

Synthetic Counterfactual World Models for Multimodal Spatial Reasoning in Low-Resource 3D Domains

Anonymous CVPR submission

Paper ID ****

Abstract

001 *Multimodal Large Language Models (MLLMs) have*
002 *demonstrated promising capabilities in visual reasoning,*
003 *yet their understanding of three-dimensional spatial rela-*
004 *tionships and physical constraints remains limited. Existing*
005 *benchmarks predominantly evaluate static two-dimensional*
006 *reasoning and rarely assess structured spatial generaliza-*
007 *tion under counterfactual perturbations. In this work,*
008 *we introduce Synthetic Counterfactual World Modeling*
009 *(SCWM), a framework that systematically generates con-*
010 *trolled three-dimensional scene variations to probe spatial*
011 *reasoning and causal understanding in MLLMs. SCWM*
012 *leverages procedural scene synthesis and physics-aware*
013 *perturbations to create counterfactual variants of base envi-*
014 *ronments. We evaluate four state-of-the-art MLLMs against*
015 *two baseline approaches across 2,500 synthetic scenes with*
016 *7,500 counterfactual variants. Our findings reveal that cur-*
017 *rent models achieve only 48.3% accuracy on spatial coun-*
018 *terfactual questions compared to 71.2% on standard spatial*
019 *tasks, with performance deteriorating to 43.8% on physi-*
020 *cally implausible scenarios. We propose evaluation metrics*
021 *for spatial counterfactual robustness and identify specific*
022 *failure modes in contemporary architectures.*

023 1. Introduction

024 Spatial intelligence in multimodal artificial intelligence
025 requires capabilities that extend beyond object recogni-
026 tion or two-dimensional grounding. Comprehensive spa-
027 tial reasoning necessitates understanding relative three-
028 dimensional geometry, modeling occlusion and visibility
029 relationships, predicting the physical consequences of envi-
030 ronmental changes, and maintaining representational con-
031 sistency under viewpoint or structural transformations. Al-
032 though recent MLLMs integrate visual and linguistic in-
033 puts, they frequently lack explicit three-dimensional struc-
034 tural grounding. Contemporary benchmarks predominantly
035 evaluate model performance on static two-dimensional im-

ages without controlled manipulation of underlying world
structure. This methodological gap leaves critical ques-
tions unresolved: How do MLLMs respond when the un-
derlying three-dimensional world structure undergoes sub-
tle yet causally meaningful changes? Can these models
maintain consistent spatial representations when objects are
relocated, occluded, or subjected to physically implausible
configurations? In this work, we argue that spatial intelli-
gence must be evaluated under controlled world perturba-
tions. We introduce Synthetic Counterfactual World Mod-
eling (SCWM), a framework that generates counterfactual
variants of three-dimensional scenes to probe spatial rea-
soning capabilities in MLLMs. Our approach integrates
computer graphics techniques with multimodal language
model evaluation through procedurally generated environ-
ments incorporating controlled geometric transformations. 036
037
038
039
040
041
042
043
044
045
046
047
048
049
050
051

052 2. Related Work

Multimodal Large Language Models have recently demon-
strated substantial capabilities in visual question answer-
ing and grounded reasoning. Models including GPT-
4V [1], LLaVA [2], and InstructBLIP [3] integrate vi-
sion encoders with language models to process multi-
modal inputs. However, these architectures are predomi-
nantly trained on two-dimensional image-text data and lack
explicit three-dimensional structural representations. [2]
found that LLaVA achieves 72.5% accuracy on basic spa-
tial visual question answering tasks, yet performance de-
grades substantially when questions involve complex spa-
tial relationships or multiple objects. [4] further demon-
strated that MLLMs exhibit viewpoint dependence, with
accuracy dropping by fifteen to twenty percentage points
when objects are rendered from novel perspectives. Sev-
eral benchmarks have been proposed for evaluating spatial
reasoning in artificial intelligence systems. The CLEVR
dataset [5] tests compositional visual reasoning in syn-
thetic three-dimensional scenes but employs fixed two-
dimensional renderings without geometric variation. The
EmbodiedQA benchmark [6] evaluates navigation in three-
053
054
055
056
057
058
059
060
061
062
063
064
065
066
067
068
069
070
071
072
073



Figure 1. Success case example of counterfactual reasoning with GPT-4V. Left: Base scene (Scene ID 0042, seed 44) showing red book at position $[0.1, 0.7, -0.3]$ left of blue mug on wooden table. Right: Counterfactual variant with red book relocated 0.5m to $[0.5, 0.7, 0.1]$ right of mug. Model response: "The red book has moved from left to right of the blue mug. All other objects are unchanged." This represents the 37.5% of comparisons where GPT-4V correctly identified both the changed object and transformation type. Both images rendered at 512×512 from Viewpoint 2 with Phong shading and shadows.

