

Imaginative Perception Tokens Enhance Spatial Reasoning in Multimodal Language Models

Anonymous CVPR submission

Paper ID ****

Abstract

001 Vision-language models (VLMs) excel at many tasks, yet
002 continue to struggle with spatial reasoning—problems
003 where the key information is not directly observable in the
004 input. Many spatial questions require imaginative percep-
005 tion: simulating an unseen viewpoint, tracing a trajectory
006 through an occluded space, or integrating partial views
007 into a coherent spatial map. Humans naturally support
008 this kind of reasoning through imagination. Prior work
009 has introduced intermediate visual representations (e.g., vi-
010 sual thoughts, depth, or box tokens), but these intermediates
011 often refine structure already visible rather than predict-
012 ing the missing spatial structure implied by the evidence.
013 We introduce **Imaginative Perception Tokens (IPT)**, inter-
014 mediate perceptual representations that externalize what a
015 VLM would perceive under an alternative spatial configura-
016 tion while remaining consistent with the observed input. To
017 study this capability, we formulate three tasks that require
018 imaginative perception: **Perspective Taking (PET)**, **Path**
019 **Tracing (PT)**, and **Multiview Counting (MVC)**. For each
020 task, we construct datasets of $\sim 20K$ examples spanning
021 simulated and real-world settings, paired with ground-truth
022 intermediate imaginations, final answers, and curated eval-
023 uation benchmarks. Using the unified VLM BAGEL [12] as
024 our backbone, IPT supervision improves spatial reasoning
025 across several settings and often outperforms textual chain-
026 of-thought training, even when no image is generated at in-
027 ference time. For example, on MVC, IPT improves accuracy
028 by 3.4% and achieves performance competitive with strong
029 closed-source models on Path Tracing. We also find that
030 mixed training with IPT and label-only data can further im-
031 prove performance. In contrast, textual chain-of-thought
032 can be detrimental on these tasks, substantially degrading
033 performance in some cases, highlighting a modality mis-
034 match when forcing spatial computation through language.
035 Overall, IPT provides a principled supervision signal for
036 reasoning over unobserved structure, yielding stronger spa-
037 tial generalization and a more interpretable intermediate

aligned with the underlying geometry of the task.

038

1. Introduction

039

040 Spatial reasoning still remains a persistent challenge
041 for vision-language models (VLMs) [1, 8, 11]. Many spa-
042 tial questions require reasoning about how objects relate
043 within a three-dimensional environment, how these rela-
044 tionships change under viewpoint transformations, or how
045 information from multiple partial observations should be in-
046 tegrated into a coherent scene representation [17, 41] for
047 vision-language models (VLMs) [1, 8, 11]. While current
048 models can often recognize objects and attributes, they fre-
049 quently struggle when reasoning requires manipulating spa-
050 tial structure, such as predicting how a scene would appear
051 from another viewpoint [21, 23] or aggregating information
052 across multiple views [37].

053 A key reason for this difficulty is that many spatial rea-
054 soning problems cannot be solved by analyzing the input
055 alone. Instead, they require constructing a spatial repre-
056 sentation that is not directly observed. Humans naturally
057 address such problems through imagination: when asked
058 what lies to the left after moving to a new position, or
059 how many objects exist in a room seen from several view-
060 points, we mentally simulate the scene from unseen per-
061 spectives or integrate partial observations into a unified spa-
062 tial map [34, 41, 42]. In other words, spatial reasoning often
063 depends on imagining missing spatial structure that proceed
064 despite incomplete observations.

065 Existing approaches provide only partial solutions. Re-
066 cent work teaches models to generate intermediate visual
067 thoughts alongside language [14, 16, 20], while others in-
068 troduce structured perceptual intermediates, such as depth
069 maps or bounding boxes represented as tokens [2, 27, 38].
070 Although these methods demonstrate that intermediate vi-
071 sual representations can support reasoning, they primarily
072 operate over information already present in the input obser-
073 vation, refining visible structures or extracting perceptual
074 attributes. However, as discussed above, many spatial rea-

