

MindBlock: Probing Spatial Assembly and Structure in Unified Multimodal Models

Anonymous CVPR submission

Paper ID

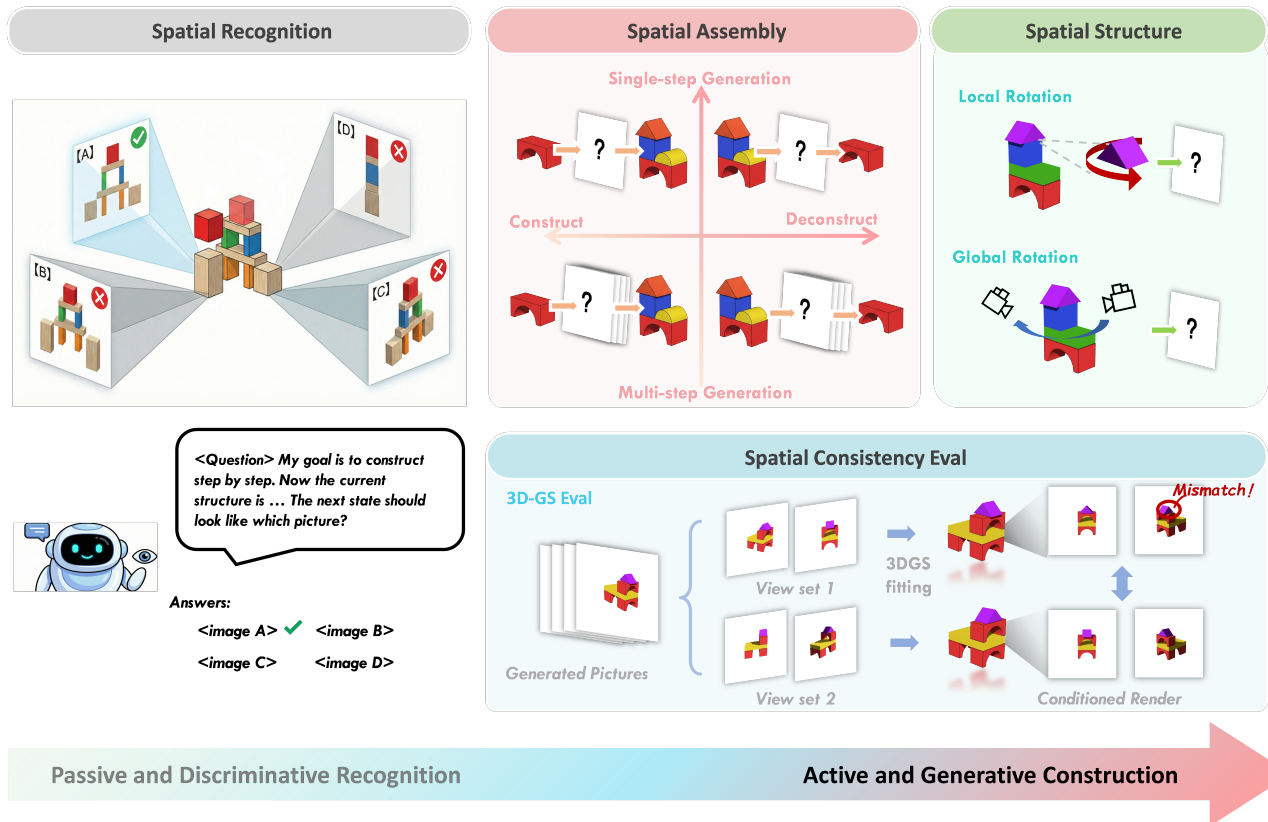


Figure 1. Children aged four to six begin constructing complex block structures from diagrams, demonstrate awareness of enclosure and decoration, and start to grasp three-dimensional spatial relationships. Can AI agents match this level of spatial intelligence? We transform visual-spatial reasoning from **passive selection (multi-choice question answering)** to **active construction (conditional image generation)**, rigorously examining its internal three-dimensional world modeling.

Abstract

001
002
003
004
005
006
007

While Unified Multimodal Models (UMMs) show remarkable reasoning capabilities, their spatial intelligence remains limited to passive 2D Question-Answering (QA). In this paper, we argue true spatial intelligence demands active construction: not only recognize a 3D structure in pixel space, but building and modifying it. We introduce **MindBlock** to challenge models' active generative con-

struction in pixel space across two primary axes: **Spatial Assembly** evaluating step-by-step compositional and causal reasoning, and **Spatial Structure** probing spatial equivariance through local sub-component rotation and global viewpoint transformation. To move beyond pixel-level metrics, we propose **3DGS-Eval**, a novel validation protocol using 3D Gaussian Splatting to reconstruct implicit scenes from model-generated multi-view images. This allows us to quantify structural consistency, verifying for the first

008
009
010
011
012
013
014
015
016

017 *time whether a model’s generative output admits a coherent*
018 *internal 3D world model. Furthermore, we conduct*
019 *a deep-dive diagnostic analysis into the representational*
020 *grounding of spatial logic, disentangling whether structural*
021 *consistency relies on textual Chain-of-Thought (CoT) as a*
022 *symbolic scaffold, or emerges as a native spatial intuition*
023 *within the generative latent space. Our findings reveal a*
024 *significant “perception-execution” gap: while current models*
025 *correctly identify the intended spatial state yet fail in active*
026 *construction. They struggle to maintain spatial equivariance*
027 *without explicit symbolic scaffolding. MindBlock provides a*
028 *rigorous foundation for the next generation of embodied,*
029 *physically-grounded multimodal AI.*

030 *UMMs Spatial Intelligence Dataset and Benchmark*

031 1. Introduction

032 The ability to assemble blocks into meaningful structures
033 unfolds gradually during human development. Toddlers as
034 young as two can stack multiple blocks, but purposeful construction
035 typically emerges around age three. By ages four to five,
036 children begin to replicate simple structures from observation
037 and to grasp spatial relationships such as *next to* and *on top of* [8].
038 The skill of *spatial mental rotation* (imagining how a structure
039 would look from a different viewpoint) develops around ages
040 five to six [4]. This trajectory, from simple stacking to
041 perspective-aware, multi-step assembly, raises a research question:

043 *Can current AI models match the spatial intelligence of a*
044 *six-year-old?*

045 Recent studies have begun to explore this question by
046 evaluating and training foundation models in block-building
047 tasks. BrickGPT [11] models the locations of LEGO bricks
048 using spatial coordinates represented as text, employing large
049 language models (LLMs) to generate physically stable
050 structures. However, this approach operates entirely within
051 text space, whereas humans build blocks using visual input.
052 PhyBlock [10] evaluates 3D block assembly in vision-language
053 models (VLMs) through visual question answering (VQA).
054 Yet multiple-choice question answering can be hacked through
055 statistical shortcuts: correct answers neither guarantee
056 genuine understanding nor translate to robotic assembly.
057 Notably, unified multimodal models (UMMs) [3, 15] have
058 recently emerged to enable both multimodal inputs and
059 outputs, typically in the form of images and text. This
060 makes UMMs uniquely suited for perceiving visual instructions
061 and generating step-by-step image sequences for planning [17].

063 In this work, we move beyond *passive, discriminative*
064 VQA in spatial intelligence: rather than asking models to
065 *select* the correct answer, we ask them to *build* the correct

066 structure. Specifically, we introduce **MindBlock**, a benchmark
067 that probes spatial intelligence through two complementary
068 axes; see Fig. 1.

- 069 • **Spatial Assembly** tests whether models can actively
070 construct block structures step by step. Given an image of a
071 target structure, the model must generate an interleaved
072 sequence of reasoning steps and intermediate images,
073 each showing the scene after one block is placed. This
074 demands *compositional reasoning* (assembling parts into
075 a coherent whole) and *causal reasoning* (a block placed
076 in mid-air will fall; a missing base block invalidates every
077 step above it).
- 078 • **Spatial Structure** tests whether models maintain a
079 consistent 3D understanding across viewpoints. Given the
080 same structure, the model must generate what it looks like
081 from a different angle, preserving every block’s relative
082 position and occlusion. We call this *spatial equivariance*:
083 a model that truly understands a structure in 3D should
084 produce a predictably transformed image when the view-
085 point rotates.

086 Standard image-similarity metrics (*e.g.*, SSIM, LPIPS,
087 CLIP-S) cannot capture a fundamental failure mode of
088 modern Unified Multimodal Models (UMMs): a model
089 may produce images that look plausible individually yet
090 imply mutually inconsistent 3D structures. For instance,
091 blocks may appear on different sides of the structure when
092 rendered from different viewpoints, yielding images that
093 are locally convincing but globally incompatible. To directly
094 measure this phenomenon, we introduce **3DGS-Eval**.
095 Given a set of model-generated multi-view images, we
096 reconstruct a scene using 3D Gaussian Splatting (3DGS),
097 a compact volumetric representation, and evaluate the
098 geometric consistency of the reconstruction. If the
099 generated views encode a coherent 3D structure, the
100 reconstruction converges to a stable geometry; otherwise,
101 contradictory views cause the volumetric fit to
102 deteriorate, revealing hidden spatial inconsistencies that
103 standard 2D metrics fail to detect.

104 Our evaluation experiments reveal that good performance
105 under standard 2D metrics can be misleading. As shown
106 in Table 1, several models achieve high perceptual
107 similarity scores while producing multi-view generations
108 that are geometrically inconsistent. Images that appear
109 plausible in isolation often contradict each other when
110 interpreted as a single 3D structure, exposing a
111 fundamental *2D–3D consistency gap* in current UMMs.

112 To better understand the origin of these failures, we
113 conduct a series of diagnostic analyses. First, we
114 examine whether spatial reasoning benefits from explicit
115 reasoning by comparing generations with and without
116 Chain-of-Thought (CoT) explanations. CoT improves
117 volumetric consistency under 3DGS-Eval, suggesting that
118 language can act as a symbolic scaffold that stabilizes
119 spatial rea-

119 soning during generation. We also extend the evaluation
120 to multi-step construction tasks. While some models suc-
121 ceed in single-step edits, they often fail to maintain con-
122 sistency across longer reasoning trajectories, exhibiting a
123 phenomenon we term *spatial drift*, where small geometric
124 errors accumulate over time.

125 Finally, by comparing generative and discriminative
126 variants of the same tasks, we identify three distinct fail-
127 ure modes: perceptual failures (incorrectly interpreting the
128 input structure), planning failures (producing an incorrect
129 sequence of operations), and execution failures (failing to
130 realize a correct plan in image generation). This analy-
131 sis reveals a pronounced *perception–execution gap*: mod-
132 els frequently recognize correct spatial configurations more
133 reliably than they can construct them.

134 Our contributions are three-fold:

- 135 • **New Tasks.** We evaluate spatial intelligence through *ac-*
136 *tive construction*: models must generate correct visual
137 structures and reasoning trajectories, rather than merely
138 selecting answers. What a model builds provides a
139 stronger probe of its internal 3D understanding than what
140 it recognizes.
- 141 • **New Benchmark.** We introduce **MindBlock**, a spa-
142 tial reasoning benchmark covering six tasks that span
143 atomic operations (single-step assembly and deletion),
144 multi-step structural reasoning (forward assembly and re-
145 verse deconstruction), and viewpoint consistency (local
146 and global equivariance). Evaluation combines conven-
147 tional 2D image similarity with **3DGS-Eval**, a volumetric
148 consistency metric that exposes hidden geometric contra-
149 dictions across views.
- 150 • **New Insights.** Our analysis reveals three systematic lim-
151 itations of current UMMs: (1) a *2D–3D consistency gap*,
152 where high perceptual similarity does not imply correct
153 geometry; (2) *spatial drift* in multi-step reasoning, where
154 structural errors accumulate over time; and (3) a pro-
155 nounced *perception–execution gap*, where models recog-
156 nize correct spatial configurations more reliably than they
157 can generate them. Together, these findings suggest that
158 while UMMs exhibit promising spatial intuition, their in-
159 ternal representations remain fragile and insufficient for
160 robust 3D reasoning.