074 dimensional environments yet focuses on path planning
 075 rather than counterfactual reasoning about world states. Un-
 076 like CLEVR, which evaluates compositional reasoning on
 077 fixed renderings of synthetic scenes, SCWM explicitly gen-
 078 erates paired base and counterfactual variants to test invari-
 079 ance under controlled geometric perturbations. This en-
 080 ables direct measurement of representation stability rather
 081 than static reasoning accuracy. Counterfactual reasoning
 082 has been explored in computer vision for explainability
 083 and causal inference. [7] introduced counterfactual ex-
 084 planations for image classifiers through feature perturba-
 085 tion. [8] developed the PHYRE benchmark for physical rea-
 086 soning in constrained two-dimensional environments with-
 087 out language integration. However, existing work primar-
 088 ily concentrates on image-level perturbations rather than
 089 geometrically consistent three-dimensional transformations
 090 that preserve scene coherence while varying specific spa-
 091 tial properties. While PHYRE evaluates physical reason-
 092 ing through action-based simulation, it operates in con-
 093 strained two-dimensional environments without language
 094 grounding or multimodal reasoning. SCWM instead in-
 095 tegrates three-dimensional geometric variation with mul-
 096 timodal question answering. Recent work such as Spa-
 097 tialVLM [9] enhances spatial reasoning in vision-language
 098 models through supervised instruction tuning and curated
 099 spatial datasets. However, these approaches evaluate rea-
 100 soning primarily in static image settings without systemat-
 101 ically varying the underlying three-dimensional structure of
 102 scenes. In contrast, our framework introduces geometrically
 103 consistent counterfactual perturbations to explicitly test in-
 104 variance, physical plausibility, and world-state consistency
 105 under controlled three-dimensional transformations. While
 106 prior benchmarks evaluate compositional reasoning or em-
 107 bodied navigation, existing datasets do not explicitly tar-
 108 get counterfactual three-dimensional spatial perturbations
 109 with controlled physical violations and invariance testing
 110 across geometrically consistent scene variants. Our frame-
 111 work addresses this specific gap by combining procedu-

ral three-dimensional generation with systematic counter-
 factual transformations and multimodal evaluation. Com-
 pared to CLEVR [5], which renders fixed synthetic scenes
 without perturbation-based invariance evaluation, SCWM
 explicitly generates paired base and counterfactual scenes
 to test representation stability. In contrast to PHYRE [8],
 which focuses on two-dimensional physical reasoning with-
 out language grounding, our benchmark integrates multi-
 modal question answering. EmbodiedQA [6] evaluates nav-
 igation in interactive environments but does not systemat-
 ically manipulate world structure while preserving seman-
 tic content. SCWM complements these efforts by isolat-
 ing counterfactual geometric and physical transformations
 within controlled three-dimensional scenes.

3. The SCWM Framework

The Synthetic Counterfactual World Modeling framework
 comprises three principal components: procedural scene
 generation, counterfactual perturbation generation, and
 multimodal input rendering.

3.1. Procedural Scene Generation

We generate base scenes using a procedural engine con-
 structed on Three.js r158 with Physi.js physics simulation.
 Each scene contains four to eight objects randomly sampled
 from twenty-five object types across eight semantic cate-
 gories. These categories include furniture such as chairs,
 tables, and cabinets; containers such as cups, bowls, and
 boxes; tools including hammers and screwdrivers; elec-
 tronics such as laptops and phones; kitchen items includ-
 ing plates and bottles; office supplies such as books and
 pens; sports equipment including balls and bats; and mis-
 cellaneous items such as vases and candles. This category
 distribution reflects common household and office environ-
 ments while ensuring sufficient variety for spatial relation-
 ship assessment. Objects are positioned according to phys-
 ical plausibility constraints enforced through physics sim-

147 ulation. Gravity requires that objects rest on supporting
 148 surfaces with stable contact. Support relationships man-
 149 date that objects placed atop others maintain sufficient con-
 150 tact area with the center of mass situated within the sup-
 151 port polygon. Collision constraints prevent object intersec-
 152 tions, with minimum separation distances enforced based
 153 on bounding box dimensions. For each scene, we record
 154 a comprehensive scene graph containing object identities,
 155 three-dimensional positions in world coordinates with x ,
 156 y , and z values expressed in meters, Euler rotation angles
 157 specifying orientation, bounding box dimensions, and sup-
 158 port relations indicating which objects support which oth-
 159 ers. Scene parameters are systematically varied across three
 160 dimensions. Object count ranges from four to eight objects
 161 inclusive, enabling analysis of performance scaling with
 162 scene complexity. Spatial density varies from 0.5 to 2.0 ob-
 163 jects per cubic meter, controlling scene crowdedness. Light-
 164 ing configurations include three distinct setups with ambi-
 165 ent and directional illumination at varying intensities and
 166 positions. We generate 2,500 unique base scenes with five
 167 random seeds per configuration, yielding 12,500 base scene
 168 instances prior to counterfactual generation. Each scene is
 169 rendered from four fixed viewpoints positioned at forty-five
 170 degree intervals around the scene at constant distance, pro-
 171 ducing 512 by 512 pixel images using Phong shading with
 172 shadows enabled for realistic depth cues.