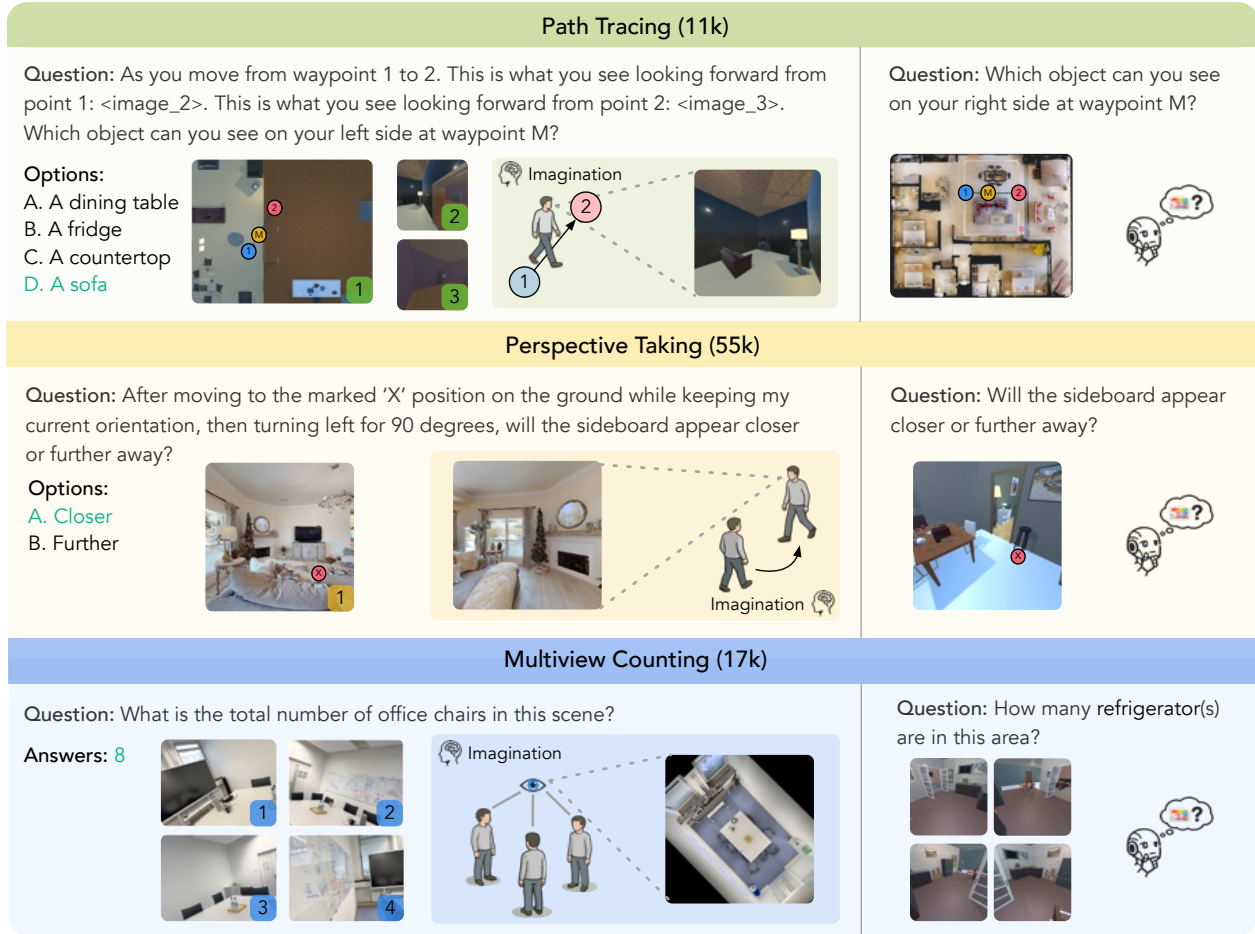


Figure 1. **Overview of the three spatial imagination tasks.** The left columns show training examples with ground-truth imaginative perception; the right columns show evaluation examples.

075 soning problems arise precisely because the required spa-
076 tial information is not directly observable, and therefore re-
077 quires imagination.

078 To address this gap, we propose **Imaginative Perceptual**
079 **Tokens** for VLMs. When VLMs are trained with them, they
080 enable intermediate reasoning steps that represent novel
081 spatial views. Unlike standard perceptual intermediates that
082 describe structures visible in the input, imaginative repre-
083 sentations correspond to what the model would perceive if
084 it were observing the input from a different spatial confi-
085 guration, such as from an unseen viewpoint or after integrat-
086 ing multiple partial observations into one. At the same time,
087 they are not unconstrained imagination: the predicted per-
088 cept must remain consistent with the observed scene. These
089 tokens externalize the model’s prediction of what would be
090 perceived given incomplete spatial evidence.

091 To study this capability, we propose three spatial rea-
092 soning tasks that fundamentally require imaginative percep-
093 tion. (1) Perspective Taking requires predicting how a scene

would appear from a new viewpoint given a single first-
person observation (“If you move to the marked position
and turn left, will the chair appear on your left or right?”);
(2) Path Tracing requires inferring what an agent would see
along a navigation path based on a top-down view (“If you
walk along the marked path, which object will you see on
your side?”); Finally, (3) Multiview Counting requires inte-
grating multiple partial observations into a top-down view
to determine the number of objects present in the scene.
These tasks would be made easy when correctly predicting
what would be perceived in a different spatial configura-
tion. For each task we construct a dataset of approximately
20k examples each drawn from both real-world and syn-
thetic simulated environments, with ground-truth interme-
diate spatial imaginations paired with final answers. Each
dataset is accompanied by a human-filtered benchmark for
evaluation. Together these constitute the first datasets de-
signed explicitly to train and evaluate visually-grounded in-
termediate spatial reasoning in models.

094
095
096
097
098
099
100
101
102
103
104
105
106
107
108
109
110
111
112

113 Empirically, we find that training with imaginative perceptual supervision can improve performance on these spatial reasoning tasks compared to answer-only supervision, and often compares favorably to textual chain-of-thought approaches. These improvements can persist even when the model does not explicitly generate intermediate images at inference time, suggesting that such supervision may help models develop stronger internal spatial representations. At the same time, we observe that the benefits vary across tasks and settings, indicating that imagination quality and task structure both play important roles.

124 Overall, our results suggest that supervising models with intermediate perceptual predictions offers a useful direction for improving spatial reasoning, particularly in settings where the required structure is not directly observable from the input.

129 2. Related Works

130 **Evaluation of VLMs’ spatial reasoning.** A growing body of benchmarks has established that spatial reasoning remains a persistent weakness of modern vision-language models. Early datasets target brittleness in basic spatial predicates: SpatialSense [35] reduces language priors through adversarial crowdsourcing, while VSR [22] scales relation types in a caption-verification format, and What’sUp [17] uses minimal-pair testing to reveal systematic failures on left/right and above/below distinctions. More recent work shifts from 2D relations to viewpoint and 3D structure. 3DSRBench [23] shows that models fail under modest changes in perspective, depth, and occlusion, and ViewSpatial-Bench [21] identifies a “perspective gap”: models often succeed in camera-centered views but break when asked to adopt human-centered viewpoints.

145 Benchmarks have also expanded to multi-image and video settings where maintaining a consistent spatial state is essential. VSI-Bench [34] tests whether models can build a persistent mental map from videos, while MMSI-Bench [37] reports large human–model gaps on cross-view scene reconstruction. MindCube [41] is closely aligned with our motivation, targeting spatial mental modeling from limited views, including perspective-taking and “what-if” scene dynamics. Counting Stacked Objects [13] studies 3D object counting under heavy occlusion across multiple views, directly analogous to our Multiview Counting setting. Finally, benchmark design work emphasizes that shortcuts remain pervasive: Brown *et al.* [3] construct VSI-Bench-Debiased by iteratively pruning samples solvable via priors, reinforcing the need for evaluations where success requires genuine spatial computation.

