iden-  
 tify comparative strengths and failure modes across crit-  
 ics and post-training regimes, showing that some text-only  
 LLMs outperform VLMs on semantic layout quality, and  
 that image-based VLM refinement is the most effective  
 critic modality.

## 2. Related Works

### LLM and VLM-Driven 3D Scene Generation.

Recent work has increasingly adopted LLMs and VLMs  
 for 3D indoor scene generation, leveraging their spa-  
 tial reasoning and world knowledge to produce structured  
 layouts from text descriptions. LayoutGPT [14] uses  
 in-context learning to directly predict numerical layouts,

150 SceneCraft [18] generates Blender code via scene graphs  
151 with VLM-based iterative refinement, and I-Design [5]  
152 leverages multi-agent LLM collaboration with scene graph  
153 mechanisms. Holodeck [42] uses GPT-4 for fully auto-  
154 mated environment generation with constraint-based place-  
155 ment, while LLplace [41] and FlairGPT [25] refine LLM-  
156 driven layout generation through supervised fine-tuning and  
157 detailed object descriptions, respectively. More recently,  
158 LayoutVLM [32] combines VLM semantic knowledge with  
159 differentiable optimization for physically plausible layouts,  
160 and Holodeck 2.0 [4] introduces vision-language-guided  
161 generation with interactive editing. A parallel line of work  
162 focuses on improving physical plausibility, PhyScene [40],  
163 integrates physics-based guidance for collision avoidance  
164 and object reachability, PhysGaussian [38] couples physics  
165 simulation with 3D Gaussian Splatting, and OptiScene [43]  
166 applies Direct Preference Optimization to align layouts with  
167 human preferences and physical constraints.

168 **VLM-Based Evaluation and Its Limitations.** Sev-  
169 eral scene generation works [5, 32] adopt VLM-based eval-  
170 uation, asking a model to score rendered views of gen-  
171 erated scenes. However, VLM judges suffer from well-  
172 documented limitations: hallucinations are pervasive across  
173 large vision-language models [3, 19, 26], with studies show-  
174 ing that even GPT-4V fails on basic visual patterns [34]  
175 and that statistical biases and language priors cause sys-  
176 tematic object hallucinations [22]. IR3D-Bench [27] fur-  
177 ther exposes limitations in VLM spatial precision through  
178 active 3D reconstruction. On the programmatic side, Lay-  
179 outGPT [14] reports distributional statistics (FID, KL diver-  
180 gence) and out-of-bound rates, PhyScene [40] introduces  
181 collision and reachability constraints, and frameworks like  
182 AutoLayout [8] and OptiScene [43] combine geometric  
183 scores with LLM ratings. The closest work to ours is  
184 SceneEval [33], which introduces a structured evaluation  
185 framework with both fidelity metrics (object counts, at-  
186 tributes, spatial relationships relative to the input text) and  
187 plausibility metrics (collision, support, navigability, acces-  
188 sibility), grounded in SceneEval-500, a dataset of human-  
189 annotated scene descriptions. However, SceneEval still re-  
190 lies on an LLM for object category matching within its  
191 pipeline, its plausibility metrics are physics-oriented and do  
192 not cover semantic layout coherence such as co-occurrence  
193 plausibility, scale correctness, or canonical orientation, and  
194 it evaluates only final scenes rather than refinement trajec-  
195 tories. *SceneCritic* addresses these gaps by grounding eval-  
196 uation in dataset-derived relational priors and producing per-  
197 object, per-relationship assessments.

### 198 3. *SceneOnto* Construction

199 *SceneCritic* derives its constraints from a structured rela-  
200 tional representation of indoor scenes. We construct *Sce-*  
201 *neOnto* to provide this representation by aggregating spa-

202 tial and semantic priors from three complementary datasets:  
203 3D-FRONT [15], which provides detailed geometric an-  
204 notations for 6,813 professionally designed indoor scenes;  
205 ScanNet [10], which contributes real-world RGB-D re-  
206 constructions of 1,511 indoor environments; and Visual  
207 Genome [21], which supplies object co-occurrence and sup-  
208 port relationships from 61,530 annotated images. From  
209 these sources, we extract four types of relational priors:  
210 object dimensions, support relations, co-occurrence statis-  
211 tics, and orientation relationships. The resulting ontology  
212 is organized as a per-room-type relational graph, where  
213 nodes represent object categories and edges encode condi-  
214 tional co-occurrence probabilities, enabling *SceneCritic*  
215 to traverse object relationships during evaluation. Figure 2 il-  
216 lustrates a sub-graph of *SceneOnto* for three common room  
217 types: bedroom, dining room, and living room.

218 **Dimension Extraction.** We extract object dimensions  
219 (width, height, and depth in meters) from 3D-FRONT and  
220 ScanNet. For 3D-FRONT, we compute bounding boxes  
221 from scene annotations with instance-specific scale param-  
222 eters applied to derive metric estimates. For ScanNet, we  
223 identify per-object mesh vertices using segmentation and  
224 aggregation annotations and compute bounding boxes from  
225 the grouped vertices. For each object category, we record  
226 percentile statistics (p5, p25, median, p75, p95) along  
227 with the mean, standard deviation, and number of obser-  
228 vations. These distributions form the geometric constraints  
229 that *SceneCritic* uses for scale verification. Figure 3 (left)  
230 visualizes the distribution of object dimensions, showing  
231 the expected scale ranges across categories.

232 **Support Relations.** Support refers to the surface on  
233 which an object rests. We extract support information from  
234 3D-FRONT and ScanNet through geometric reasoning: ob-  
235 jects are processed in ascending order of their vertical posi-  
236 tion, and each is evaluated against previously placed objects  
237 whose top surface aligns with the current object’s bottom  
238 within a 5 cm tolerance. Objects touching the ground are  
239 labeled as floor-supported. For Visual Genome, we retain  
240 annotated predicates corresponding to physical support in-  
241 teractions (e.g., *on*, *sitting on*, *standing on*) and apply addi-  
242 tional filtering to reduce annotation noise, including remov-  
243 ing impossible supporting surfaces, rejecting most same-  
244 category supports, enforcing that heavier objects cannot be  
245 supported by lighter ones, and applying a spatial sanity  
246 check based on bounding box positions.