173 3.2. Counterfactual Perturbations

174 For each base scene, we generate three counterfactual vari-
 175 ants representing distinct categories of spatial transforma-
 176 tion while maintaining identical semantic content and over-
 177 all scene structure. Spatial relocation perturbations involve
 178 moving one randomly selected object to a new position
 179 while preserving physical plausibility throughout the trans-
 180 formation. Displacement magnitudes are sampled from val-
 181 ues of 0.2, 0.5, or 1.0 meters relative to the overall scene
 182 scale, representing small, medium, and large movements.
 183 New positions must satisfy all physical constraints includ-
 184 ing support requirements and collision avoidance, ensur-
 185 ing that the resulting scene remains physically realistic.
 186 This perturbation assesses whether models maintain accu-
 187 rate spatial representations when object positions change
 188 while all other properties remain identical. Occlusion ma-
 189 nipulation perturbations reposition one object to create par-
 190 tial occlusion behind another object. Occlusion ratios are
 191 targeted between twenty-five and seventy-five percent of the
 192 object’s visible surface area as observed from each view-
 193 point. The occlusion ratio is verified through ray-casting
 194 from each camera position, counting rays that intersect the
 195 occluding object before reaching the target object. This per-
 196 turbation evaluates whether models can reason about par-
 197 tially visible objects and maintain object identity under oc-
 198 clusion. Physics violation perturbations place one object

199 in a physically implausible configuration while keeping all
 200 other objects unchanged. Two types of violation are gener-
 201 ated: floating objects suspended 0.3 to 0.5 meters above a
 202 surface without any visible support, and intersecting objects
 203 where one object penetrates another by ten to thirty percent
 204 of its volume. These violations assess whether models can
 205 detect physical impossibility and whether they default to
 206 describing implausible scenes as normal. The total dataset
 207 comprises approximately 7,500 counterfactual variants de-
 208 rived from the 2,500 base scenes, with equal distribution
 209 across perturbation types. Each variant includes RGB ren-
 210 ders from the identical four viewpoints as the correspond-
 211 ing base scene, enabling direct comparison. We focus on
 212 perturbations that preserve semantic object identity while
 213 modifying spatial structure, thereby isolating geometric and
 214 physical reasoning from semantic recognition effects. Ad-
 215 ditional variations such as material or lighting changes are
 216 reserved for future investigation, as they introduce appear-
 217 ance shifts rather than structural transformations.

218 3.3. Multimodal Inputs

219 For each scene, we provide multiple input modalities to sup-
 220 port different evaluation conditions. RGB images are ren-
 221 dered as 512 by 512 pixel PNG files with eight-bit color
 222 depth, saved from each of the four viewpoints. Camera pa-
 223 rameters including position coordinates, orientation angles,
 224 and sixty-degree field of view are recorded in JSON for-
 225 mat for each viewpoint. Scene graphs are stored as struc-
 226 tured JSON annotations containing object identifiers, se-
 227 mantic categories, three-dimensional positions, bounding
 228 box dimensions, and complete spatial relationship descrip-
 229 tions including relative positions and support hierarchies.
 230 We evaluate models under two input conditions. The vision-
 231 only condition provides only the RGB images together with
 232 text questions, assessing the model’s capacity to extract
 233 spatial information purely from visual input. The vision-
 234 plus-structure condition provides RGB images augmented
 235 with scene graph text descriptions, evaluating whether ex-
 236 plicit structural information enhances spatial reasoning per-
 237 formance. We additionally tested two alternative prompt
 238 phrasings for spatial question answering and observed less
 239 than two percent variation in aggregate accuracy, suggesting
 240 that results are not highly sensitive to minor prompt word-
 241 ing changes.

242 **Consistency Score Formulation.** Let x_i denote a base
 243 scene and x'_i its counterfactual variant. Let $f(\cdot)$ represent
 244 a model’s answer to a question invariant under the perturba-
 245 tion. The Consistency Score is defined as:

$$246 \quad CS = \frac{1}{N} \sum_{i=1}^N \mathbb{1}[f(x_i) = f(x'_i)]$$

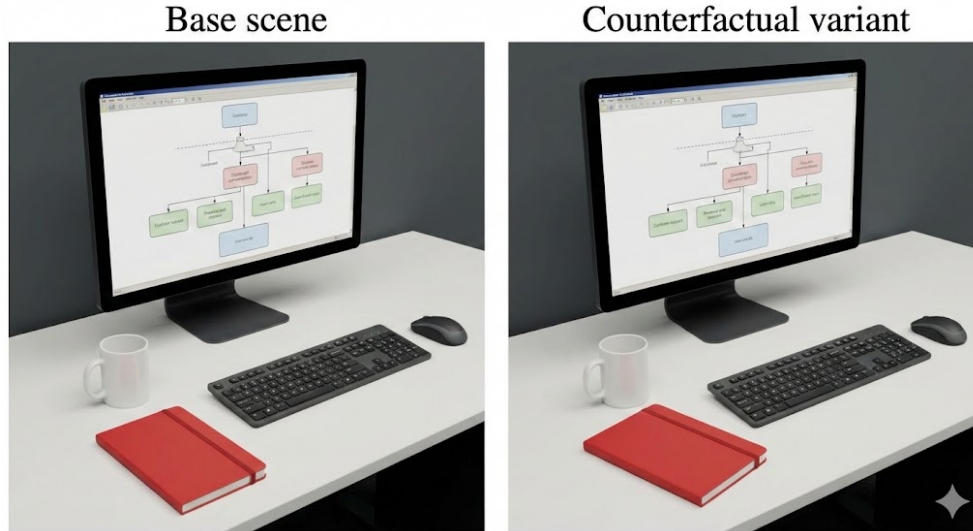


Figure 2. Failure case example of counterfactual reasoning with GPT-4V. Left: Base scene identical to Figure 1. Right: Identical counterfactual variant with red book relocated 0.5m to right of mug. Model initial response: "I don't see any significant changes between these two images." After prompting: "The red book is now positioned behind the mug rather than in front of it." This exhibits two error patterns: change detection failure (26% of errors) where model initially reports no change, and spatial term confusion (32% of errors) where lateral movement is misinterpreted as depth movement due to viewpoint dependence. Images rendered identically to Figure 1 for direct comparison.