161 2. MindBlock: Probing Spatial Assembly and 162 Structure

163 The MindBlock benchmark is designed to evaluate two fun-
164 damental dimensions of spatial intelligence: the capacity
165 for *assembly*—the goal-directed construction and decon-
166 struction of physical entities—and the mastery of *struc-*
167 *ture*—the geometric understanding of 3D space.

2.1. Evaluating Spatial Assembly 168

We define spatial assembly as the model’s capacity to exe-
169 cute precise state transitions within a 3D coordinate frame.
170 This is evaluated through two complementary single-step
171 tasks and a long-horizon loop.
172

Single-Step Assembly We focus on the precision of
173 atomic manipulations through two distinct paradigms:
174

- **Atomic Deletion: Targeted Deletion.** This task requires
175 the model to precisely remove a specific target entity (e.g.,
176 “delete only the purple cube”) from a given scene. It
177 serves as a rigorous test for *structural persistence* and
178 *occlusion recovery*: the model must preserve the exact
179 coordinates of all remaining elements and the camera
180 perspective while hallucinating the previously occluded
181 background. This ensures the model treats the structure
182 as a collection of independent 3D entities rather than a
183 monolithic image.
184
- **Atomic Completion: Goal-Directed Completion.** Given
185 a partial assembly and a target state image, the model
186 must autonomously infer and generate the single most
187 critical next step. This evaluates the model’s proactive
188 planning and its ability to bridge the gap between current
189 and desired physical configurations.
190

Multi-Step Assembly-Deconstruction Loop To evalu-
191 ate long-horizon consistency, we extend the aforementioned
192 mappings into a bi-directional loop:
193

- **Iterative Completion: Forward Assembly.** The iterative
194 execution of atomic completions, where the model must
195 progressively build toward a target configuration across
196 multiple steps.
197
- **Iterative Deletion: Reverse Deconstruction.** The contin-
198 uous application of atomic deletions, requiring the model
199 to sequentially dismantle a structure while maintaining
200 spatial coherence at each step.
201

202 Unlike single-step tasks that probe atomic precision, the
203 multi-step loop challenges the model to manage *structural*
204 *evolution*. As the sequence progresses, the model must
205 navigate dynamic occlusion and maintain the integrity of
206 hierarchical support. Successfully deconstructing a com-
207 plex structure requires an awareness of physical causal-
208 ities—such as identifying load-bearing blocks before re-
209 moval—that far exceeds the requirements of single-frame
210 editing, providing a definitive test for the stability of the in-
211 ternal world model.

2.2. Evaluating Spatial Structure 212

To probe the model’s internal 3D world model, we test
213 for *geometric equivariance*—the principle that spatial rep-
214 resentations should transform predictably under geometric
215 operations.
216

Local Equivariance: Local Structure Manipulation:
217 We task the model with rotating a specific sub-component
218 (e.g., “Rotate the top triangular prism 90 degrees clock-
219