247 **Co-occurrence Statistics.** We compute pairwise object  
248 co-occurrence across all scenes, both globally and per room  
249 type, for all three datasets. For each scene, we extract the  
250 set of unique object categories present and count how many  
251 scenes contain both in the ordered pair  $(a, b)$  with a cer-  
252 tain threshold distance. We derive the conditional proba-  
253 bility  $P(b|a)$ , and normalized pointwise mutual informa-  
254 tion ( $nPMI$ ), defined as  $PMI(a, b)/(-\log P(a, b))$ , which

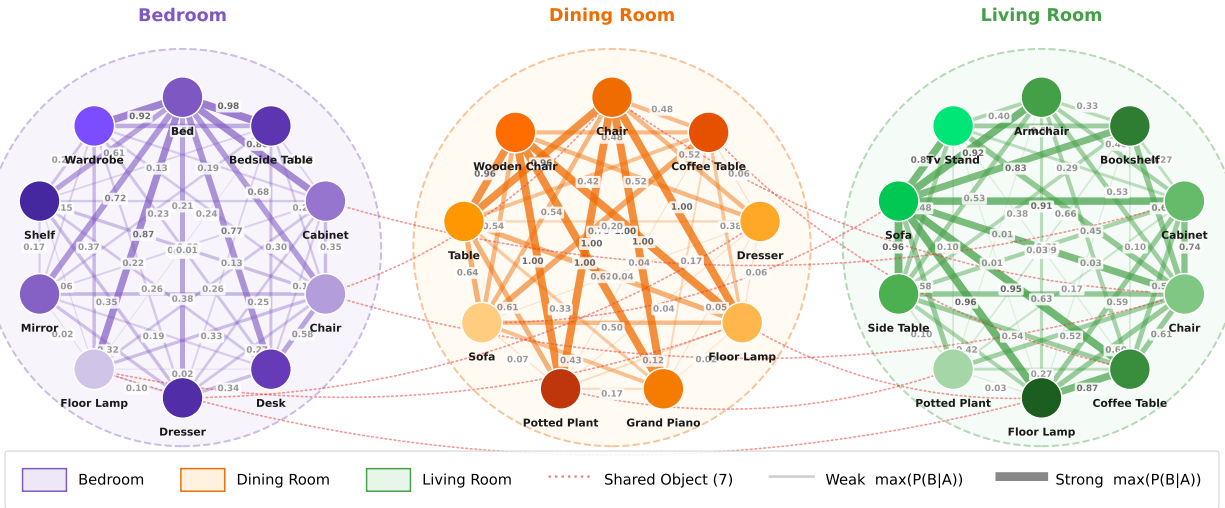


Figure 2. **Relational sub-graph from *SceneOnto* for three common room types.** Nodes represent object categories, with size proportional to frequency. Edges encode conditional co-occurrence probabilities  $\max(P(B|A))$ , with thicker edges indicating stronger associations. Dotted red lines connect shared object categories across room types. *SceneCritic* traverses this structure during evaluation to verify whether object co-occurrences in a generated layout are plausible for the specified room type.

255 ranges over  $[-1, 1]$ . These statistics form the edges in *SceneOnto*'s relational graph and are what *SceneCritic* uses to  
 256 assess co-occurrence plausibility and proximity. Figure 3  
 257 (center) plots co-occurrence frequency against semantic as-  
 258 sociation (nPMI), revealing that 73.6% of the 489 object-  
 259 pair edges have positive association.  
 260

261 **Orientation Relationships.** We extract orientation  
 262 statistics exclusively from 3D-FRONT, which provides per-  
 263 object quaternion rotations and room floor geometry. We  
 264 compute three types of orientation relationships. **Back-to-**  
 265 **wall:** for categories such as sofa, bed, dresser, and book-  
 266 shelf, we compute the angular difference between the ob-  
 267 ject's back direction and the inward normal of the nearest  
 268 wall segment. **Faces-center:** for categories such as sofa,  
 269 chair, and TV, we compute the angular deviation between  
 270 the object's facing direction and the direction toward the  
 271 room centroid. **Faces-pair:** for each object, we compute  
 272 the angular deviation between its facing direction and the  
 273 direction toward every other object within a 5 m radius. The  
 274 facing direction is derived from the object's yaw. For each  
 275 relationship type, the ontology records the fraction of obser-  
 276 vations satisfying the constraint, along with the mean angu-  
 277 lar deviation, mean distance, and number of samples. Fig-  
 278 ure 3 (right) shows the placement strategy landscape, with  
 279 objects distributed along back-to-wall and faces-center frac-  
 280 tions.

#### 281 4. *SceneCritic* as a Symbolic Evaluator

282 *SceneCritic* evaluates generated layouts by traversing  
 283 the relational structure encoded in *SceneOnto*, verifying

284 whether the spatial arrangement of objects is consistent  
 285 with the priors derived from real indoor scenes. Follow-  
 286 ing the evaluation axes adopted by prior scene generation  
 287 work [4, 29, 32], *SceneCritic* assesses layouts along three  
 288 dimensions: (a) spatial semantics, (b) orientation correct-  
 289 ness, and (c) overlap. The key difference is that where prior  
 290 methods rely on VLM judges to score these criteria from  
 291 rendered views, *SceneCritic* resolves them symbolically by  
 292 traversing *SceneOnto*'s constraints at the object and rela-  
 293 tionship level.

294 **Spatial Semantics.** Spatial semantics evaluates whether  
 295 the objects in a scene are consistent with their expected  
 296 semantic relationships. This includes whether objects are  
 297 positioned in ways that align with their functional roles  
 298 and typical real-world usage patterns. Compositionally, it  
 299 evaluates whether multiple objects combine to form coher-  
 300 ent functional groupings through their joint spatial arrange-  
 301 ment, such as a desk and chair forming a workspace or a  
 302 bed, nightstand, and lamp forming a bedside arrangement.  
 303 Our spatial semantics verifier evaluates five sub-criteria:  
 304 *scale*, *co-occurrence*, *plausibility*, *proximity*, and *complete-*  
 305 *ness*.

306 **Orientation Verification.** The orientation verifier  
 307 checks whether each object's orientation aligns with its ex-  
 308 pected placement as defined by *SceneOnto*. It evaluates  
 309 three aspects: whether an object has its back to the near-  
 310 est wall, whether it faces the room center, and whether it  
 311 faces another object it is expected to be oriented toward.  
 312 Each check compares the object's yaw-derived facing di-  
 313 rection against the target direction, with the ontology spec-  
 314 ifying which checks apply to which object categories and