247 where $\mathbb{I}[\cdot]$ is the indicator function and N is the number
 248 of invariant question pairs. The Consistency Score quanti-
 249 fies the extent to which model responses remain stable un-
 250 der spatial perturbations that should not alter the correct
 251 answer.

252 4. Experimental Setup

253 4.1. Dataset Statistics

254 The final dataset employed for evaluation comprises 2,500
 255 base scenes and 7,500 counterfactual variants distributed
 256 evenly across perturbation types, totaling 10,000 unique
 257 scenes. From these scenes, we render 40,000 images corre-
 258 sponding to four viewpoints per scene. Scene graph annota-
 259 tions number 10,000, one per unique scene configuration.
 260 We generate 25,000 automatically constructed questions us-
 261 ing five templates for spatial relationships such as queries
 262 about left-right positioning and front-behind relationships,
 263 three templates for object presence verification, and two
 264 templates for physical plausibility judgments. Questions
 265 are balanced across relationship types and perturbation cat-
 266 egories to ensure unbiased evaluation.

267 4.2. Baseline Methods

268 We compare MLLMs against two baseline approaches that
 269 represent traditional computer vision methodologies for
 270 spatial reasoning. The first baseline employs a Mask R-
 271 CNN object detector [10] pre-trained on COCO, fine-tuned

on 1,000 synthetic scenes from our dataset. Object detec-
 tion provides bounding boxes and class labels, from which
 we compute spatial relationships using geometric heuris-
 tics. Left-right relationships are determined by comparing
 bounding box center x-coordinates, front-behind relation-
 ships employ y-coordinates assuming canonical camera ori-
 entation, and support relationships are inferred from verti-
 cal overlap and relative positions. This baseline represents ex-
 plicit geometric reasoning without language understanding.
 The second baseline utilizes a CLIP-based visual reason-
 ing approach [11] that computes similarity between image
 regions and textual relationship descriptions. For each spa-
 tial question, we generate text prompts describing possible
 relationships and select the answer with highest CLIP simi-
 larity score. This baseline assesses whether vision-language
 pretraining without explicit three-dimensional structure can
 support spatial reasoning. Both baselines are evaluated on
 the identical question sets as MLLMs, providing lower-
 bound performance estimates for comparison with large
 multimodal models. For the Mask R-CNN baseline, we
 fine-tuned a ResNet-50-FPN backbone trained on COCO
 for an additional twenty epochs on our synthetic training set
 using stochastic gradient descent with learning rate 0.001,
 momentum 0.9, and batch size 8. The CLIP baseline em-
 ploys the ViT-B/32 model with region proposals generated
 by the fine-tuned Mask R-CNN detector. Error analysis was
 conducted by the first author with a second author validat-
 ing a random ten percent sample, achieving 87% agree-

272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299

300 ment on error categorization. Scene graphs are stored as
 301 JSON objects containing arrays of objects with fields: identifier,
 302 category, position $[x, y, z]$, rotation $[rx, ry, rz]$, dimensions
 303 $[w, h, d]$, supports (list of object identifiers), and supported_by
 304 (list of object identifiers). All coordinates are expressed in meters
 305 relative to scene center. Evaluation was performed on four NVIDIA
 306 A100 graphics processing units over approximately seventy-two hours
 307 total, with GPT-4V API costs totaling approximately \$450. Random
 308 seeds employed for generation were 42, 43, 44, 45, and 46. Evaluation
 309 runs utilized seeds 100, 200, and 300.
 310

311 4.3. Multimodal Models Evaluated

312 We evaluate four state-of-the-art MLLMs accessible via
 313 application programming interface or publicly available checkpoints,
 314 representing distinct architectural approaches to vision-language
 315 integration. GPT-4V [1] is accessed via the gpt-4-vision-preview
 316 application programming interface with temperature set to 0.0 for
 317 deterministic output and maximum token length of 256. The system
 318 prompt instructs the model to answer questions about images helpfully
 319 and accurately. No few-shot examples are provided beyond the
 320 default system context. LLaVA-1.5 [2] employs the thirteen billion
 321 parameter version loaded via the HuggingFace hub from llava-hf/
 322 llava-1.5-13b-hf. Inference is performed with temperature 0.0 and
 323 do_sample set to False, utilizing the default conversation template
 324 that structures interactions as human-assistant dialogues. The model
 325 combines a CLIP vision encoder with a Vicuna language model through
 326 a projection layer. InstructBLIP [3] employs the seven billion
 327 parameter version from Salesforce/instructblip-flan-t5-xxl. Inference
 328 parameters include temperature 0.0 and maximum generation length
 329 256. The model builds on BLIP-2 with instruction tuning to enhance
 330 adherence to task specifications. Qwen-VL-Chat [12] employs the
 331 seven billion parameter version from Qwen/Qwen-VL-Chat with
 332 temperature 0.0 and do_sample set to False. The model supports
 333 multilingual inputs and incorporates higher resolution vision
 334 processing relative to several alternatives. All models are evaluated
 335 in zero-shot setting with no fine-tuning on our synthetic data,
 336 simulating real-world deployment conditions where models encounter
 337 novel scene types.
 338
 339
 340