161 Collectively, these benchmarks diagnose *where* VLMs fail spatially, but they typically evaluate *discriminative* understanding—reading off a relation from an observed view—rather than *constructive* spatial imagination. Our

work is complementary: we isolate *imaginative perception* as a standardized intermediate substrate, and pair each task with a ground-truth intermediate spatial imagination rather than only a final answer label.

Intermediate representations for spatial reasoning. Chain-of-thought prompting [31] can improve multi-step reasoning, but serializing viewpoint transformations, occlusions, and geometric constraints into language is often awkward and error-prone. This motivates intermediate representations in modalities better aligned with spatial computation. One direction externalizes reasoning into explicit visual buffers: Visual Sketchpad [16] equips models with drawing actions for iterative refinement, and MVoT [20] introduces visualization-of-thought traces that help on dynamic spatial tasks where text CoT struggles. ThinkMorph [14] studies interleaved text–image reasoning traces, and OpenAI describes o3/o4-mini as using chains-of-thought that include simple image transformations during reasoning [25]. A complementary line introduces latent visual scratchpads: Mirage [38] frames latent tokens as “machine mental imagery,” and Mull-Tokens [27] generalizes to modality-agnostic latent thinking tokens.

Our work differs in what the intermediate is meant to represent. Many prior approaches treat intermediate images or latents as optional visualizations of *visible* structure. We instead target *imaginative perception*: predicting what would be perceived under an unobserved spatial configuration (e.g., a rotated viewpoint or a top-down path state), a representation constrained by the input but not present in it. This framing provides a principled criterion for when intermediate visual thoughts are necessary and a controlled way to supervise them.

Unified multimodal models for interleaved understanding and generation. Producing imaginative perceptual intermediates within a single model requires the ability to both understand and generate images. Unified decoder-only architectures treat image tokens as first-class sequence elements, enabling arbitrary text–image interleaving. Chameleon [5] is an early example, while Show-o2 [33] and Janus [7] offer alternative unified designs that balance understanding and generation. We build on BAGEL [12], a unified model pretrained on large interleaved corpora that exhibits strong spatial capabilities, making it a natural substrate for producing intermediate spatial imaginations. Crucially, however, a unified architecture alone does not guarantee that intermediate images are *used* in a way that supports reasoning; our work provides task constructions and supervision that make imaginative perception the relevant computational substrate.

3. Spatial Imagination: Tasks and Datasets

We introduce three spatial reasoning tasks that require constructing a missing spatial representation from incomplete

Table 1. **Dataset and benchmark statistics.** All experiments use the AI2-THOR subset for training; additional data sources are released for future research. †: human-verified subset.

Task	Source (samples)	AI2-THOR (20,531) + Habitat (19,998) + Real images (15,000)
Perspective Taking	Train	55,529
	Eval	AI2-THOR† (238), Habitat† (300)
	IPT Format	Novel-viewpoint image
Path Tracing	Source (samples)	AI2-THOR (11,204)
	Train	11,204
	Eval	AI2-THOR† (329), Real† (332)
IPT Format	Sideview image	
Multiview Counting	Source (samples)	AI2-THOR (17,079) + MessyTable (1,880) + ScanNet (540)
	Train	19,499
	Eval	AI2-THOR† (260)
IPT Format	Top-down BEV map	

217 inputs (single-view, partial-view, or map inputs). For each
218 task, we build a 10k–50k training set with paired *ground-*
219 *truth spatial imaginations* (task-specific intermediate visual
220 supervision) and final answers, and we release a human-
221 filtered benchmark for controlled evaluation. All datasets
222 will be released publicly; for consistency across tasks, we
223 train our models on the AI2-THOR [19] subset of each
224 training set. Table 1 and Fig 1 summarize the training data
225 and evaluation benchmarks.

226 3.1. Perspective Taking

227 Given a first-person view of an indoor scene with target po-
228 sitions marked, the model must answer a spatial question
229 (e.g., “After moving to ‘X’ and turning left 90°, will the
230 {object} be on your left or right?”) about the scene from
231 the new viewpoint. Since the target view is never provided,
232 the model must mentally simulate the spatial transformation
233 rather than read off the answer directly.

234 **Sub-categories.** Questions span two spatial relation types
235 across six balanced sub-categories. *Distance change* asks
236 whether a target object becomes closer or further after the
237 viewpoint shift: (1) *closer* and (2) *further*. *Relative po-*
238 *sition* asks whether the object falls to the left or right in
239 the new view, defined by the object’s lateral position before
240 and after the transformation: (3) *left→left*; (4) *left→right*;
241 (5) *right→left*; (6) *right→right*. Overall accuracy is the un-
242 weighted mean across all six sub-categories so that each
243 spatial relationship contributes equally, preventing models
244 from gaming the metric by over-predicting common cases.

245 **Imaginative perception target.** A novel-viewpoint render-
246 ing of the scene from the target position, directly supervised
247 against ground-truth renders from the 3D scene.

248 **Data.** Synthetic data is generated from AI2-THOR [19]
249 and Habitat [24, 26, 29] by sampling source/target cam-
250 era pairs, rendering first-person views, and annotating the
251 source view with a red “X” marking the target. Questions
252 cover two relation types (distance change, relative posi-
253 tion) across six balanced sub-categories. A *mixed* training
254 data variant additionally incorporates real-world examples
255 from the Visual Spatial Tuning dataset [36] (camera mo-
256 tion subset) as a synthetic-to-real bridge. The base training

set contains 20,531 AI2-THOR examples; the mixed vari-
ant totals 55,529. We evaluate on held-out human-verified
AI2-THOR (238) and Habitat (300) benchmarks. Full sub-
category breakdowns and data generation details are pro-
vided in the Appendix.