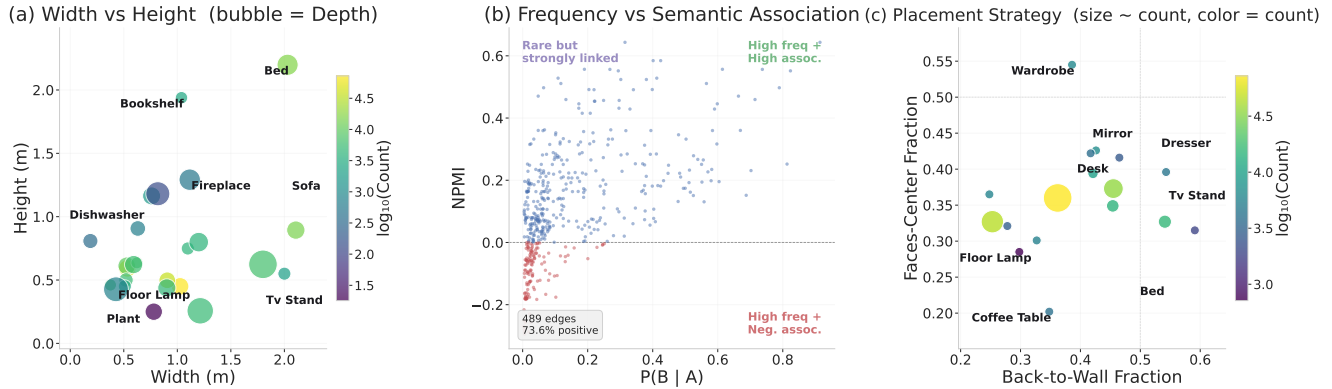


Figure 3. Summary statistics from *SceneOnto*. (Left) Object dimension distributions, with bubble size encoding depth. (Center) Co-occurrence frequency vs. semantic association (nPMI): 73.6% of 489 object-pair edges show positive association. (Right) Object placement strategies, showing the fraction of instances with back-to-wall vs. faces-center orientation (e.g., objects like beds and coffee tables cluster toward back-to-wall, while wardrobes and mirrors favor faces-center).

315 the expected angular tolerances.

316 **Overlap Detection.** The overlap verifier checks whether  
 317 objects violate basic physical constraints. It detects two  
 318 types of violations: Proximity Overlap and True Overlap.  
 319 Proximity Overlap detects when coarse spatial footprints  
 320 of objects encroach on each other, including cases where  
 321 objects are merely in close proximity without true geomet-  
 322 ric contact. True Overlap reports only cases where objects  
 323 genuinely occupy shared space, filtering out false positives  
 324 from loose bounding approximations.

## 325 5. Experiments Setup

326 **Baseline Methods.** We test *SceneCritic* on three represen-  
 327 tative 3D scene layout generation methods that span differ-  
 328 ent generation strategies: LayoutGPT [14], which uses in-  
 329 context learning to predict layouts directly; Holodeck [42],  
 330 which applies constraint-based optimization over LLM-  
 331 generated placements; and LayoutVLM [32], which jointly  
 332 optimizes semantic and physical plausibility via a VLM.  
 333 For each method, we use two generation backbones:  
 334 Gemini-2.5-Flash and Qwen2.5-VL-72B.

335 **Baseline Evaluator.** As the VLM-based baseline eval-  
 336 uator, we select Gemini-2.5-Pro, a stronger proprietary  
 337 model than the generation backbones. We render the gener-  
 338 ated layouts in Blender<sup>1</sup> with Objaverse assets [12] scaled  
 339 to match the object dimensions specified in the layout.  
 340 From each rendered scene, we capture 2D images from mul-  
 341 tiple viewpoints (side, front, rear, and top) and provide them  
 342 to the evaluator along with a prompt describing the evalua-  
 343 tion criteria: semantic consistency, orientation correctness,  
 344 and overlap.

345 **Evaluation Scenes.** We evaluate on two categories of  
 346 rooms. *Pivotal rooms*: bedroom and living room, are the  
 347 standard room types used across all baselines and prior

work, enabling direct comparison. *Extended rooms*: cover  
 348 less common settings including bookstore, buffet restau-  
 349 rant, children’s room, classroom, computer room, deli, din-  
 350 ing room, game room, and florist room, which test general-  
 351 ization beyond the typical evaluation scope. 352

## 353 6. Evaluator Reliability Analysis

### 354 6.1. VLM Instability

355 We first examine the reliability of VLM-based evaluation  
 356 by measuring score variance across viewpoints and repeated  
 357 evaluations using Gemini-2.5-Pro as the evaluator. Table 1  
 358 reports the per-view variance relative to the top-view base-  
 359 line across all generation methods. Score variance is largest  
 360 on the metrics where the generated scene contains the most  
 361 errors. For example, LayoutVLM+Gemini produces fre-  
 362 quent overlap errors, which explains its consistently high  
 363 overlap variance across all views (11.2 Left, 28.1 Right,  
 364 17.5 Front). Additionally, left and right views represent  
 365 mirrored perspectives of the same scene and should receive  
 366 similar scores, yet orientation variance differs substantially  
 367 between them. LayoutVLM shows high variance in the  
 368 right view but low variance in the left, while other meth-  
 369 ods exhibit the opposite pattern. This asymmetry suggests  
 370 the evaluator is reacting to 2D visual appearance rather than  
 371 extracting 3D spatial relationships. Re-evaluating the same  
 372 top-view image under identical conditions produces sub-  
 373 stantial variance: the semantic score varies by 12.46 for  
 374 LayoutGPT+Gemini and 19.32 for LayoutGPT+Qwen72B.  
 375 If an evaluator cannot produce consistent scores for the  
 376 same input, downstream comparisons between generation  
 377 methods become unreliable. We find that method rankings  
 378 reverse depending on which viewpoint is chosen for evalua-  
 379 tion. 380

<sup>1</sup>3D modeling and rendering package: <http://www.blender.org>

Table 1. **Per-view variance relative to the top-view baseline.** Each entry is computed as  $\frac{(\text{View}-\text{Top1})^2}{4}$  for Semantic (Sem), Orientation (Ori), and Overlap (Ovlp). The *Re-eval* column measures variance when the same top-view image is evaluated twice under identical conditions. The remaining columns measure variance between the top-view and alternative viewpoints (Left, Right, Front).

Method	Backbone	<i>Re-eval</i>			Left			Right			Front		
		Sem.	Ori.	Ovlp	Sem.	Ori.	Ovlp	Sem.	Ori.	Ovlp	Sem.	Ori.	Ovlp
Holodeck	Qwen72B	0.11	6.58	3.48	<b>55.06</b>	<b>88.27</b>	1.89	<b>30.25</b>	<b>75.43</b>	1.44	<b>47.68</b>	<b>37.33</b>	0.00
LayoutGPT	Qwen72B	<b>19.32</b>	1.09	0.18	3.65	25.00	0.41	2.39	9.46	4.08	16.61	0.24	<b>46.04</b>
LayoutVLM	Qwen72B	0.49	4.06	0.38	0.44	0.61	1.31	0.05	6.07	2.09	1.90	0.04	0.37
Holodeck	Gemini	0.30	1.97	0.00	4.69	8.56	0.00	0.42	0.09	7.02	1.22	1.97	1.56
LayoutGPT	Gemini	<b>12.46</b>	1.90	1.92	0.12	35.64	<b>11.36</b>	0.58	5.25	10.30	1.66	0.02	2.43
LayoutVLM	Gemini	2.07	4.95	<b>19.27</b>	4.95	1.39	11.16	0.18	9.74	<b>28.09</b>	1.19	0.05	17.47

380

## 6.2. Comparison with *SceneCritic*

381

382

383

384

385

386

387

388

389

390

391

392

393

394

Given these instabilities, we evaluate the same generated layouts using *SceneCritic*, which operates directly on the 3D layout representation without rendered views. Table 2 compares the two evaluators. Interestingly, while VLM evaluation ranks Holodeck+Qwen72B highest (59.46 average), *SceneCritic* ranks LayoutVLM+Qwen72B highest (80.32). Figure 4 provides qualitative results; *SceneCritic* produces per-object, per-constraint assessments that identify the specific violations responsible for each score, such as desks not facing nearby chairs or cabinets not backed against walls. These violations are clearly visible in the 3D layout but are inconsistently detected by the VLM evaluator from 2D renders, which explains the divergence in scores between the two evaluators.