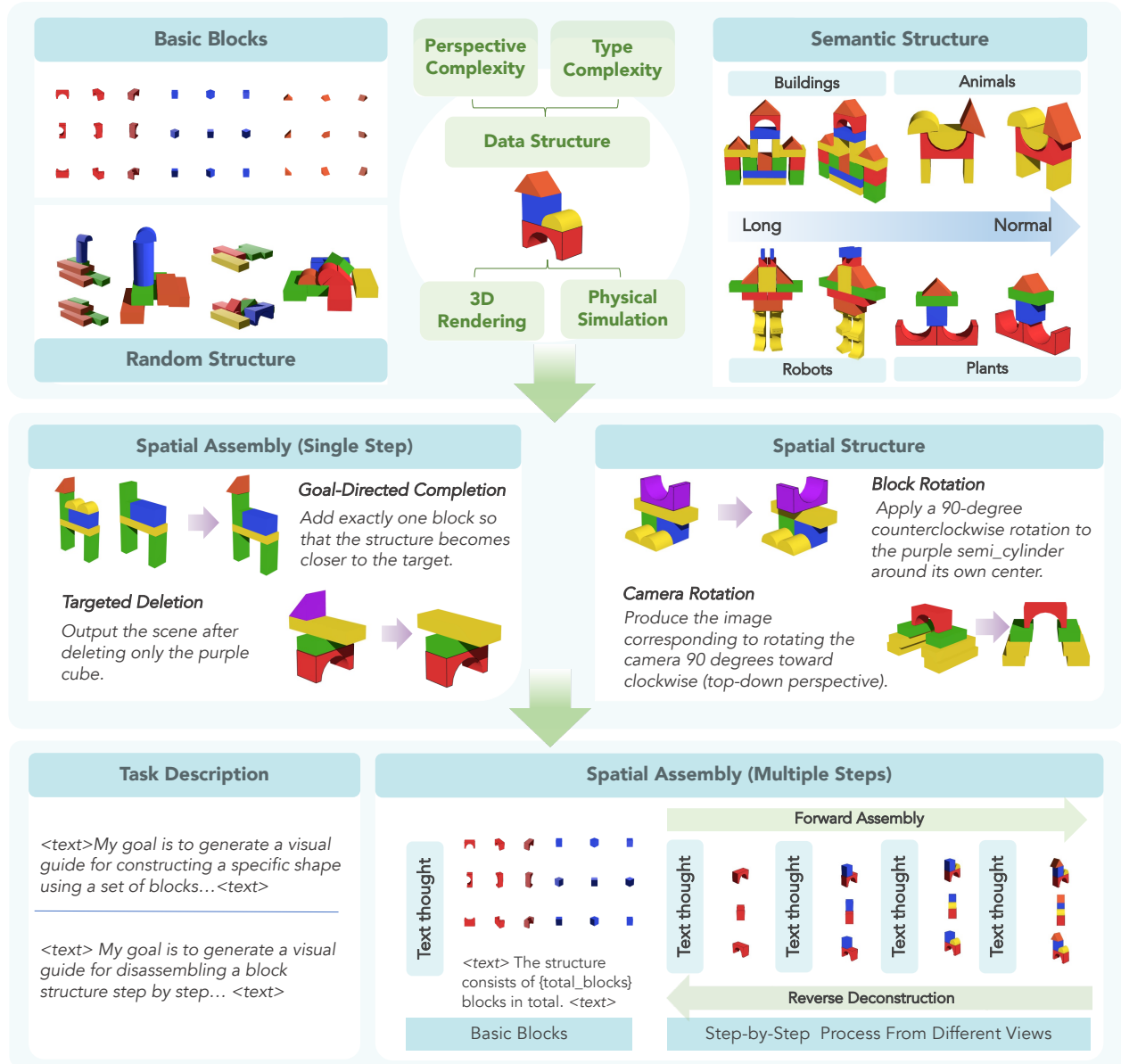


Figure 2. **Overview of the MindBlock benchmark and dataset.** The dataset is constructed from basic blocks via a 3D rendering and physical simulation pipeline (top). The benchmark evaluates spatial intelligence through two primary dimensions: **Spatial Assembly** (left and bottom), covering single-step and long-horizon construction/deconstruction, and **Spatial Structure** (right), testing geometric equivariance through local and global transformations.

220 wise”) while preserving the rest of the structure. This evaluates
 221 the precision of the model’s local coordinate system and its ability to disentangle individual parts within a structured
 222 whole.
 223

224 **Global Equivariance: Global Structure Consistency:**
 225 We assess the model’s ability to generate a globally trans-
 226 formed view (e.g., a 180-degree rotation) of the entire
 227 scene. This serves as a test for global 3D consistency;
 228 a model possessing a true structural understanding must

229 maintain the relative positions and occlusions of all blocks
 230 across arbitrary viewpoint shifts, a capability we quantita-
 231 tively verify using our **3DGS-Eval** protocol.

3. Dataset Construction and Evaluation 232

233 MindBlocks bridges the gap between 2D perception and 3D
 234 execution by grounding structural assembly in a rigorous
 235 physical and geometric protocol.

236	3.1. Physically-Grounded Data Synthesis		
237	We utilize the Genesis physics engine [5] to generate a		
238	large-scale dataset of stable brick structures with precise		
239	control over geometry, lighting conditions, and camera		
240	placement. This simulation-based pipeline enables physi-		
241	cally consistent rendering and fine-grained supervision sig-		
242	nals for spatial reasoning tasks. Our structure library		
243	spans two complementary categories: (i) <i>semantic struc-</i>		
244	<i>tures</i> (e.g., furniture, vehicles, animals), which leverage		
245	common-sense priors about real-world objects; (ii) <i>random</i>		
246	<i>structures</i> , which remove semantic cues to emphasize pure		
247	geometric reasoning and compositional spatial understand-		
248	ing. The overall distribution of categories and task com-		
249	plexity is illustrated in Fig. 3.		
250	The construction of our dataset follows a rigorous		
251	simulation-in-the-loop pipeline. For the random structures,		
252	we employ a constrained random generation process com-		
253	bined with physical simulation to ensure the structural in-		
254	tegrity of each step. For the semantic structures, we uti-		
255	lize a hybrid approach of manual modeling and physical		
256	simulation, encompassing a diverse set of scenes includ-		
257	ing buildings, furniture, and vehicles. A portion of the se-		
258	manitic models is adapted from PhyBlocks [10]. As shown		
259	in the block count distributions (Fig. 3, middle), the com-		
260	plexity of these structures is parameterized by their as-		
261	sembly sequence lengths, spanning from simple 5-block		
262	builds to complex "Long" sequences exceeding 40 blocks,		
263	thereby providing a comprehensive benchmark for evalu-		
264	ating model scalability in spatial reasoning.		
265	For each structure, we generate a bottom-up assembly		
266	sequence that simulates the incremental construction		
267	process. Each construction trajectory consists of multiple		
268	discrete steps, where exactly one block is added at each		
269	step . This sequential formulation exposes the interme-		
270	diate structural states of the scene and provides supervi-		
271	sion for step-by-step reasoning about stability, place-		
272	ment, and support relationships.		
273	At every step of the construction sequence, we render		
274	rich multimodal supervision signals:		
275	• Multi-view Observations: We render images from mul-		
276	tiple viewpoints covering a full 360-degree perspective.		
277	In practice, we sample eight predefined egocentric cam-		
278	era viewpoints around the structure, providing the multi-		
279	angle observations necessary for volumetric understand-		
280	ing and spatial disambiguation.		
281	• Structured Reasoning Traces: We generate template-		
282	synthesized Chain-of-Thought (CoT) reasoning traces.		
283	These traces explicitly describe the intended assembly		
284	rationale, including block attributes, precise (x, y, z) co-		
285	ordinates, and hierarchical support relationships (e.g.,		
286	"Place block A on top of block B").		
	3.2. 3DGS-Eval: Measuring Volumetric Consistency	287	288
	To move beyond pixel-level similarity, we introduce 3DGS-	289	
	Eval , which employs 3D Gaussian Splatting (3D-GS) to	290	
	reconstruct a volumetric field from generated images. Fol-	291	
	lowing the pipeline in MVGBench [16], we initialize 3D	292	
	Gaussians using generated views and their corresponding	293	
	camera poses. Each primitive—parameterized by position,	294	
	opacity, anisotropic covariance, and color—is optimized via	295	
	differentiable rasterization to minimize photometric re-	296	
	construction error against the target images.	297	
	To capture fine geometric details, an adaptive density	298	
	control mechanism iteratively refines the field by split-	299	
	ting, merging, or pruning Gaussians. This process forces the	300	
	representation to converge into a coherent 3D entity. Con-	301	
	sequently, 3DGS-Eval acts as a " geometric polygraph ": if	302	
	the generated views lack underlying 3D consistency, the op-	303	
	timization fails to form a stable structure, leading to struc-	304	
	tural collapse.	305	
	Finally, the optimized 3D representation is re-rendered	306	
	into 2D viewpoints. By comparing these renderings against	307	
	the original generated images using PSNR and SSIM , we	308	
	quantify the volumetric consistency and reconstruction fi-	309	
	delity of the multi-view content.	310	
	4. Experiments	311	
	4.1. Evaluation Setup	312	
	Models. We evaluate a diverse set of models across three	313	
	primary categories to benchmark the current landscape of	314	
	spatial intelligence:	315	
	• General-purpose UMMS: We include state-of-the-	316	
	art autoregressive and diffusion-based unified multi-	317	
	modal models, namely Emu3.5 [2], Ovis-U1 [13], IL-	318	
	LUME [7], BAGEL [3] and Omni-Gen [14], which are	319	
	capable of native image-text generation.	320	
	• Specialized and Editing Models: We assess Blip3-o-	321	
	next [1] and Step1X-Edit [9], which are optimized for	322	
	conditional image manipulation and structural preserva-	323	
	tion.	324	
	Metrics. We employ a dual-pronged evaluation strategy:	325	
	• Image Similarity (2D): We evaluate semantic alignment	326	
	between generated and ground truth images using the	327	
	CLIP Similarity score (CLIP-S), computed with a pre-	328	
	trained CLIP encoder.	329	
	• 3DGS-Eval (3D): Our core contribution. We use model-	330	
	generated multi-view images to reconstruct a 3D scene	331	
	via 3D Gaussian Splatting. Metrics include PSNR , SSIM ,	332	
	and depth for rendering quality, and Chamfer Distance	333	
	(CD) for geometric fidelity against the ground truth.	334	
	Implementation Details. For 3DGS-Eval, we generate	335	
	8 views per structure. All baseline models are evaluated	336	