3.2. Path Tracing

Given a top-down map with a marked path $1 \rightarrow 2$, a mid-
point M_1 , and egocentric forward views at waypoints 1
and 2, the model must identify which object is visible on
a queried side at M_1 . Neither the top-down map nor the
endpoint views reveal first-person visibility at the midpoint,
requiring the model to imagine what the agent would see
from ground level.

We evaluate under three input settings of increasing spa-
tial cues: *Path* (map only), *PathArr* (map + query direction
arrow), and *EgoDir* (map + egocentric endpoint views).

Imaginative perception target. A sideview image — a
first-person rendering from M_1 — that externalizes the
3D visibility reasoning the top-down input cannot support.
Ground-truth sideviews are rendered directly from the sim-
ulator at M_1 .

Data. Synthetic data is generated from AI2-THOR [19] and
ProcTHOR [10], sampling feasible two-waypoint paths bal-
anced across room types and distance bins. Questions are
template-generated with four answer choices and quality-
filtered via TIFA-style verification [15] using GPT-4.1 ma-
jority voting; samples answerable from endpoint views
alone are removed to ensure genuine imagination is re-
quired. The synthetic training set contains 11,204 exam-
ples. A real-world benchmark of 332 human-verified ques-
tions is constructed from Matterport3D [6] top-down views
and evaluated on Path and PathArr settings only. Full filter-
ing criteria and real-world annotation pipeline details are in
the Appendix.

3.3. Multiview Counting

Given several first-person frames of the same environment,
the model must select the correct count of a queried object
(e.g., “How many chairs are in this area?”). Since no sin-
gle view reveals the full layout, and the same object often
appears across multiple frames, the model must construct a
unified spatial representation that resolves both occlusions
and cross-view duplicates.

Imaginative perception target. A top-down bird’s-eye
view (BEV) map aggregating all input views, making de-
duplication explicit by mapping each object to a single spa-
tial location. Ground-truth BEV maps are rendered from an
overhead camera in the 3D scene.

Data. Synthetic examples are generated via multi-camera
and rotation trajectory types. Real-world data is sourced
from MessyTable [4] (fixed multi-camera rig; overhead im-
age as ground-truth BEV) and ScanNet++ [40] (point-cloud

308 BEV maps converted to photorealistic overhead images via
 309 Qwen Edit [32]). Questions are four-choice MCQ with dis-
 310 tractors sampled near the true count. The base training set
 311 contains 17,079 synthetic examples; the mixed variant to-
 312 tals 19,499. We evaluate on a human-verified benchmark of
 313 260 samples. Details on trajectory types, BEV rendering,
 314 and distractor sampling are in the Appendix.

315 4. Method: Imaginative Perception Tokens

316 The core of our approach is to enable Multimodal Language
 317 Models (MLLMs) to externalize spatial reasoning through
 318 **Imaginative Perception Tokens**. Unlike standard textual
 319 chain-of-thought or methods outsourcing visual imagina-
 320 tion with an external visual generation model, our method
 321 requires the model to generate a visual representation of
 322 a *non-observed* spatial configuration—such as an unseen
 323 viewpoint or an integrated top-down map—as a functional
 324 prerequisite for answering a spatial query.

325 4.1. Problem Formalization

326 Given an input context \mathcal{C} consisting of one or more ob-
 327 served images $\mathcal{I}_{obs} = \{I_1, \dots, I_k\}$ and a spatial lan-
 328 guage query Q , the goal is to predict the correct answer
 329 A . We decompose this into a two-stage generative pro-
 330 cess. First, the model generates **imaginative perception**
 331 **tokens** \hat{I}_{imag} , representing the implied spatial structure re-
 332 quested by the task (e.g., the view from a new coordinate):
 333 $P(\hat{I}_{imag} | \mathcal{I}_{obs}, Q)$ Second, the conditioned on this imagina-
 334 tive perception tokens \hat{I}_{imag} , the model produces the final
 335 answer: $P(A | \mathcal{I}_{obs}, Q, \hat{I}_{imag})$.

336 4.2. Architecture

337 We implement this approach using BAGEL [12], a uni-
 338 fied decoder-only transformer that natively supports inter-
 339 leaved multimodal understanding and generation. BAGEL
 340 employs a Mixture-of-Transformer-Experts (MoT) design:
 341 the model utilizes two transformer experts, one optimized
 342 for multimodal understanding and another for generation.
 343 Both operate on the same token sequence through shared
 344 self-attention at every layer. Images are represented via two
 345 distinct paths. *Understanding tokens* (U) are extracted via
 346 a SigLIP2 [30] ViT encoder to capture semantic content,
 347 while *Generation tokens* (G) are latent representations from
 348 a FLUX VAE used for high-fidelity synthesis. Because all
 349 tokens (text, U , and G) coexist in a single shared context
 350 window, the model maintains lossless interaction between
 351 understanding and generation modules.

352 While BAGEL’s standard generation tokens are typically
 353 used for open-ended text-to-image generation or editing, we
 354 repurpose this generative capacity for **spatial reasoning**. In
 355 our framework, the generation target is not a stylistic output
 356 but a precise **view imagination**—a visually grounded inter-

mediate that represents the unobserved 3D structure of the
 357 scene. 358

4.3. Training and Inference 359

Training Objective. We optimize the framework using a
 360 multi-task loss $\mathcal{L}_{total} = \lambda_{fm}\mathcal{L}_{fm} + \lambda_{lm}\mathcal{L}_{lm}$. The model
 361 is trained to jointly produce the imaginative perception and
 362 the final answer: 363

- 364 **1. Flow-Matching Loss (\mathcal{L}_{fm}):** For the imaginative in-
 365 termediate, BAGEL adopts the **Rectified Flow** method.
 366 The model learns to predict the velocity field v_t required
 367 to transform Gaussian noise into the target latent G_{gt}
 368 representing the unobserved view, conditioned on the
 369 preceding context \mathcal{C} : 370

$$\mathcal{L}_{fm} = \mathbb{E}_{t, G_0, \mathcal{C}} [\|v_t(G_t | \mathcal{C}) - (G_{gt} - G_0)\|^2] \quad (1) \quad 370$$