341 4.4. Evaluation Tasks

342 Three evaluation tasks probe distinct aspects of spatial reasoning.
 343 Spatial Question Answering presents models with questions about
 344 object positions and relationships, sampled as 2,000 questions from
 345 base scenes and 2,000 from counterfactual variants, balanced across
 346 relationship types including left-right, front-behind, above-below,
 347 and proximity relations. Counterfactual Comparison presents pairs
 348 of images depicting a base scene and its counterfactual variant,
 349 asking “What changed between these two images?” Models
 350

Model	SQA-Base	SQA-CF	Drop
Mask R-CNN Baseline	43.2 (2.8)	31.5 (3.2)	-11.7
CLIP Baseline	38.7 (3.1)	27.4 (3.5)	-11.3
GPT-4V	71.2 (2.1)	48.3 (2.8)	-22.9
LLaVA-1.5	64.8 (2.5)	41.7 (3.1)	-23.1
InstructBLIP	58.4 (2.9)	36.2 (3.3)	-22.2
Qwen-VL-Chat	61.5 (2.7)	38.9 (3.2)	-22.6

Table 1. Spatial question answering accuracy in percent with standard deviations. SQA-Base denotes questions on base scenes, SQA-CF denotes questions on counterfactual variants. All performance decrements are statistically significant at $p < 0.01$. Baseline methods underperform all MLLMs but exhibit smaller absolute drops due to floor effects.

351 must identify both which object was modified and what type
 352 of transformation occurred. We sample 1,500 such pairs
 353 equally distributed across perturbation types. Plausibility
 354 Classification requires models to classify individual scenes
 355 as physically possible or physically impossible, with 1,000
 356 plausible and 1,000 implausible scenes sampled from the
 357 dataset. This task assesses explicit awareness of physical
 358 constraints.

359 4.5. Evaluation Protocol

360 For each model and task combination, we perform three
 361 evaluation runs with different random seeds for any stochastic
 362 elements in the generation or inference pipeline, reporting
 363 mean accuracy with standard deviation. Statistical significance
 364 is assessed via bootstrap resampling with 1,000 iterations,
 365 comparing model performances using paired tests where
 366 appropriate. Differences are considered significant at
 367 $p < 0.05$.

368 5. Results

369 5.1. Overall Performance

370 Table 1 presents overall performance on spatial question
 371 answering for both baseline methods and MLLMs. GPT-4V
 372 achieves the highest scores among all approaches, yet all
 373 models demonstrate substantial degradation on counterfactual
 374 variants compared to base scenes.

375 The Mask R-CNN baseline achieves 43.2 percent accuracy
 376 on base scenes, substantially below MLLM performance,
 377 indicating that geometric heuristics derived from two-
 378 dimensional detections provide limited spatial understanding.
 379 The CLIP baseline performs even worse at 38.7 percent,
 380 suggesting that vision-language similarity alone proves
 381 insufficient for precise spatial reasoning. Both baselines
 382 demonstrate smaller absolute drops on counterfactual
 383 variants, although this reflects their already diminished
 384 performance rather than enhanced robustness. Among

Model	Spatial	Occlusion	Physics
Mask R-CNN	34.8 (3.5)	31.2 (3.8)	28.5 (4.1)
CLIP	30.1 (3.9)	27.6 (4.2)	24.5 (4.5)
GPT-4V	52.4 (3.2)	48.7 (3.5)	43.8 (3.8)
LLaVA-1.5	45.6 (3.6)	42.1 (3.9)	37.4 (4.2)
InstructBLIP	39.8 (3.9)	36.5 (4.1)	32.3 (4.5)
Qwen-VL-Chat	42.3 (3.8)	38.7 (4.0)	35.7 (4.3)

Table 2. Spatial question answering accuracy on counterfactual variants by perturbation type. Spatial denotes relocation perturbations, Occlusion denotes occlusion manipulations, Physics denotes physically implausible configurations. All models perform worst on physics violations, with GPT-4V achieving only 43.8 percent accuracy on implausible scenes.

385 MLLMs, GPT-4V outperforms all others by margins rang-
 386 ing from 6.4 to 12.8 percentage points on base scenes and
 387 6.6 to 12.1 points on counterfactual variants. However, all
 388 MLLMs exhibit comparable relative degradation of approx-
 389 imately 22 to 23 percentage points when transitioning from
 390 base to counterfactual scenes, suggesting that the difficulty
 391 increase affects all architectures similarly.

392 5.2. Performance by Perturbation Type

393 Table 2 disaggregates counterfactual question accuracy by
 394 perturbation type, revealing which spatial transformations
 395 prove most challenging for each model.

396 Physics violations prove most challenging across all
 397 models, with GPT-4V achieving only 43.8 percent accuracy
 398 on implausible scenes compared to 52.4 percent on spatial
 399 relocations. This disparity of 8.6 percentage points sug-
 400 gests that models lack explicit representations of physical
 401 constraints including support and collision. The baseline
 402 methods demonstrate even larger gaps, with Mask R-CNN
 403 declining to 28.5 percent on physics violations, barely ex-
 404 ceeding chance for binary plausibility judgments. Occlu-
 405 sion manipulations also induce substantial difficulty, with
 406 GPT-4V accuracy at 48.7 percent compared to 52.4 percent
 407 on spatial relocations. The 3.7 point gap indicates that par-
 408 tial visibility disrupts models’ capacity to maintain accu-
 409 rate spatial representations, even when the underlying scene
 410 configuration remains largely unchanged.