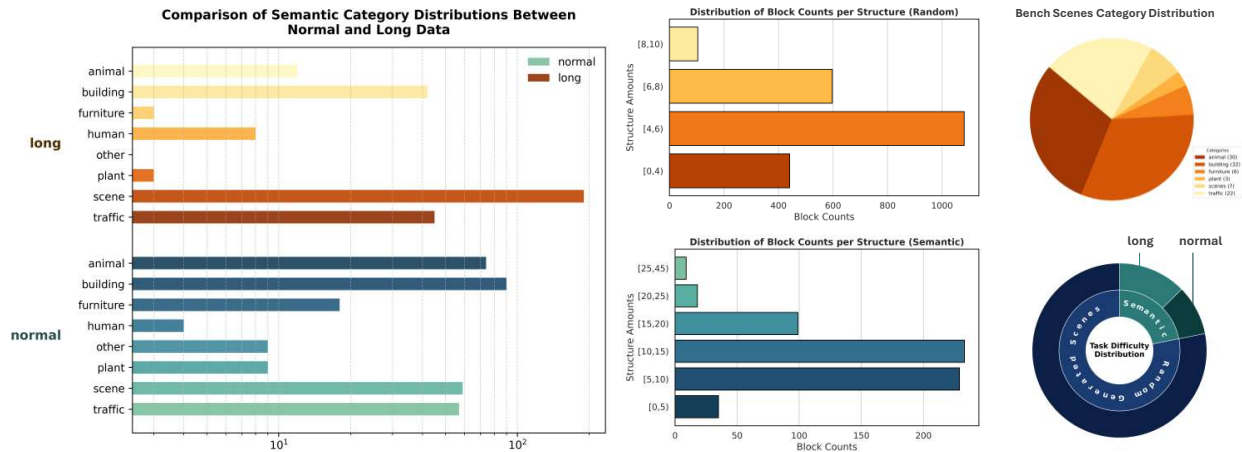


Figure 3. **Statistics of the synthesized dataset.** (Left) Comparison of semantic category distributions between "Normal" and "Long" data subsets. (Middle) Distribution of block counts (assembly steps) for Random and Semantic structures, illustrating the complexity gradient. (Right) Category distribution of bench scenes and the overall task difficulty ratio (Normal vs. Long) across Semantic and Random subsets.

Table 1. **Quantitative Comparison on MindBlock Benchmark.** We evaluate the single-step spatial intelligence of state-of-the-art Unified Multimodal Models (UMMs). The results are categorized into *Image Similarity* (2D heuristic) and our proposed *3D Gaussian Splatting* (volumetric consistency). \uparrow denotes higher is better; \downarrow denotes lower is better. Best results are **bolded**, and second best are underlined.

Task	Model	Image Similarity	3D Gaussian Splatting (3DGS-Eval)			
		CLIP-S \uparrow	PSNR \uparrow	SSIM \uparrow	Depth Err. \downarrow	CD \downarrow
<i>Spatial Assembly (Single-step Construction)</i>						
	BAGEL [3]	0.77	13.83	0.73	75.97	11.72
	Blip3-o-next [1]	<u>0.81</u>	11.90	0.69	83.20	13.05
	ILLUME [7]	0.68	9.40	0.61	101.30	16.20
	Emu3.5 [2]	0.75	13.10	0.71	79.20	12.10
	Step1X-Edit [9]	0.86	15.20	0.78	68.50	9.10
	Ovis-U1 [13]	0.76	<u>14.10</u>	<u>0.75</u>	<u>73.10</u>	<u>10.40</u>
	Omni-Gen [14]	0.70	8.70	0.59	108.40	17.80
<i>Spatial Structure (Geometric Equivariance)</i>						
	BAGEL [3]	0.74	12.26	0.76	82.64	10.73
	Blip3-o-next [1]	<u>0.79</u>	10.80	0.70	90.50	12.90
	ILLUME [7]	0.66	8.30	0.60	110.20	17.40
	Emu3.5 [2]	0.72	11.50	0.73	86.70	11.80
	Step1X-Edit [9]	0.83	13.80	0.80	74.20	9.30
	Ovis-U1 [13]	0.73	<u>12.90</u>	<u>0.77</u>	<u>79.60</u>	<u>10.20</u>
	Omni-Gen [14]	0.68	7.90	0.57	118.60	19.30

337 using their default inference hyper-parameters with a tem-
338 perature of 0.2.