- 371 **2. Language Modeling Loss (\mathcal{L}_{lm}):** We minimize the neg-
 372 ative log-likelihood of the final VQA answer tokens A ,
 373 conditioned on the observed context and the ground-
 374 truth imaginative tokens: 374

$$\mathcal{L}_{lm} = - \sum_{i=1}^{|A|} \log P(a_i | \mathcal{C}, U_{gt}, G_{gt}, a_{<i}) \quad (2) \quad 375$$

Inference. At inference time, the model operates in one
 376 of two modes depending on the task and configuration. 377
 378 In the **text-only** mode, the model produces only a tex-
 379 tual answer without generating any visual intermediate
 380 $A \sim P(A | \mathcal{C})$, serving as a baseline. In the **imag-**
 381 **ination** mode, the model first performs iterative denois-
 382 ing over VAE tokens to produce the imaginative latent:
 383 $\hat{G}_{imag} = \int_0^1 v_t(G_t | \mathcal{C}) dt$ The decoded image \hat{I}_{imag} is
 384 immediately re-encoded and appended to the context as
 385 both ViT understanding tokens and VAE generation tokens:
 386 $\mathcal{C}' = [\mathcal{C}, \text{ViT}(\hat{I}_{imag}), \text{VAE}(\hat{I}_{imag})]$ The model then at-
 387 tends to its own imagination to predict the final answer
 388 $A \sim P(A | \mathcal{C}')$.

5. Experiments 389

390 We evaluate *imaginative perception tokens* on the three
 391 spatial reasoning tasks introduced in Sec. 3: **Perspec-**
 392 **tive Taking (PET)**, **Path Tracing (PT)**, and **Multiview**
 393 **Counting (MVC)**. To enable controlled comparisons, we
 394 train all task-specific models on the AI2-THOR subset of
 395 each dataset. We additionally report transfer to cross-
 396 environment benchmarks (Habitat), real-world images, and
 397 external datasets. All tasks use multiple-choice evaluation
 398 with balanced answer distributions. 398

399 PT is evaluated under three input variants that provide
 400 increasing spatial cues: **EgoDir** (egocentric direction only),
 401 **Path** (top-down path overlay), and **PathArr** (path with di-
 402 rectional arrows) and average accuracy reported. Unless 402

Table 2. **Main results.** Accuracy (%) on AI2-THOR (in-domain) and different-environment (out-of-domain) benchmarks. PT reports the average across input settings (EgoDir/Path/PathArr for AI2-THOR; Real/Real+Arr for different environments). Text CoT generates a textual chain-of-thought before answering. IPT (Imaginative Perception Token) generates an intermediate image before answering. For our models, accuracy reports the **maximum** between answer-only and free-generation inference. Best per group in **bold**.

Model	AI2-THOR			Different Env.	
	PET	PT	MVC	PET	PT
<i>VQA Models</i>					
GPT-5	79.8	60.2	53.5	69.3	80.9
GPT-5.2	45.5	32.9	44.2	54.0	63.0
Gemini 2.5 Flash	51.0	41.5	30.8	66.3	71.4
Gemini 3 Flash	55.0	42.3	56.9	51.3	83.2
InternVL3.5-8B	51.5	35.8	44.6	47.7	47.4
Qwen2.5-VL-7B	50.7	37.3	38.8	54.3	44.8
Qwen3-VL-8B	52.0	35.9	43.8	46.7	64.1
<i>Unified Models</i>					
Janus-Pro-7B	51.8	33.5	33.1	44.7	35.3
Chameleon 7B	34.3	16.3	5.4	47.3	24.5
<i>Ours (fine-tuned BAGEL)</i>					
Bagel (base)	40.3	29.9	35.4	62.7	42.7
Bagel (label-only)	97.5	65.7	63.9	82.0	54.7
+ Text CoT	83.1	49.7	62.3	70.3	52.2
+ IPT	96.8	49.0	67.3	87.0	57.5
+ Mixed Training	97.8	66.7	62.3	87.7	58.6

403 otherwise stated, we report accuracy (%) and use the same
404 prompt formatting across baselines and our models.

405 5.1. Setup

406 **Baselines.** We compare against two groups of models,
407 evaluated zero-shot with task-specific prompts. *VQA mod-*
408 *els* include GPT-5, GPT-5.2, Gemini 2.5 Flash, Gemini 3
409 Flash, InternVL3.5-8B, Qwen2.5-VL-7B, and Qwen3-VL-
410 8B. *Unified models* that support both understanding and
411 generation include Janus-Pro-7B and Chameleon 7B.

412 **Our model variants.** We fine-tune BAGEL [12] under sev-
413 eral configurations to isolate the contribution of imagina-
414 tion supervision. Each fine-tuned model is task-specific and
415 trained on a single task using AI2-THOR data only (unless
416 noted otherwise):

- 417 • **Bagel (base):** pretrained model with no task-specific fine-
418 tuning.
- 419 • **Bagel (label-only):** fine-tuned with answer supervision
420 only, with no intermediate thought.
- 421 • **+ Text CoT:** trained to generate a textual chain-of-
422 thought describing the imagined spatial configuration be-
423 fore answering. Training CoTs are generated by GPT-5.1
424 using simulator ground-truth scene metadata.
- 425 • **+ IPT:** trained to generate an intermediate image (the
426 imaginative perception token) before answering.
- 427 • **+ Mixed Training:** trained on a mixture of IPT exam-
428 ples (image-generation targets) and label-only examples

(answer supervision only).

Training details. We fine-tune BAGEL-7B-MoT with
AdamW ($lr 1 \times 10^{-5}$, 2,000 warmup steps) on 8 GPUs using
FSDP bf16, following BAGEL [12] and ThinkMorph [14].
For multi-image inputs, each image is resized to 512×512 .
Unless noted, IPTs use Latent-64 resolution.