395

## 6.3. Human Alignment Verification

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

To validate which evaluator better reflects human judgment, we conduct a pairwise evaluation study with 16 annotators providing 594 judgments across three criteria (semantic consistency, orientation correctness, and overlap) at two difficulty levels (easy and complex rooms). For reference, prior human evaluations in this space typically use 5 annotators [5, 32, 43], while several prominent methods rely entirely on automated metrics [14, 33, 37].

*SceneCritic* achieves 94.44% agreement with human judgments for easy scenes and 83.33% for complex scenes. The VLM evaluator achieves 58.82% for easy scenes and 47.06% for complex scenes, only marginally above chance. Critically, when *SceneCritic* disagrees with human judgments, the score difference is small (e.g., less than 7 points in the cases where rankings differ). When the VLM evaluator disagrees, the difference is large: Gemini assigns LayoutGPT a 75 points higher score than Holodeck in easy scenes, while human annotators rank Holodeck higher with strong inter-annotator agreement. This pattern is consistent across both difficulty levels. We also observe that the

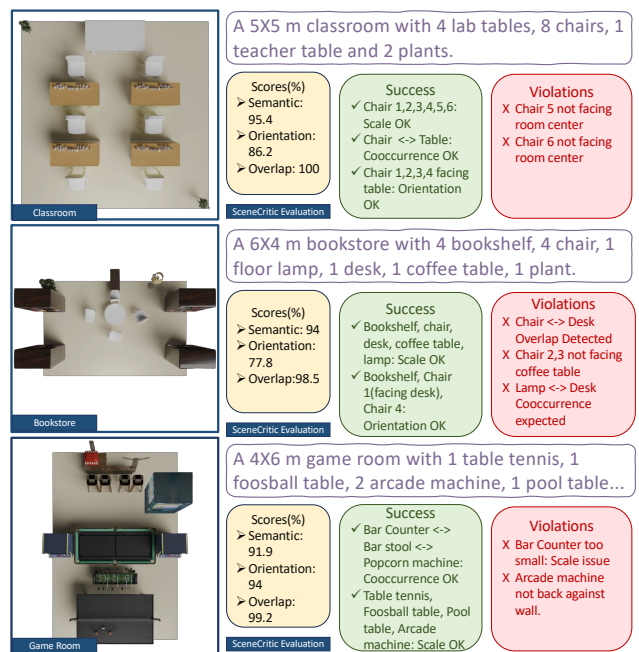


Figure 4. **Qualitative results of *SceneCritic*:** *SceneCritic* assessment for extended scenes, reporting semantic, orientation, and overlap scores while identifying specific violations (e.g., desk not backed against wall, cabinets not backed against wall) alongside successful placements across scale, orientation, co-occurrence, proximity, and completeness.

largest VLM disagreements with humans occur on the overlap metric, whereas *SceneCritic* achieves 100% agreement with human judgments on this criterion across all comparisons. Since overlap is the most visually salient spatial error, the VLM’s inability to consistently detect it from renders confirms that it is not reliably extracting geometric relationships from 2D views.

We also compute inter-annotator agreement to assess the trustworthiness of our human evaluation study. Fleiss’

Table 2. Comparison of Baseline methods across VLM-based evaluator and *SceneCritic* on Semantic (Sem), Orientation (Ori), Overlap (Ovlp) and Average (Avg) scores. In *SceneCritic*, Overlap is decomposed into ProxOvlp (Axis-Aligned Bounding Box overlap) and TrueOvlp (Oriented Bounding Box overlap).

Method	Backbone	VLM-Evaluation				<i>SceneCritic</i>				
		Sem.	Ori.	Ovlp.	Avg.	Sem.	Ori.	Ovlp.		Avg.
								ProxOvlp	TrueOvlp	
Holodeck	Qwen72B	42.61	<b>60.61</b>	<b>75.17</b>	<b>59.46</b>	76.62	52.75	<b>95.97</b>	<b>95.99</b>	75.12
LayoutGPT	Qwen72B	<b>51.09</b>	56.94	69.02	59.02	75.03	16.83	73.54	73.34	55.10
LayoutVLM	Qwen72B	50.82	45.48	71.06	55.79	<b>89.46</b>	<b>59.38</b>	91.92	92.31	<b>80.32</b>
Holodeck	Gemini	46.88	51.33	<b>80.91</b>	<b>59.71</b>	74.37	45.05	<b>97.13</b>	<b>97.21</b>	72.20
LayoutGPT	Gemini	48.18	<b>53.82</b>	77.68	59.89	74.70	37.52	76.99	77.79	63.20
LayoutVLM	Gemini	<b>50.97</b>	50.09	64.70	55.25	<b>89.27</b>	<b>58.65</b>	92.45	91.47	<b>79.96</b>

425 kappa values fall within the moderate range for easy rooms  
 426 (0.30 for semantic, 0.46 for orientation, and 0.39 for over-  
 427 lap) and within the substantial range for complex rooms  
 428 (0.64 for semantic, 0.59 for orientation, and 0.72 for over-  
 429 lap). In addition, the inter-annotator percent agreement is  
 430 81.60% for easy rooms and 91.32% for complex rooms,  
 431 supporting the reliability of our human evaluation.

## 432 7. Probing Spatial Reasoning through Iterative 433 Refinement

434 To study how models build and revise spatial structure under  
 435 iterative feedback, we design a three-stage placement  
 436 pipeline that includes *planning*, *placement*, and *refinement*.

### 437 7.1. Our Test Bed

438 In the **planning stage**, the model receives a placement con-  
 439 dition  $C_{env} = \{C_{desc}, C_{range}, C_{objects}\}$ , where  $C_{desc}$  is a nat-  
 440 ural language scene description,  $C_{range}$  defines the spatial  
 441 boundaries, and  $C_{objects}$  lists the objects to place. The model  
 442 outputs a step-by-step placement plan informed by spatial  
 443 placement principles: hierarchical ordering by object size,  
 444 semantic asset grouping (e.g., placing bed, nightstand, and  
 445 lamp together), and navigable space preservation.

446 In the **placement stage**, the model executes the plan by  
 447 producing a concrete scene layout. The model is prompted  
 448 with geometric constraints, including the spatial boundaries  
 449 of the environment, maintaining realistic object size ratios,  
 450 and assigning meaningful orientations. The output is a lay-  
 451 out specifying each object’s scale-aware bounding box, po-  
 452 sition, orientation, and category label.