339 4.2. Analysis of Single-step Assembly and Spatial 340 Structure

341 As shown in Table 1, we observe several important findings
342 regarding the spatial reasoning capabilities of current Uni-

fied Multimodal Models (UMMs).

The 2D–3D Consistency Gap. A clear discrepancy
343 emerges between 2D perceptual similarity and volumetric
344 consistency. For example, *Blip3-o-next* [1] achieves rela-
345 tively strong 2D alignment with a CLIP-S score higher than
346 *BAGEL* [3]. However, its 3D reconstruction metrics are
347 substantially worse, exhibiting lower PSNR/SSIM and sig-
348
349

Table 2. Performance during forward assembly and reverse deconstruction phases. \uparrow / \downarrow indicates direction of improvement.

Model	Image Similarity		3DGS-Eval (3D)		
	CLIP-S \uparrow	PSNR \uparrow	SSIM \uparrow	Depth Err. \downarrow	CD \downarrow
<i>Phase I: Forward Assembly</i>					
BAGEL-zebra-CoT [3]	0.78	13.5	0.74	71.3	11.8
ThinkMorph [6]	0.76	15.1	0.79	77.4	9.6
MindBlock-Random (Ours)	0.84	17.8	0.86	62.5	7.4
MindBlock-Semantic (Ours)	0.90	20.4	0.91	55.8	5.9
<i>Phase II: Reverse Deconstruction</i>					
BAGEL-zebra-CoT [3]	0.73	12.4	0.71	83.6	13.2
ThinkMorph [6]	0.75	14.0	0.76	76.8	10.7
MindBlock-Random (Ours)	0.82	16.3	0.83	67.9	8.5
MindBlock-Semantic (Ours)	0.88	19.1	0.88	60.2	6.8

350 significantly higher depth and Chamfer errors. This indicates
 351 that although individual generated views appear semanti-
 352 cally plausible, they fail to form a geometrically coherent
 353 3D scene when aggregated across viewpoints. In contrast,
 354 models such as *Ovis-UI* [13] achieve competitive or even
 355 better 3D reconstruction quality despite having comparable
 356 or slightly lower CLIP-S scores. These results highlight that
 357 strong 2D image similarity does not necessarily imply cor-
 358 rect spatial reasoning.

359 **Spatial Assembly Remains Challenging but Partially**
 360 **Solvable.** For the *Spatial Assembly* task, most models
 361 show moderate performance, with *StepIX-Edit* [9] achiev-
 362 ing the best overall results across both 2D and 3D met-
 363 rics. Models such as *Emu3.5* [2] and *Ovis-UI* [13] remain
 364 close to *BAGEL* [3], suggesting that current UMMs can par-
 365 tially reason about single-step object placement. However,
 366 large performance gaps remain between stronger models
 367 and weaker ones such as *ILLUME* [7] and *Omni-Gen* [14],
 368 which fail to produce geometrically stable structures.

369 **Equivariance is a Major Bottleneck.** All evaluated
 370 models show a consistent performance drop on the *Spa-*
 371 *tial Structure* task. Compared to *Spatial Assembly*, both
 372 2D and 3D metrics deteriorate across nearly all models,
 373 indicating that reasoning about viewpoint transformations
 374 is substantially more difficult. Even strong models fre-
 375 quently produce inconsistent internal layouts when asked
 376 to imagine the same structure from a new viewpoint, lead-
 377 ing to degraded reconstruction quality and larger geometric
 378 errors. This suggests that current UMMs still lack robust
 379 3D-equivariant internal representations.

380 **3D Evaluation Reveals Hidden Failures.** Importantly,
 381 several models that appear competitive under 2D metrics
 382 reveal substantial failures under our 3DGS-Eval protocol.
 383 This demonstrates that traditional image-based evaluation
 384 may significantly overestimate spatial reasoning ability. By
 385 explicitly reconstructing a volumetric scene from model-
 386 generated views, our evaluation exposes inconsistencies
 387 that remain invisible under standard 2D metrics.

4.3. Multi-step Assembly: The Challenge of Interleaved Reasoning

388 Multi-step assembly is a more rigorous probe of spatial
 389 intelligence, as it requires the model to generate an interleaved
 390 sequence of text (reasoning/instructions) and images (inter-
 391 mediate states). Most contemporary UMMs are limited to
 392 single-turn image generation and fail this task entirely. We
 393 evaluate the few models capable of long-context interleaved
 394 output: *BAGEL-zebra-CoT*, *Thinkmorph*, and our trained
 395 *MindBlock* variants.

396 **Model Specialization for Interleaved Reasoning.** Since
 397 the ability to generate long-context, interleaved text-image
 398 sequences is not yet a standard feature in most off-the-
 399 shelf UMMs, we perform continual training on the *BAGEL-*
 400 *zebra-CoT* [3] backbone to establish a competitive base-
 401 line for multi-step evaluation. We develop two specialized
 402 variants: **MindBlock-Random** and **MindBlock-Semantic**,
 403 fine-tuned on the respective partitions of our dataset². Dur-
 404 ing the training phase, the models are jointly supervised on
 405 two core assembly operators: *forward assembly* and *reverse*
 406 *deconstruction*. To enforce robust 3D structural awareness,
 407 each state transition is supervised across four views, requir-
 408 ing the model to maintain viewpoint-invariant information
 409 while progressing through the interleaved reasoning trajec-
 410 tory. This specialization ensures that the performance gains
 411 reported in Table 2 reflect the models’ internalized spatial
 412 logic rather than a mere failure to adhere to the interleaved
 413 format.

414 The results in Table 2 underscore the difficulty of main-
 415 taining structural integrity over multiple steps. General-
 416 purpose models like *Thinkmorph* suffer from “spatial drift,”
 417 where small errors in early steps accumulate, leading to
 418 physically impossible final structures. In contrast, our
 419 *MindBlock-Semantic* model, which benefits from a two-
 420 stage curriculum training on meaningful structures, demon-
 421 strates superior stability, effectively learning a native spatial
 422 intuition that goes beyond simple symbolic matching.