5.2. Main results

Table 2 reports results on our benchmarks.

Spatial reasoning remains difficult for current VLM and unified models.

Among the zero-shot baselines, GPT-5 is the strongest
across nearly all settings, yet still trails our best fine-tuned
variants on multiple in-distribution tasks. Smaller open
VLM models (InternVL3.5-8B, Qwen2.5-VL-7B, Qwen3-
VL-8B) hover near chance on PET (50–52%) and struggle
on PT, indicating that these tasks are not solvable through
superficial cues. Unified models perform worse overall:
Chameleon 7B drops to 34.3% on PET and 5.4% on MVC,
suggesting that current unified designs often trade away un-
derstanding robustness in exchange for generation capabil-
ity.

Answer supervision alone yields large gains and trans- fers across environments.

Bagel (label-only) substantially improves over Bagel
(base) across all tasks, rising from 40.3% to 97.5% on AI2-
THOR PET, from 29.9% to 65.7% on PT, and from 35.4%

Table 3. **Ablation on latent size.** Accuracy (%) with w/ Thought inference mode at different imagination resolutions. Best per column in **bold**.

Latent Size	Resolution	PET		MVC
		AI2-THOR	Habitat	AI2-THOR
Latent-4	64 × 64	87.4	73.3	53.5
Latent-16	256 × 256	95.3	81.0	56.2
Latent-32	512 × 512	95.0	87.0	58.9
Latent-64	1024 × 1024	96.8	83.3	63.1

to 63.9% on MVC. These improvements transfer: label-only reaches 82.0% on Habitat PET, showing that spatial reasoning can be learned in simulation and generalized to new environments.

Imagination supervision helps most when language is a poor interface.

On MVC, IPT achieves the best accuracy (67.3%), outperforming label-only (63.9%) and Text CoT (62.3%). On different-environment PET (Habitat), IPT reaches 87.0% (vs. 82.0% for label-only), and Mixed Training improves further to 87.7%. On PT, Mixed Training achieves the best results on both synthetic (66.7%) and real (58.6%) benchmarks, outperforming label-only (65.7% / 54.7%) and all baselines. IPT also improves real-world PT transfer (57.5%) over label-only (54.7%) and Text CoT (52.2%). Notably, IPT models are evaluated in *answer-only* mode: the model does not generate an image at inference, yet the imagination targets during training strengthen internal spatial representations that transfer across environments.

Text CoT underperforms label-only and IPT.

Text CoT typically falls behind label-only (e.g., PET 83.1% vs. 97.5%, PT 49.7% vs. 65.7%) and also behind IPT (e.g., MVC 62.3% vs. 67.3%, PET 83.1% vs. 96.8%). Compared to label-only, the Text CoT objective forces the model to allocate capacity to generating long spatial descriptions during fine-tuning, which competes with answer prediction. Compared to IPT, the gap reflects a modality mismatch: viewpoint changes, occlusions, and cross-view correspondences are difficult to serialize into natural language, and the resulting textual traces introduce noise rather than useful structure. IPT represents these relationships directly in the visual modality where they are naturally expressed.

5.3. Ablations

Latent resolution controls imagination quality and downstream accuracy.

Table 3 and Fig. 2 ablate IPT resolution on PET and MVC. At Latent-4 (64 × 64), imaginations are blurry and lose spatial detail; at Latent-64 (1024 × 1024), imaginations become sharper and more spatially faithful, preserving object identities and relative positions. Quantitatively, increasing resolution from Latent-4 to Latent-64 improves AI2-

Table 4. **Ablation on thought modality and inference mode.** Accuracy (%) on AI2-THOR benchmarks (PT uses EgoDir variant). We compare Text CoT vs. IPT training and vary inference mode: generate thought then answer (w/ text/image), answer directly (answer-only), or condition on ground-truth (w/ GT image).

Training	Inference	PET	PT	MVC
Text CoT	w/ text	83.1	53.1	61.5
Text CoT	answer-only	78.3	55.8	62.3
IPT	w/ image	96.8	50.4	63.1
IPT	answer-only	96.8	61.1	62.3
IPT	w/ GT image	96.7	86.7	67.3

THOR PET from 87.4% to 96.8% and MVC from 53.5% to 63.1%. Habitat PET peaks at Latent-32 (87.0%) and drops slightly at Latent-64 (83.3%), suggesting mild overfitting to AI2-THOR appearance statistics at the highest resolution.

Thought modality and inference mode. Table 4 ablates the training signal (Text CoT vs. IPT) and inference mode (generate thought, answer-only, or oracle GT).

IPT training builds stronger spatial representations than Text CoT.

On PT, IPT with answer-only inference (61.1%) outperforms Text CoT with answer-only inference (55.8%) by 5.3 points. On MVC, IPT with image generation (63.1%) outperforms Text CoT with text generation (61.5%).

Imagination supervision is useful, but explicit generation is not required at inference.

For IPT models, answer-only mostly outperforms generating the imagination explicitly: on PT, answer-only reaches 61.1% vs. 50.4% with generation. For Text CoT, generating the chain-of-thought also slightly underperforms answer-only (53.1% vs. 55.8% on PT), though the gap is smaller than for IPT. This asymmetry suggests that producing faithful imaginations is harder than producing text descriptions, and imperfect generations can mislead downstream reasoning. However, training with imagination targets remains valuable: answer-only IPT matches GPT-5 on PT (61.1%).

Ground-truth imaginations reveal headroom.

When given ground-truth imaginations instead of model-generated ones, PT accuracy jumps from 50.4% to 86.7% (+36.3) and MVC rises from 63.1% to 67.3% (+4.2). The large PT gap indicates that imagination quality is the dominant bottleneck for path tracing; for PET, model-generated imaginations nearly match GT (96.8% vs. 96.7%), leaving little room for improvement.

IPT transfers to aligned external benchmarks.