453 In the **refinement stage**, the model iteratively improves  
 454 the layout under one of several critic modalities. We evalu-  
 455 ate three critic types: (a) a **heuristic critic** that provides  
 456 feedback based on three constraint objectives, including  
 457 whether objects lie within the spatial boundary, whether all

Table 3. Models categorized by post-training strategy and model parameters.

Model	Category	Post-Training Details	Params
Qwen3-14B [39]	General RL	GRPO-style reinforcement learning	14B
Qwen3-235B [39]	General RL	GRPO-style reinforcement learning	235B
UI-Venus-Navi-72B [17]	General RL	GRPO-based reasoning optimization	72B
Gemini-2.5-flash [9]	RLAIF + RLHF	SFT + Reward Model + RL	N/A
Qwen2.5-VL-7B-MM-RLHF [13]	RLHF	PPO-style human feedback alignment	7B
Qwen2.5-VL-72B-VL [2]	RLHF	SFT + DPO (preference optimization)	72B
LLaMA4-Maverick [1]	RLHF	SFT + Online RL + DPO	~17B (active)
Qwen3-14B-Intuitior-MATH-1EPOCH [45]	RLIF	Iterative feedback RL (Intuitior)	14B
Qwen3-14B-GRPO-MATH-1EPOCH [45]	RLIF	GRPO under RLIF objective	14B
Qwen2.5-14B-GRPO [31]	RLVR	GRPO with verifiable reward	14B
VL-Reasoner-72B [35]	RLVR	GRPO + Self-Selective Revision (SSR)	72B
Holo1.5-72B	RLVR	GRPO-based verifiable reward RL	72B
DeepSeek-3.2V	RLVR	GRPO-style reasoning RL	37B(active)

458 required objects are placed, and whether bounding boxes  
 459 overlap; (b) an **LLM critic** that receives the layout as text  
 460 and provides feedback in natural language; and (c) a **VLM**  
 461 **critic** that receives rendered observations of the scene. We  
 462 further divide VLM-based refinement into three input vari-  
 463 ants: image-only, image+text, and semantics+text. For the  
 464 heuristic critic, the three objectives are combined into a  
 465 single reward score with per-object feedback, allowing the  
 466 model to iteratively update the layout toward a valid con-  
 467 figuration. For LLM and VLM critics, feedback is derived  
 468 from the corresponding model and the layout is modified  
 469 accordingly. In all cases, *SceneCritic* scores every interme-  
 470 diate layout, producing refinement trajectories that reveal  
 471 how each model responds to criticism across successive cor-  
 472 rection steps.

### 473 7.2. Models Evaluated

474 Table 3 summarizes the 13 models evaluated, spanning four  
 475 post-training categories: general RL (GRPO-style training),  
 476 RLHF (PPO or DPO-based preference alignment), RLIF  
 477 (iterative feedback RL), and RLVR (RL with verifiable re-  
 478 wards). Model sizes range from 7B to 235B parameters.

Table 4. Comparison across refinement variants using MLLM backbones with *SceneCritic*. Overlap is decomposed into ProxOvlp (Axis-Aligned Bounding Box overlap) and TrueOvlp (Oriented Bounding Box overlap).

Method	Semantic Verifier Evaluation				Avg.
	Sem.	Ori.	Overlap		
			ProxOvlp	TrueOvlp	
<i>Heuristic</i>					
Gemini-2.5-flash	76.3	43.7	<b>99.3</b>	<b>98.9</b>	73.03
Qwen2.5-72B-VL	72.2	48.1	91.3	91.3	70.53
DeepSeek-3.2V	<b>82.6</b>	64.4	95.3	95.7	<b>80.83</b>
LLaMA4 Maverick	72.6	48.7	95.3	94.1	72.00
Qwen3-235B	75.0	47.1	98.2	97.8	73.37
Qwen3-14B	73.4	44.6	91.4	91.8	69.87
Qwen2.5-VL-7B-MM-RLHF	64.0	47.6	88.3	88.7	66.70
Qwen2.5-14B-GRPO	67.3	46.6	85.2	85.6	66.43
Qwen3-14B-GRPO-MATH-1EPOCH	63.8	<b>65.1</b>	86.2	87.4	71.90
Qwen3-14B-Intuitior-MATH-1EPOCH	65.8	34.3	87.4	87.9	62.58
Holo1.5-72B	68.3	48.7	90.0	90.0	69.00
VL-Reasoner-72B	68.9	42.1	95.6	95.4	68.83
UI-Venus-Navi-72B	66.5	36.6	91.7	91.7	64.93
<i>LLM</i>					
Gemini-2.5-flash	74.2	50.7	<b>98.3</b>	<b>98.4</b>	74.42
Qwen2.5-72B-VL	72.5	48.6	89.5	89.3	70.17
DeepSeek-3.2V	<b>78.9</b>	62.8	97.1	96.9	<b>79.57</b>
LLaMA4 Maverick	69.9	44.7	89.9	89.5	68.10
Qwen3-235B	76.9	54.6	96.5	95.7	75.87
Qwen3-14B	73.3	49.2	93.4	93.4	71.97
Qwen2.5-VL-7B-MM-RLHF	65.5	56.9	90.4	90.6	70.97
Qwen2.5-14B-GRPO	67.9	50.7	86.2	86.6	68.33
Qwen3-14B-GRPO-MATH-1EPOCH	63.8	<b>65.5</b>	68.2	68.9	65.95
Qwen3-14B-Intuitior-MATH-1EPOCH	68.4	43.3	82.6	82.7	64.78
Holo1.5-72B	66.9	42.1	80.9	81.4	63.38
VL-Reasoner-72B	72.8	47.4	91.0	91.0	70.40
UI-Venus-Navi-72B	72.6	45.0	89.6	89.6	69.07
<i>Image</i>					
Gemini-2.5-flash	<b>78.7</b>	<b>55.8</b>	<b>97.5</b>	<b>96.6</b>	<b>77.18</b>
Qwen2.5-72B-VL	71.2	49.1	89.0	88.2	69.63
Holo1.5-72B	68.5	50.2	84.8	85.1	67.88
VL-Reasoner-72B	71.0	44.0	87.2	87.2	67.40
UI-Venus-Navi-72B	72.2	49.1	87.7	87.1	69.57
<i>Img+Text</i>					
Gemini-2.5-flash	<b>76.7</b>	48.4	<b>95.7</b>	<b>96.3</b>	<b>73.70</b>
Qwen2.5-72B-VL	71.0	41.8	89.4	89.0	67.33
Holo1.5-72B	67.9	49.1	83.6	84.3	66.98
VL-Reasoner-72B	72.4	44.1	91.1	90.9	69.17
UI-Venus-Navi-72B	72.0	<b>51.8</b>	90.8	90.6	71.50
<i>Sem+Text</i>					
Gemini-2.5-flash	<b>75.8</b>	<b>53.0</b>	<b>98.6</b>	<b>98.6</b>	<b>75.80</b>
Qwen2.5-72B-VL	72.3	44.6	90.4	89.9	69.02
Holo1.5-72B	67.6	45.8	77.3	77.8	63.65
VL-Reasoner-72B	72.8	50.2	88.2	88.4	70.43
UI-Venus-Navi-72B	70.7	48.6	90.2	90.3	69.85