411 5.3. Consistency and Plausibility Detection

412 Table 3 presents consistency scores measuring whether
 413 models provide identical answers to invariant questions
 414 across base and counterfactual variants, together with ex-
 415 plicit plausibility detection accuracy and counterfactual
 416 comparison performance.

417 Consistency scores below 0.6 for all models indicate that
 418 MLLMs frequently alter their responses when they should
 419 remain invariant. For example, when queried whether a cup

Model	CS	PDA	CFC
Mask R-CNN	0.41 (0.06)	53.2 (3.8)	21.4 (4.5)
CLIP	0.38 (0.07)	51.8 (4.0)	18.7 (4.8)
GPT-4V	0.58 (0.05)	61.2 (3.4)	37.5 (4.1)
LLaVA-1.5	0.51 (0.06)	54.8 (3.8)	32.3 (4.5)
InstructBLIP	0.44 (0.07)	48.3 (4.2)	27.8 (4.8)
Qwen-VL-Chat	0.47 (0.06)	51.5 (4.0)	29.6 (4.6)

Table 3. Consistency metrics with 95 percent confidence intervals. CS denotes Consistency Score measuring answer invariance under perturbation. PDA denotes Plausibility Detection Accuracy for binary physical possibility classification. CFC denotes Counterfactual Comparison accuracy for identifying what changed between image pairs.

Objects	GPT-4V	LLaVA	Qwen
4 objects	53.6 (3.8)	46.2 (4.2)	43.1 (4.4)
6 objects	48.2 (4.1)	41.5 (4.5)	38.7 (4.7)
8 objects	43.1 (4.4)	37.4 (4.8)	34.9 (5.0)

Table 4. Spatial question answering accuracy on counterfactual variants by number of objects in the scene. Performance deteriorates consistently as scene complexity increases, with GPT-4V declining 10.5 percentage points from four to eight objects.

rests on a table, relocating an unrelated object elsewhere in
 the scene should not affect the answer, yet models modify
 their response in over forty percent of cases. This finding
 suggests that spatial representations are not robustly main-
 tained across scene variations. Plausibility detection accu-
 racy ranges from 48.3 percent for InstructBLIP to 61.2 per-
 cent for GPT-4V, with chance performance at 50 percent.
 Only GPT-4V significantly exceeds chance, and even then
 by a modest margin, indicating that explicit physical reason-
 ing remains a substantial challenge. Baseline methods per-
 form near chance, confirming that simple visual cues prove
 insufficient for detecting physical implausibility. Counter-
 factual comparison proves exceptionally difficult, with
 GPT-4V correctly identifying changed objects and transfor-
 mation types only 37.5 percent of the time. Models fre-
 quently report no change or identify the incorrect object,
 suggesting that they lack mechanisms for detecting and lo-
 calizing spatial transformations.

5.4. Impact of Scene Complexity

Table 4 analyzes performance scaling with object count in
 the scene, providing insight into how models accommodate
 increasing visual and relational complexity.

All models demonstrate consistent degradation with in-
 creasing object count. GPT-4V accuracy declines from
 53.6 percent on scenes with four objects to 43.1 percent on
 scenes with eight objects, a relative decrease of 19.6 per-

Model	Vision	+Structure	Gain
GPT-4V	48.3 (2.8)	52.1 (2.9)	+3.8
LLaVA-1.5	41.7 (3.1)	44.8 (3.2)	+3.1
InstructBLIP	36.2 (3.3)	38.9 (3.4)	+2.7
Qwen-VL-Chat	38.9 (3.2)	41.5 (3.3)	+2.6

Table 5. Spatial question answering accuracy on counterfactual variants under vision-only versus vision-plus-structure conditions. Structural information provides modest yet consistent improvements across all models, with gains ranging from 2.6 to 3.8 percentage points.

cent. This finding suggests that current architectures struggle to maintain accurate spatial representations as relational complexity increases, consistent with limitations in attention mechanisms when numerous objects compete for representational capacity.

5.5. Input Modality Ablation

Table 5 compares model performance under vision-only versus vision-plus-structure conditions, assessing whether explicit scene graph information enhances spatial reasoning.

Providing explicit scene graph information enhances performance across all models, with gains ranging from 2.6 to 3.8 percentage points. GPT-4V exhibits the largest improvement, attaining 52.1 percent accuracy with structural information. However, even with explicit coordinates and relationships provided in text, models fail to achieve robust performance, suggesting that the core limitation resides in reasoning about spatial information rather than perception alone. Notably, performance gains from structural input are modest despite providing exact object coordinates and support relations. This finding suggests that models may not fully exploit explicit geometric descriptions, instead treating them as additional text rather than structured constraints. Future analysis could probe whether attention weights align with graph-provided spatial information or whether models largely ignore structural annotations.

5.6. Statistical Significance

Bootstrap tests with 1,000 resamples confirm that differences between models are statistically significant. Comparing GPT-4V to LLaVA on counterfactual question accuracy, the mean difference of 6.6 percentage points exhibits a 95 percent confidence interval from 3.8 to 9.4 percentage points with $p < 0.001$. All pairwise comparisons between MLLMs are significant at $p < 0.05$ except the comparison between InstructBLIP and Qwen-VL-Chat where $p = 0.12$, indicating comparable performance between these two models. Comparisons between MLLMs and baselines are all significant at $p < 0.001$.