Table 5 evaluates transfer to external benchmarks that test similar spatial capabilities: SAT [28] (perspective-taking subset) and MessyTable [4] (multiview counting). On SAT, Bagel (label-only) improves from 34.9% to 59.1%



Figure 2. Qualitative examples of model-generated imaginative perception tokens. Top two rows: MVC example showing imagined top-down BEV maps. Bottom: PET examples showing imagined novel viewpoints. From left to right, imagination resolution increases from Latent-4 (64×64) to Latent-64 (1024×1024). Higher resolution produces sharper and more spatially faithful imaginations, preserving object identities and relative positions needed for downstream reasoning.

Table 5. **Transfer to similar external benchmarks.** Accuracy (%). SAT tests perspective taking and MessyTable tests multiview counting, both in domains unseen during training. Best in **bold**.

Model	PET	MVC
	SAT (66)	MessyTable (200)
Bagel (base)	34.9	29.0
Bagel (label-only)	59.1	32.5
+ Text CoT	50.0	30.0
+ IPT	57.6	28.5
+ Mixed Training	63.6	37.0

529 over Bagel (base), and Mixed Training further improves to
530 63.6%. On MessyTable, Mixed Training reaches 37.0%, up
531 from 29.0% for Bagel (base).

Training with our data improves performance on other spatial benchmarks.

532 Finally, we test whether our training data improves spa-
533 tial reasoning on tasks with different structures: Scan-
534 Net [9] (in-the-wild multiview counting), MindCube [41]
535 (abstract geometric reasoning), and All-Angles-Bench [39]
536 (cross-view matching on EgoHumans [18]). Because IPTs
537 are task-specific by construction (e.g., rotated views for
538 PET, bird’s-eye paths for PT), they do not directly trans-
539 fer to these settings. We therefore fine-tune on AI2-THOR
540 MVC using answer supervision only. Bagel (fine-tuned)
541 consistently improves over Bagel (base) across all three
542 benchmarks (40.5%→52.0% on ScanNet, 39.5%→47.5%
543 on MindCube, 40.0%→50.0% on All-Angles), indicating
544 that our simulator data builds broadly useful spatial repre-
545 sentations even when the specific imaginative token target
546 changes.
547

Table 6. **Does our data help on other spatial tasks?** Accuracy (%) on benchmarks beyond our training task categories. Fine-tuning on our AI2-THOR MVC data consistently improves over Bagel (base), suggesting that the spatial reasoning learned from our datasets transfers broadly. Best per column in **bold**.

Model	ScanNet (200)	MindCube (200)	All-Angles (170)
<i>VQA Models</i>			
GPT-5	58.5	67.3	67.9
GPT-5.2	48.5	37.5	29.4
Gemini 2.5 Flash	48.0	50.3	37.9
Gemini 3 Flash	62.5	56.5	64.2
InternVL3.5-8B	53.5	42.1	54.8
Qwen2.5-VL-7B	63.5	47.8	51.8
Qwen3-VL-8B	62.5	34.5	42.3
<i>Unified Models</i>			
Janus-Pro-7B	39.5	42.0	45.0
Chameleon 7B	5.5	25.4	17.7
<i>Ours (fine-tuned BAGEL)</i>			
Bagel (base)	40.5	39.5	40.0
Bagel (fine-tuned)	52.0	47.5	50.0

6. Conclusion

548 We introduced Imaginative Perception Tokens (IPTs), inter-
549 mediate visual representations that externalize spatial rea-
550 soning about unobserved structure in multimodal language
551 models, and designed three tasks: Perspective Taking, Path
552 Tracing, and Multiview Counting, with ground-truth inter-
553 mediate imaginations. Training with imagination supervi-
554 sion consistently outperforms label-only and text chain-of-
555 thought baselines, even without explicit imagination at in-
556 ference, and ablations confirm that imagination quality di-
557 rectly governs downstream accuracy.
558

559
560
561
562
563
564
565
566
567
568
569
570
571
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
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615