479

### 7.3. Results

480

481

482

483

484

485

486

487

488

**Text-only LLMs can outperform VLMs on spatial reasoning:** DeepSeek-3.2V achieves the highest average scores under both heuristic refinement (80.83) and LLM-based refinement (79.57), outperforming proprietary VLMs such as Gemini-2.5-Flash. This is driven by strong semantic scores (82.6 heuristic, 78.9 LLM), suggesting that text-based spatial reasoning can be more effective than visual-semantic knowledge for layout generation. **Orientation is the most challenging metric:** Across all models

and refinement methods, orientation scores remain consistently low relative to semantic and overlap scores. Notably, Qwen2.5-14B-GRPO, a relatively small 14B model trained with a math-oriented GRPO objective, achieves the best orientation scores (65.1 heuristic, 65.5 LLM) by a significant margin, outperforming much larger models. We hypothesize this is because orientation verification requires precise angular computation, which may benefit from math-oriented training objectives. **Image-based refinement is the most effective critic modality:** Using Gemini-2.5-Flash as a consistent backbone across all refinement methods, image-based VLM refinement yields the strongest performance on semantic and orientation metrics. Orientation scores increase from 43.70 (heuristic) to 50.7 (LLM) to 55.8 (image-based), an 12-point improvement. A similar pattern holds for UI-Venus-Navi-72B, whose orientation score rises from 36.6 (heuristic) to 49.1 (image-based). These results suggest that visual feedback provides richer spatial information for correcting placement and orientation than text-only inputs. Combining text with images can degrade performance, as the additional textual input may cause the model to prioritize textual cues over spatial information available in the rendered views.

## 8. Conclusion

We introduced *SceneCritic*, a symbolic evaluator for indoor floor-plan layouts, and *SceneOnto*, a structured spatial ontology constructed from 3D-FRONT, ScanNet, and Visual Genome. Together, they provide stable, interpretable, object-level evaluation of semantic coherence, orientation correctness, and overlap, without relying on rendered views or model-based judges. *SceneOnto* is inherently extensible; by providing a simple mapping from a new dataset into our ontology format, novel room types and object categories can be readily integrated, and *SceneCritic* immediately supports these additions. Our experiments show that *SceneCritic* aligns substantially better with human judgments than VLM-based evaluators, particularly on complex scenes where VLM agreement with humans drops to 47.06%. Through our iterative refinement testbed, we find that text-only LLMs can outperform VLMs on semantic layout quality, that math-oriented training objectives improve orientation reasoning even at small model scales, and that image-based VLM refinement is the most effective critic modality for semantic and orientation correction.

## References

- [1] Meta AI. Introducing llama 4: Advancing multimodal intelligence, 2024. 7
- [2] Shuai Bai, Keqin Chen, Xuejing Liu, Jialin Wang, Wenbin Ge, Sibao Song, Kai Dang, Peng Wang, Shijie Wang, Jun Tang, Humen Zhong, Yuanzhi Zhu, Mingkun Yang, Zhaohai Li, Jianqiang Wan, Pengfei Wang, Wei Ding, Zheren Fu,