5. Further Analyses and Discussions

5.1. Is Spatial Reasoning Carried by CoT or Intuition?

423 A fundamental question in unified multimodal modeling is
 424 whether spatial structural awareness is an emergent prop-
 425 erty of generative “intuition” (pixel-level patterns) or if it
 426 relies on explicit symbolic reasoning. We investigate this by
 427 comparing UniCoT [12]—which generates a textual Chain-
 428 of-Thought (CoT) prior to the image—against the base
 429 *BAGEL* model [3] which performs direct image synthesis.

430 **The Divergence Metric.** We formalize the impact of rea-
 431 soning through the *Spatial Reasoning Gain* (Δ_{SR}), measur-
 432 ing the improvement in 3D structural integrity when CoT
 433

Table 3. Comparison of Spatial Intelligence between Direct Generation (BAGEL) and CoT-Augmented Generation (UniCoT)[12]. 3DGS-Eval metrics represent the quality of a 3D reconstruction from 8 generated views. CoT-based approach significantly reduces geometric errors (Chamfer Distance) while maintaining high semantic alignment.

Model	CLIP-S ↑	3DGS-PSNR ↑	3DGS-SSIM ↑	Chamfer Dist. ↓
BAGEL (Direct)	0.782	13.42	0.712	12.4
UniCoT[12]	0.814	20.88	0.785	6.5
Gain (Δ_{SR})	+0.032	+7.46	+0.073	-47.6%

438 is enabled. Let \mathcal{G}_{cot} and \mathcal{G}_{dir} denote the image generation
439 processes for UniCoT and the baseline, respectively. The
440 divergence is defined as:

$$441 \quad \Delta_{SR} = \mathbb{E}_{s \sim \mathcal{S}} [\Phi(\mathcal{G}_{cot}(I, T | \text{CoT})) - \Phi(\mathcal{G}_{dir}(I, T))] \quad (1)$$

442 where Φ represents the 3DGS-Eval score, and \mathcal{S} is the set
443 of single step assembly instructions in MindBlock.

444 **Quantitative Analysis: CoT as a Symbolic Scaffold.** As
445 shown in Table 3, direct generation (BAGEL) achieves rea-
446 sonable 2D similarity (CLIP-S) but fails under 3DGS-Eval.
447 Without the “textual anchor” of a CoT, the model often ne-
448 glects occluded block faces or misinterprets the depth of
449 “next-to” relations. In contrast, UniCoT shows a improve-
450 ment in 3D Reconstructibility. This suggests that the CoT
451 acts as a **symbolic scaffold**: by explicitly predicting the co-
452 ordinates and support relations (e.g., “*Block A is placed at*
453 *(x,y,z) on top of B*”), the model constrains the latent dif-
454 fusion process to respect 3D geometry rather than merely
455 hallucinating a 2D projection.

456 5.2. The Consistency Gap: Perception vs. Execu- 457 tion

458 To further point out whether the failure of spatial reasoning
459 originates from perceptual misunderstanding or generative
460 incapacity, we conduct a focused diagnostic analysis on the
461 **BAGEL** model. Specifically, we decouple the *Spatial As-*
462 *sembly* task into two paradigms: (1) **Discriminative Recog-**
463 **nitition**, where BAGEL must select the correct next-step im-
464 age from a 4-way Multiple-Choice Question (MCQ), and
465 (2) **Constructive Execution**, where the model generates the
466 next-step image in pixel space.

467 To objectively evaluate the semantic fidelity of the gen-
468 erated images, we introduce a **GPT-Matching** protocol. A
469 GPT-5 judge is presented with the model-generated image
470 alongside the four ground-truth and distractor options from
471 the MCQ. The judge identifies which of the four options the
472 generated image most closely resembles. If the judge se-
473 lects the ground-truth option, the execution is deemed suc-
474 cessful.

Table 4. **Diagnostic Analysis of BAGEL: The Knowledge-Action Gap.** We contrast the model’s accuracy in selecting the correct spatial step against its ability to generate a semantically aligned image of that same step.

Task Setting	Perception Acc. ↑	Execution Match ↑	The Gap (Δ) ↓
Single-step Assembly	52.5%	28.2%	24.3%

475 **The Generative Bottleneck.** As shown in Table 4, a dis-
476 crepancy exists between BAGEL’s internal spatial knowl-
477 edge and its constructive output. While the model correctly
478 identifies the intended spatial state in nearly half of the
479 cases (52.5%), it successfully executes that state in pixel
480 space only 28.2% of the time. This 24.3% gap reveals that
481 spatial intelligence in current UMMs is often *superficial*:
482 models may possess the discriminative “intuition” to recog-
483 nize correct 3D relationships, yet they lack the **generative**
484 **grounding** required to translate this symbolic understand-
485 ing into consistent, physically-plausible visual structures.
486 Our GPT-Matching results further indicate that failed gener-
487 ations frequently drift toward visual distractors, even when
488 the model’s discriminative head had correctly rejected them
489 during the MCQ phase.

490 6. Conclusion

491 In this work, we introduced **MindBlock**, a diagnostic
492 benchmark for evaluating spatial intelligence in unified
493 multimodal models through *active construction*. Moving
494 beyond passive visual question answering, MindBlock re-
495 quires models to generate step-by-step visual assemblies
496 and maintain structural consistency under spatial transfor-
497 mations. To rigorously evaluate geometric coherence, we
498 further proposed **3DGS-Eval**, a volumetric consistency pro-
499 tocol that reconstructs a 3D representation from model-
500 generated multi-view images.