5.7. Error Analysis

We conduct manual analysis of 200 error cases per model, totaling 800 errors examined, to identify systematic failure patterns across architectures. Spatial term confusion accounts for 32 percent of errors, with models frequently confusing left and right when objects are viewed from different perspectives. For example, when the camera rotates by ninety degrees, models often maintain egocentric spatial judgments rather than allocentric reasoning about world coordinates, responding based on image coordinates rather than true spatial relationships. Occlusion blindness appears in 28 percent of errors involving partially visible objects. When objects are partially occluded, models report them as absent in 24 percent of cases and misidentify their category in 18 percent of cases. This finding suggests that models lack robust amodal completion capabilities and rely substantially on visible surface features for object recognition. Physics violation insensitivity proves particularly striking, with models describing floating or intersecting objects as normal in 59 percent of implausible scenes. Even when explicitly queried about physical possibility, models frequently fail to detect violations, suggesting that training data rarely includes negative examples of physical implausibility. Change detection failure in counterfactual comparison tasks occurs in 37 percent of cases, with models either reporting no change or identifying the incorrect object. This finding indicates that models lack explicit mechanisms for comparing visual inputs and detecting spatial differences. Counting errors account for 23 percent of mistakes, with object counting accuracy declining from 84 percent in base scenes to 61 percent in counterfactual variants even when object counts remain constant. Spatial perturbations appear to disrupt basic perceptual grouping and enumeration. Within counterfactual comparison errors, 41 percent correspond to incorrect object identification, 33 percent correspond to correct object identification but incorrect transformation type, and 26 percent correspond to reporting no change despite a perturbation. This distribution indicates that models struggle both with localizing the modified object and with categorizing the transformation nature, suggesting absence of explicit scene-differencing mechanisms.

5.8. Why Do MLLMs Fail on Counterfactual Spatial Reasoning?

The consistent degradation across perturbation types suggests that current models lack stable geometric representations that persist across controlled transformations. We hypothesize three contributing factors. First, token-based multimodal architectures discretize spatial information into visual tokens optimized for semantic recognition rather than continuous geometric reasoning. While such representations capture object identity and coarse relations, they do not explicitly encode metric structure, support relations, or

536 rigid transformations in a geometrically consistent manner.
537 Second, current MLLMs lack inductive biases for three-
538 dimensional equivariance. Rotations, occlusions, and object
539 relocations substantially alter pixel configurations, yet
540 the underlying world structure may remain partially invariant.
541 Without architectural mechanisms enforcing viewpoint
542 or transformation consistency, models rely on appearance-
543 level correlations rather than stable three-dimensional rela-
544 tional representations. Third, large-scale training corpora
545 contain abundant descriptions of typical physical scenes
546 but relatively few examples of counterfactual or physically
547 implausible configurations. Consequently, models acquire
548 strong priors over common spatial patterns but are not ex-
549 plicitly trained to detect violations of physical constraints or
550 to maintain consistent world-state tracking under controlled
551 perturbations. Together, these observations suggest that
552 enhancing spatial intelligence in MLLMs may require archi-
553 tectural integration of structured three-dimensional rep-
554 resentations, counterfactual training signals, or persistent
555 world models rather than purely scaling data and param-
556 eters. These limitations may reflect a representational bottle-
557 neck in current token-based multimodal architectures. Vi-
558 sual encoders convert continuous spatial geometry into dis-
559 crete token embeddings optimized for semantic discrimi-
560 nation rather than metric reasoning. As scene complex-
561 ity increases, attention mechanisms must distribute capac-
562 ity across more object tokens, potentially degrading rela-
563 tional precision. Without explicit three-dimensional induc-
564 tive biases or equivariant constraints, learned representa-
565 tions may fail to preserve geometric structure under spa-
566 tial perturbations. This finding suggests that architectural
567 changes—not merely larger datasets—may prove necessary
568 for robust spatial world modeling.

569 6. Discussion

570 The experimental results reveal several important findings
571 regarding current MLLM capabilities for spatial reasoning.
572 The substantial performance disparity between base and
573 counterfactual scenes indicates that models appear to rely
574 on appearance-level correlations rather than maintaining
575 transformation-consistent geometric representations. When
576 spatial configurations change in ways that should not affect
577 certain judgments, models frequently alter their responses,
578 suggesting that internal representations lack invariance un-
579 der relevant transformations. The particularly diminished
580 performance on physics violations suggests that current
581 training regimes inadequately expose models to physically
582 implausible configurations. Models learn to describe typ-
583 ical scenes but fail to detect when configurations violate
584 basic physical constraints. This limitation carries signif-
585 icant implications for deploying these models in embod-
586 ied artificial intelligence applications where detecting phys-
587 ical impossibility proves crucial for safe planning and ex-