- 540 Yiheng Xu, Jiabo Ye, Xi Zhang, Tianbao Xie, Zesen Cheng,  
541 Hang Zhang, Zhibo Yang, Haiyang Xu, and Junyang Lin.  
542 Qwen2.5-vl technical report, 2025. 7
- [3] Zechen Bai, Pichao Wang, Tianjun Xiao, Tong He, Zongbo  
543 Han, Zheng Zhang, and Mike Zheng Shou. Hallucination of  
544 multimodal large language models: A survey. *arXiv preprint*  
545 *arXiv:2404.18930*, 2024. 3
- [4] Zixuan Bian, Ruohan Ren, Yue Yang, and Chris Callison-  
546 Burch. Holodeck 2.0: Vision-language-guided 3d world gener-  
547 ation with editing. *arXiv preprint arXiv:2508.05899*, 2025.  
548 1, 3, 4
- [5] Ata Çelen, Guo Han, Konrad Schindler, Luc Van Gool, Iro  
549 Armeni, Anton Obukhov, and Xi Wang. I-design: Personal-  
550 ized llm interior designer. In *European Conference on Com-*  
551 *puter Vision*, 2024. 3, 6
- [6] Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa  
552 Sadigh, Leonidas Guibas, and Fei Xia. Spatialvlm: Endow-  
553 ing vision-language models with spatial reasoning capabili-  
554 ties. In *Proceedings of the IEEE/CVF Conference on Com-*  
555 *puter Vision and Pattern Recognition*, 2024. 1
- [7] ChunTeng Chen, YiChen Hsu, YiWen Liu, WeiFang Sun,  
556 TsaiChing Ni, ChunYi Lee, Min Sun, and YuanFu Yang.  
557 Scenefoundry: Generating interactive infinite 3d worlds.  
558 *arXiv preprint arXiv:2601.05810*, 2026. 1
- [8] Weixing Chen, Dafeng Chi, Yang Liu, Yuxi Yang, Yexin  
559 Zhang, Yuzheng Zhuang, Xingyue Quan, Jianye Hao, Guan-  
560 bin Li, and Liang Lin. Autolayout: Closed-loop layout syn-  
561 thesis via slow-fast collaborative reasoning. *arXiv preprint*  
562 *arXiv:2507.04293*, 2025. 1, 3
- [9] Gheorghie Comanici, Eric Bieber, Mike Schaekermann, Ice  
563 Pasapat, Noveen Sachdeva, Inderjit Dhillon, Marcel Blis-  
564 tein, Ori Ram, Dan Zhang, Evan Rosen, et al. Gemini 2.5:  
565 Pushing the frontier with advanced reasoning, multimodality,  
566 long context, and next generation agentic capabilities. *arXiv*  
567 *preprint arXiv:2507.06261*, 2025. 7
- [10] Angela Dai, Angel X Chang, Manolis Savva, Maciej Hal-  
568 ber, Thomas Funkhouser, and Matthias Nießner. Scannet:  
569 Richly-annotated 3d reconstructions of indoor scenes. In  
570 *Proceedings of the IEEE conference on computer vision and*  
571 *pattern recognition*, 2017. 3
- [11] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs,  
572 Jordi Salvador, Kiana Ehsani, Winson Han, Eric Kolve,  
573 Ali Farhadi, Aniruddha Kembhavi, and Roozbeh Mottaghi.  
574 ProcTHOR: Large-Scale Embodied AI Using Procedural  
575 Generation. In *NeurIPS*, 2022. Outstanding Paper Award.  
576 1, 2
- [12] Matt Deitke, Dustin Schwenk, Jordi Salvador, Luca Weihs,  
577 Oscar Michel, Eli VanderBilt, Ludwig Schmidt, Kiana  
578 Ehsani, Aniruddha Kembhavi, and Ali Farhadi. Objaverse:  
579 A universe of annotated 3d objects. In *Proceedings of*  
580 *the IEEE/CVF conference on computer vision and pattern*  
581 *recognition*, 2023. 5
- [13] Hanze Dong, Wei Xiong, Bo Pang, Haoxiang Wang, Han  
582 Zhao, Yingbo Zhou, Nan Jiang, Doyen Sahoo, Caiming  
583 Xiong, and Tong Zhang. Rlhf workflow: From reward mod-  
584 eling to online rlhf. *arXiv preprint arXiv:2405.07863*, 2024.  
585 7
- [14] Weixi Feng, Wanrong Zhu, Tsu-jui Fu, Varun Jampani, Ar-  
586 jun Akula, Xuehai He, Sugato Basu, Xin Eric Wang, and  
587 William Yang Wang. Layoutgpt: Compositional visual plan-  
588 ning and generation with large language models. *Advances*  
589 *in Neural Information Processing Systems*, 2023. 1, 2, 3, 5,  
590 601 602 603 604 605 606 607
- [15] Huan Fu, Bowen Cai, Lin Gao, Ling-Xiao Zhang, Jiaming  
608 Wang, Cao Li, Qixun Zeng, Chengyue Sun, Rongfei Jia, Bin-  
609 qiang Zhao, et al. 3d-front: 3d furnished rooms with layouts  
610 and semantics. In *Proceedings of the IEEE/CVF Interna-*  
611 *tional Conference on Computer Vision*, 2021. 3
- [16] Qiao Gu, Ali Kuwajerwala, Sacha Morin, Krishna Murthy  
612 Jatavallabhula, Bipasha Sen, Aditya Agarwal, Corban  
613 Rivera, William Paul, Kirsty Ellis, Rama Chellappa, et al.  
614 Conceptgraphs: Open-vocabulary 3d scene graphs for per-  
615 ception and planning. In *2024 IEEE International Confer-*  
616 *ence on Robotics and Automation (ICRA)*, 2024. 2
- [17] Zhangxuan Gu, Zhengwen Zeng, Zhenyu Xu, Xingran Zhou,  
617 Shuheng Shen, Yunfei Liu, Beitong Zhou, Changhua Meng,  
618 Tianyu Xia, Weizhi Chen, et al. Ui-venus technical report:  
619 Building high-performance ui agents with rft. *arXiv preprint*  
620 *arXiv:2508.10833*, 2025. 7
- [18] Ziniu Hu, Ahmet Iscen, Aashi Jain, Thomas Kipf, Yisong  
621 Yue, David A Ross, Cordelia Schmid, and Alireza Fathi.  
622 Scenecraft: An llm agent for synthesizing 3d scenes as  
623 blender code. In *Forty-first International Conference on Ma-*  
624 *chine Learning*, 2024. 3
- [19] Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong,  
625 Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua  
626 Peng, Xiaocheng Feng, Bing Qin, et al. A survey on hal-  
627 lucination in large language models: Principles, taxonomy,  
628 challenges, and open questions. *ACM Transactions on Infor-*  
629 *mation Systems*, 43(2):1–55, 2025. 3
- [20] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt,  
630 Luca Weihs, Alvaro Herrasti, Daniel Gordon, Yuke Zhu, Ab-  
631 hinav Gupta, and Ali Farhadi. AI2-THOR: An Interactive 3D  
632 Environment for Visual AI. *arXiv*, 2017. 1 633
- [21] Ranjay Krishna, Yuke Zhu, Oliver Groth, Justin Johnson,  
634 Kenji Hata, Joshua Kravitz, Stephanie Chen, Yannis Kalan-  
635 tidis, Li-Jia Li, David A Shamma, et al. Visual genome:  
636 Connecting language and vision using crowdsourced dense  
637 image annotations. *International journal of computer vision*,  
638 2017. 3 639
- [22] Sicong Leng, Hang Zhang, Guanzheng Chen, Xin Li, Shijian  
640 Lu, Chunyan Miao, and Lidong Bing. Mitigating object hal-  
641 lucinations in large vision-language models through visual  
642 contrastive decoding. In *Proceedings of the IEEE/CVF Con-*  
643 *ference on Computer Vision and Pattern Recognition*, 2024.  
644 3 645
- [23] Chengzu Li, Caiqi Zhang, Han Zhou, Nigel Collier, Anna  
646 Korhonen, and Ivan Vulić. TopViewRS: Vision-language  
647 models as top-view spatial reasoners. In *Proceedings of*  
648 *the 2024 Conference on Empirical Methods in Natural Lan-*  
649 *guage Processing*, 2024. 1 650
- [24] Lu Ling, Chen-Hsuan Lin, Tsung-Yi Lin, Yifan Ding, Yu  
651 Zeng, Yichen Sheng, Yunhao Ge, Ming-Yu Liu, Aniket Bera,  
652 and Zhaoshuo Li. Scenethesis: A language and vision  
653