501 Our key findings are three-fold. First, a *2D–3D con-*
502 *sistency gap* pervades current UMMs: models that achieve
503 high perceptual similarity scores often produce multi-view
504 generations that are geometrically incoherent, a failure re-
505 vealed only through our **3DGS-Eval** protocol. Second, ge-
506 ometric equivariance remains a critical bottleneck, with all
507 evaluated models degrading substantially when asked to
508 reason about the same structure from a novel viewpoint.
509 Third, a pronounced *perception–execution gap* (24.3% for
510 BAGEL) demonstrates that models frequently understand
511 what to build yet fail to build it, suggesting that discrim-
512 inative spatial knowledge does not automatically transfer
513 to generative spatial grounding. Explicit Chain-of-Thought
514 reasoning partially bridges this gap, acting as a symbolic
515 scaffold that stabilizes 3D structure during generation.

516

References

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

- [1] Jiu hai Chen, Le Xue, Zhiyang Xu, Xichen Pan, Shusheng Yang, Can Qin, An Yan, Honglu Zhou, Zeyuan Chen, Lifu Huang, Tianyi Zhou, Junnan Li, Silvio Savarese, Caiming Xiong, and Ran Xu. Blip3o-next: Next frontier of native image generation, 2025. 5, 6
- [2] Yufeng Cui, Honghao Chen, Haoge Deng, Xu Huang, Xinghang Li, Jirong Liu, Yang Liu, Zhuoyan Luo, Jinsheng Wang, Wenxuan Wang, Yueze Wang, Chengyuan Wang, Fan Zhang, Yingli Zhao, Ting Pan, Xianduo Li, Zecheng Hao, Wenxuan Ma, Zhuo Chen, Yulong Ao, Tiejun Huang, Zhongyuan Wang, and Xinlong Wang. Emu3.5: Native multimodal models are world learners, 2025. 5, 6, 7
- [3] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, et al. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025. 2, 5, 6, 7
- [4] Andrea Frick, Katrina Ferrara, and Nora S Newcombe. Using a touch screen paradigm to assess the development of mental rotation between 31/2 and 51/2 years of age. *Cognitive processing*, 14(2):117–127, 2013. 2
- [5] Genesis Authors. Genesis: A universal and generative physics engine for robotics and beyond. <https://github.com/Genesis-Embodied-AI/Genesis>, 2024. 5
- [6] Jiawei Gu, Yunzhuo Hao, Huichen Will Wang, Linjie Li, Michael Qizhe Shieh, Yejin Choi, Ranjay Krishna, and Yu Cheng. Thinkmorph: Emergent properties in multimodal interleaved chain-of-thought reasoning, 2026. 7
- [7] Runhui Huang, Chunwei Wang, Junwei Yang, Guansong Lu, Yunlong Yuan, Jianhua Han, Lu Hou, Wei Zhang, Lanqing Hong, Hengshuang Zhao, and Hang Xu. Illume+: Illuminating unified mllm with dual visual tokenization and diffusion refinement, 2025. 5, 6, 7
- [8] Barbara Landau, E Emory Davis, Cathryn S Cortesa, Zihan Wang, Jonathan D Jones, and Amy L Shelton. Young children’s copying of block constructions: Significant constraints in a highly complex task. *Cognitive Development*, 71:101463, 2024. 2
- [9] Shiyu Liu, Yucheng Han, Peng Xing, Fukun Yin, Rui Wang, Wei Cheng, Jiaqi Liao, Yingming Wang, Honghao Fu, Chunrui Han, Guopeng Li, Yuang Peng, Quan Sun, Jingwei Wu, Yan Cai, Zheng Ge, Ranchen Ming, Lei Xia, Xianfang Zeng, Yibo Zhu, Binxing Jiao, Xiangyu Zhang, Gang Yu, and Daxin Jiang. Step1x-edit: A practical framework for general image editing, 2025. 5, 6, 7
- [10] Liang Ma, Jiajun Wen, Min Lin, Rongtao Xu, Xiwen Liang, Bingqian Lin, Jun Ma, Yongxin Wang, Ziming Wei, Haokun Lin, Mingfei Han, Meng Cao, Bokui Chen, Ivan Laptev, and Xiaodan Liang. Phyblock: A progressive benchmark for physical understanding and planning via 3d block assembly, 2025. 2, 5
- [11] Ava Pun, Kangle Deng, Ruixuan Liu, Deva Ramanan, Changliu Liu, and Jun-Yan Zhu. Generating physically stable and buildable brick structures from text. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 14798–14809, 2025. 2
- [12] Luozheng Qin, Jia Gong, Yuqing Sun, Tianjiao Li, Mengping Yang, Xiaomeng Yang, Chao Qu, Zhiyu Tan, and Hao Li. Uni-cot: Towards unified chain-of-thought reasoning across text and vision, 2026. 7, 8
- [13] Guo-Hua Wang, Shanshan Zhao, Xinjie Zhang, Liangfu Cao, Pengxin Zhan, Lunhao Duan, Shiyin Lu, Minghao Fu, Xiaohao Chen, Jianshan Zhao, Yang Li, and Qing-Guo Chen. Ovis-u1 technical report, 2025. 5, 6, 7
- [14] Shitao Xiao, Yueze Wang, Junjie Zhou, Huaying Yuan, Xingrun Xing, Ruiran Yan, Chaofan Li, Shuting Wang, Tiejun Huang, and Zheng Liu. Omnigen: Unified image generation, 2024. 5, 6, 7
- [15] Jinheng Xie, Weijia Mao, Zechen Bai, David Junhao Zhang, Weihao Wang, Kevin Qinghong Lin, Yuchao Gu, Zhijie Chen, Zhenheng Yang, and Mike Zheng Shou. Show-o: One single transformer to unify multimodal understanding and generation. *arXiv preprint arXiv:2408.12528*, 2024. 2
- [16] Xianghui Xie, Chuhang Zou, Meher Gitika Karumuri, Jan Eric Lenssen, and Gerard Pons-Moll. Mvgbench: Comprehensive benchmark for multi-view generation models, 2025. 5
- [17] Jianke Zhang, Yanjiang Guo, Yucheng Hu, Xiaoyu Chen, Xiang Zhu, and Jianyu Chen. Up-vla: A unified understanding and prediction model for embodied agent. In *International Conference on Machine Learning*, pages 74911–74922. PMLR, 2025. 2

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

587

588

589

590

591

592

593

594

595

596

597

598

599