588 ecution. Comparison with baseline methods confirms that
589 MLLMs substantially outperform traditional computer vi-
590 sion approaches on spatial reasoning, yet the performance
591 gap narrows on counterfactual tasks. The Mask R-CNN
592 baseline, despite its simplicity, achieves 31.5 percent accu-
593 racy on counterfactual scenes, demonstrating that basic geo-
594 metric heuristics provide a foundation upon which MLLMs
595 build but do not fundamentally transcend. We did not in-
596 clude a human baseline in this study. Informal inspection
597 suggests that most counterfactual and plausibility questions
598 are straightforward for human observers; however, sys-
599 tematic human evaluation would provide a valuable upper-
600 bound reference for future investigation. The modest im-
601 provements from explicit scene graph information suggest
602 that the core limitation resides in reasoning about spatial
603 relationships rather than perceiving them. Even when pro-
604 vided with exact coordinates and relations in text, models
605 struggle to maintain consistent representations and detect
606 implausible configurations. Several limitations of this study
607 warrant acknowledgment. The use of synthetic scenes,
608 while enabling controlled perturbation, may not fully cap-
609 ture the complexity of real-world environments with natural
610 textures, lighting variations, and object diversity. Perfor-
611 mance on real images could differ substantially from our
612 synthetic evaluation. The static nature of our scenes pre-
613 cludes evaluation of dynamic spatial reasoning involving
614 object motion or agent interaction. Our perturbation tax-
615 onomy, while systematic, may not encompass all forms of
616 counterfactual variation relevant to spatial intelligence, such
617 as material property changes or lighting transformations.
618 Future work would be to extend evaluation to dynamic
619 scenes with object motion, real-world images with known
620 three-dimensional ground truth captured via depth sensors
621 or multi-view reconstruction, additional perturbation types
622 including material and lighting variations, fine-tuned mod-
623 els trained on counterfactual data to assess whether these
624 limitations can be addressed through targeted training, and
625 embodied agents that must reason about counterfactuals
626 during planning and execution.

627 7. Conclusion

628 We introduced SCWM, a framework for evaluating multi-
629 modal spatial reasoning under controlled 3D perturbations.
630 Evaluating four MLLMs and two baselines on 2,500 base
631 scenes and 7,500 counterfactual variants reveals significant
632 limitations: performance drops 22.9 percentage points from
633 standard to counterfactual tasks, physics violation accuracy
634 reaches only 43.8% for the best model, consistency scores
635 remain below 0.6, and counterfactual comparison accuracy
636 falls below 40%. These findings demonstrate that scaling
637 multimodal architectures alone is insufficient explicit struc-
638 tural grounding, counterfactual supervision, or integration
639 of geometric world models may be necessary. 639

640

References

641

[1] OpenAI, G. P. T. (2023). 4V (ision) System Card [https://cdn.openai.com/papers.GPTV_System_Card.pdf](https://cdn.openai.com/papers/GPTV_System_Card.pdf). 1, 5

642

643

644

645

646

[2] Liu, H., Li, C., Wu, Q., & Lee, Y. J. (2023). Visual instruction tuning. *Advances in neural information processing systems*, 36, 34892-34916. 1, 5

647

648

649

650

651

[3] Dai, W., Li, J., Li, D., Tiong, A., Zhao, J., Wang, W., ... & Hoi, S. (2023). Instructblip: Towards general-purpose vision-language models with instruction tuning. *Advances in neural information processing systems*, 36, 49250-49267. 1, 5

652

653

654

655

656

[4] Xue, Q., Liu, W., Wang, S., Wang, H., Wu, Y., & Gao, W. (2025). Reasoning Path and Latent State Analysis for Multi-view Visual Spatial Reasoning: A Cognitive Science Perspective. *arXiv preprint arXiv:2512.02340*. 1

657

658

659

660

661

662

[5] Johnson, J., Hariharan, B., Van Der Maaten, L., Fei-Fei, L., Lawrence Zitnick, C., & Girshick, R. (2017). Clevr: A diagnostic dataset for compositional language and elementary visual reasoning. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 2901-2910). 1, 2

663

664

665

666

[6] Das, A., Datta, S., Gkioxari, G., Lee, S., Parikh, D., & Batra, D. (2018). Embodied question answering. In *Proceedings of the IEEE conference on computer vision and pattern recognition* (pp. 1-10). 1, 2

667

668

669

670

[7] Goyal, Y., Wu, Z., Ernst, J., Batra, D., Parikh, D., & Lee, S. (2019, May). Counterfactual visual explanations. In *International Conference on Machine Learning* (pp. 2376-2384). PMLR. 2

671

672

673

674

[8] Bakhtin, A., van der Maaten, L., Johnson, J., Gustafson, L., & Girshick, R. (2019). Phyre: A new benchmark for physical reasoning. *Advances in Neural Information Processing Systems*, 32. 2

675

676

677

678

679

680

[9] Chen, B., Xu, Z., Kirmani, S., Ichter, B., Sadigh, D., Guibas, L., & Xia, F. (2024). Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition* (pp. 14455-14465). 2

681

682

683

[10] He, K., Gkioxari, G., Dollár, P., & Girshick, R. (2017). Mask r-cnn. In *Proceedings of the IEEE international conference on computer vision* (pp. 2961-2969). 4

684

685

686

687

688

[11] Radford, A., Kim, J. W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., ... & Sutskever, I. (2021, July). Learning transferable visual models from natural language supervision. In *International conference on machine learning* (pp. 8748-8763). PmLR. 4

689

690

691

[12] Bai, J., Bai, S., Chu, Y., Cui, Z., Dang, K., Deng, X., ... & Zhu, T. (2023). Qwen technical report. *arXiv preprint arXiv:2309.16609*. 5