- 654 agentic framework for 3d scene generation. *arXiv preprint*  
655 *arXiv:2505.02836*, 2025. 1 712
- 656 [25] Gabrielle Littlefair, Niladri Shekhar Dutt, and Niloy J Mitra. 713  
657 Flairgpt: Repurposing llms for interior designs. In *Computer*  
658 *Graphics Forum*, 2025. 3 714
- 659 [26] Hanchao Liu, Wenyuan Xue, Yifei Chen, Dapeng Chen, Xiu- 715  
660 tian Zhao, Ke Wang, Liping Hou, Rongjun Li, and Wei Peng. 716  
661 A survey on hallucination in large vision-language models. 717  
662 *arXiv preprint arXiv:2402.00253*, 2024. 3 718
- 663 [27] Parker Liu, Chenxin Li, Zhengxin Li, Yipeng Wu, Wuyang 719  
664 Li, Zhiqin Yang, Zhenyuan Zhang, Yunlong Lin, Sirui Han, 720  
665 and Brandon Y Feng. Ir3d-bench: Evaluating vision- 721  
666 language model scene understanding as agentic inverse ren- 722  
667 dering. *arXiv preprint arXiv:2506.23329*, 2025. 3 723
- 668 [28] Xinhang Liu, Yu-Wing Tai, and Chi-Keung Tang. Agentic 3d 724  
669 scene generation with spatially contextualized vlms. *arXiv*  
670 *preprint arXiv:2505.20129*, 2025. 1 725
- 671 [29] Xingjian Ran, Yixuan Li, Linning Xu, Mulin Yu, and 726  
672 Bo Dai. Direct numerical layout generation for 3d in- 727  
673 door scene synthesis via spatial reasoning. *arXiv preprint*  
674 *arXiv:2506.05341*, 2025. 1, 4 728
- 675 [30] Manolis Savva, Abhishek Kadian, Oleksandr Maksymets, 729  
676 Yili Zhao, Erik Wijmans, Bhavana Jain, Julian Straub, Jia 730  
677 Liu, Vladlen Koltun, Jitendra Malik, et al. Habitat: A 731  
678 platform for embodied ai research. In *Proceedings of*  
679 *the IEEE/CVF international conference on computer vision*,  
680 2019. 1 732
- 681 [31] Zhihong Shao, Peiyi Wang, Qihao Zhu, Runxin Xu, Junxiao 733  
682 Song, Xiao Bi, Haowei Zhang, Mingchuan Zhang, YK Li, 734  
683 Yang Wu, et al. Deepseekmath: Pushing the limits of math- 735  
684 ematical reasoning in open language models. *arXiv preprint*  
685 *arXiv:2402.03300*, 2024. 7 736
- 686 [32] Fan-Yun Sun, Weiyu Liu, Siyi Gu, Dylan Lim, Goutam 737  
687 Bhat, Federico Tombari, Manling Li, Nick Haber, and Jiajun 738  
688 Wu. LayoutVLM: Differentiable optimization of 3d layout via 739  
689 vision-language models. In *Proceedings of the Computer Vi-*  
690 *sion and Pattern Recognition Conference*, 2025. 1, 3, 4, 5,  
691 6 740
- 692 [33] Hou In Ivan Tam, Hou In Derek Pun, Austin T. Wang, An- 741  
693 gel X. Chang, and Manolis Savva. SceneEval: Evaluating 742  
694 semantic coherence in text-conditioned 3D indoor scene syn- 743  
695 thesis. 2025. 2, 3, 6 744
- 696 [34] Shengbang Tong, Zhuang Liu, Yuexiang Zhai, Yi Ma, Yann 745  
697 LeCun, and Saining Xie. Eyes wide shut? exploring the 746  
698 visual shortcomings of multimodal llms. In *Proceedings of*  
699 *the IEEE/CVF conference on computer vision and pattern*  
700 *recognition*, 2024. 3 747
- 701 [35] Haozhe Wang, Chao Qu, Zuming Huang, Wei Chu, 748  
702 Fangzhen Lin, and Wenhui Chen. V1-rethinker: Incentivizing 749  
703 self-reflection of vision-language models with reinforcement 750  
704 learning. *arXiv preprint arXiv:2504.08837*, 2025. 7 751
- 705 [36] Xinjie Wang, Liu Liu, Yu Cao, Ruiqi Wu, Wenkang Qin, 752  
706 Dehui Wang, Wei Sui, and Zhizhong Su. Embodiedgen: 753  
707 Towards a generative 3d world engine for embodied intel- 754  
708 ligence. *arXiv preprint arXiv:2506.10600*, 2025. 1 755
- 709 [37] Hongchi Xia, Xuan Li, Zhaoshuo Li, Qianli Ma, Jiashu Xu, 756  
710 Ming-Yu Liu, Yin Cui, Tsung-Yi Lin, Wei-Chiu Ma, Shen- 757  
711 long Wang, et al. Sage: Scalable agentic 3d scene generation 758  
for embodied ai. *arXiv preprint arXiv:2602.10116*, 2026. 1,  
6
- [38] Tianyi Xie, Zeshun Zong, Yuxing Qiu, Xuan Li, Yutao Feng,  
Yin Yang, and Chenfanfu Jiang. Physgaussian: Physics-  
integrated 3d gaussians for generative dynamics. In *Proceed-  
ings of the IEEE/CVF Conference on Computer Vision and  
Pattern Recognition*, 2024. 3
- [39] An Yang, Anfeng Li, Baosong Yang, Beichen Zhang,  
Binyuan Hui, Bo Zheng, Bowen Yu, Chang Gao, Chengen  
Huang, Chenxu Lv, et al. Qwen3 technical report. *arXiv  
preprint arXiv:2505.09388*, 2025. 7
- [40] Yandan Yang, Baoxiong Jia, Peiyuan Zhi, and Siyuan  
Huang. Physcene: Physically interactable 3d scene synthesis  
for embodied ai. In *Proceedings of the IEEE/CVF Confer-  
ence on Computer Vision and Pattern Recognition*, 2024. 1,  
3
- [41] Yixuan Yang, Junru Lu, Zixiang Zhao, Zhen Luo, James JQ  
Yu, Victor Sanchez, and Feng Zheng. Llplace: The 3d in-  
door scene layout generation and editing via large language  
model. *arXiv preprint arXiv:2406.03866*, 2024. 3
- [42] Yue Yang, Fan-Yun Sun, Luca Weihs, Eli VanderBilt, Al-  
varo Herrasti, Winson Han, Jiajun Wu, Nick Haber, Ranjay  
Krishna, Lingjie Liu, et al. Holodeck: Language guided gen-  
eration of 3d embodied ai environments. In *Proceedings of  
the IEEE/CVF Conference on Computer Vision and Pattern  
Recognition*, 2024. 1, 3, 5
- [43] Yixuan Yang, Zhen Luo, Tongsheng Ding, Junru Lu, Mingqi  
Gao, Jinyu Yang, Victor Sanchez, and Feng Zheng. Op-  
tiscene: Llm-driven indoor scene layout generation via  
scaled human-aligned data synthesis and multi-stage prefer-  
ence optimization. *arXiv preprint arXiv:2506.07570*, 2025.  
3, 6
- [44] Yixuan Yang, Zhen Luo, Tongsheng Ding, Junru Lu, Mingqi  
Gao, Jinyu Yang, Victor Sanchez, and Feng Zheng. Op-  
tiscene: Llm-driven indoor scene layout generation via  
scaled human-aligned data synthesis and multi-stage prefer-  
ence optimization, 2025. 1
- [45] Xuandong Zhao, Zhewei Kang, Aosong Feng, Sergey  
Levine, and Dawn Song. Learning to reason without external  
rewards. *arXiv preprint arXiv:2505.19590*, 2025. 7
- [46] Kaizhi Zheng, Ruijian Zha, Zishuo Xu, Jing Gu, Jie Yang,  
and Xin Eric Wang. Constructing a 3d scene from a single  
image. *arXiv preprint arXiv:2505.15765*, 2025. 1
- [47] Yang Zhou, Zachary While, and Evangelos Kalogerakis.  
Scenegrphnet: Neural message passing for 3d indoor scene  
augmentation. In *Proceedings of the IEEE/CVF Interna-  
tional Conference on Computer Vision*, 2019. 2