References

- [1] Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao Deng, Wei Ding, Chang Gao, Chunjiang Ge, Wenbin Ge, Zhifang Guo, Qidong Huang, Jie Huang, Fei Huang, Binyuan Hui, Shutong Jiang, Zhaohai Li, Mingsheng Li, Mei Li, Kaixin Li, Zicheng Lin, Junyang Lin, Xuejing Liu, Jiawei Liu, Chenglong Liu, Yang Liu, Dayiheng Liu, Shixuan Liu, Dunjie Lu, Ruilin Luo, Chenxu Lv, Rui Men, Lingchen Meng, Xuancheng Ren, Xingzhang Ren, Sibao Song, Yuchong Sun, Jun Tang, Jianhong Tu, Jianqiang Wan, Peng Wang, Pengfei Wang, Qiuyue Wang, Yuxuan Wang, Tianbao Xie, Yiheng Xu, Haiyang Xu, Jin Xu, Zhibo Yang, Mingkun Yang, Jianxin Yang, An Yang, Bowen Yu, Fei Zhang, Hang Zhang, Xi Zhang, Bo Zheng, Humen Zhong, Jingren Zhou, Fan Zhou, Jing Zhou, Yuanzhi Zhu, and Ke Zhu. Qwen3-vl technical report. *arXiv preprint arXiv:2511.21631*, 2025.
- [2] Mahtab Bigverdi, Zelun Luo, Cheng-Yu Hsieh, Ethan Shen, Dongping Chen, Linda G. Shapiro, and Ranjay Krishna. Perception tokens enhance visual reasoning in multimodal language models. *arXiv preprint arXiv:2412.03548*, 2024.
- [3] Ellis Brown, Jihan Yang, Shusheng Yang, Rob Fergus, and Saining Xie. Benchmark designers should “train on the test set” to expose exploitable non-visual shortcuts. *arXiv preprint arXiv:2511.04655*, 2025.
- [4] Zhongang Cai, Junzhe Zhang, Daxuan Ren, Cunjun Yu, Haiyu Zhao, Shuai Yi, Chai Kiat Yeo, and Chen Change Loy. Messytable: Instance association in multiple camera views. In *European Conference on Computer Vision*, pages 1–16. Springer, 2020.
- [5] Chameleon Team. Chameleon: Mixed-modal early-fusion foundation models. *arXiv preprint arXiv:2405.09818*, 2024.
- [6] Angel Chang, Angela Dai, Thomas Funkhouser, Maciej Halber, Matthias Nießner, Manolis Savva, Shuran Song, Andy Zeng, and Yinda Zhang. Matterport3d: Learning from rgb-d data in indoor environments, 2017.
- [7] Xiaokang Chen, Zhiyu Wu, Xingchao Liu, Zizheng Pan, Wen Liu, Zhenda Xie, Xingkai Yu, and Chong Ruan. Januspro: Unified multimodal understanding and generation with data and model scaling. *arXiv preprint arXiv:2501.17811*, 2025.
- [8] Christopher Clark, Jieyu Zhang, Zixian Ma, Jae Sung Park, Mohammadreza Salehi, Rohun Tripathi, Sangho Lee, Zhongzheng Ren, Chris Dongjoo Kim, YINUO Yang, Vincent Shao, Yue Yang, Weikai Huang, Ziqi Gao, Taira Anderson, Jianrui Zhang, Jitesh Jain, George Stoica, Winson Han, Ali Farhadi, and Ranjay Krishna. Molmo2: Open weights and data for vision-language models with video understanding and grounding. *arXiv preprint arXiv:2601.10611*, 2026.
- [9] Angela Dai, Angel X. Chang, Manolis Savva, Maciej Halber, Thomas Funkhouser, and Matthias Nießner. Scannet: Richly-annotated 3d reconstructions of indoor scenes, 2017.
- [10] Matt Deitke, Eli VanderBilt, Alvaro Herrasti, Luca Weihs, Jordi Salvador, Kiana Ehsani, Winson Han, Eric Kolve, Ali Farhadi, Aniruddha Kembhavi, and Roozbeh Mottaghi. Procthor: Large-scale embodied ai using procedural generation, 2022.
- [11] Matt Deitke, Christopher Clark, Sangho Lee, Rohun Tripathi, Yue Yang, Jae Sung Park, Mohammadreza Salehi, Niklas Muennighoff, Kyle Lo, Luca Soldaini, Jiasen Lu, Taira Anderson, Erin Bransom, Kiana Ehsani, Huong Ngo, Yensung Chen, Ajay Patel, Mark Yatskar, Chris Callison-Burch, Andrew Head, Rose Hendrix, Favyen Bastani, Eli VanderBilt, Nathan Lambert, Yvonne Chou, Arnavi Chheda, Jenna Sparks, Sam Skjonsberg, Michael Schmitz, Aaron Sarnat, Byron Bischoff, Pete Walsh, Chris Newell, Piper Wolters, Tanmay Gupta, Kuo-Hao Zeng, Jon Borchardt, Dirk Groeneveld, Crystal Nam, Sophie Lebrecht, Caitlin Wittlif, Carissa Schoenick, Oscar Michel, Ranjay Krishna, Luca Weihs, Noah A. Smith, Hannaneh Hajishirzi, Ross Girshick, Ali Farhadi, and Aniruddha Kembhavi. Molmo and pixmo: Open weights and open data for state-of-the-art vision-language models. In *CVPR*, 2025.
- [12] Chaorui Deng, Deyao Zhu, Kunchang Li, Chenhui Gou, Feng Li, Zeyu Wang, Shu Zhong, Weihao Yu, Xiaonan Nie, Ziang Song, Guang Shi, and Haoqi Fan. Emerging properties in unified multimodal pretraining. *arXiv preprint arXiv:2505.14683*, 2025.
- [13] Corentin Dumery, Noa Etti, Aoxiang Fan, Ren Li, Jingyi Xu, Hieu Le, and Pascal Fua. Counting stacked objects. In *ICCV*, 2025.
- [14] 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. *arXiv preprint arXiv:2510.27492*, 2025.
- [15] Yushi Hu, Benlin Liu, Jungo Kasai, Yizhong Wang, Mari Ostendorf, Ranjay Krishna, and Noah A Smith. Tifa: Accurate and interpretable text-to-image faithfulness evaluation with question answering, 2023.
- [16] Yushi Hu, Weijia Shi, Xingyu Fu, Dan Roth, Mari Ostendorf, Luke Zettlemoyer, Noah A. Smith, and Ranjay Krishna. Visual sketchpad: Sketching as a visual chain of thought for multimodal language models. *arXiv preprint arXiv:2406.09403*, 2024.
- [17] Amita Kamath, Jack Hessel, and Kai-Wei Chang. What’s “up” with vision-language models? Investigating their struggle with spatial reasoning. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing (EMNLP)*, 2023.
- [18] Rawal Khirodkar, Aayush Bansal, Lingni Ma, Richard Newcombe, Minh Vo, and Kris Kitani. Egohumans: An egocentric 3d multi-human benchmark, 2023.
- [19] Eric Kolve, Roozbeh Mottaghi, Winson Han, Eli VanderBilt, Luca Weihs, Alvaro Herrasti, Matt Deitke, Kiana Ehsani, Daniel Gordon, Yuke Zhu, et al. Ai2-thor: An interactive 3d environment for visual ai. *arXiv preprint arXiv:1712.05474*, 2017.
- [20] Chengzu Li, Wenshan Wu, Huanyu Zhang, Yan Xia, Shaoguang Mao, Li Dong, Ivan Vulić, and Furu Wei. Imagine while reasoning in space: Multimodal visualization-of-thought. *arXiv preprint arXiv:2501.07542*, 2025.
- [21] Linnan Li, Xiaoyu Chen, Peng Chen, et al. ViewSpatial-Bench: Evaluating multi-perspective spatial under-

- 673 standing of vision-language models. *arXiv preprint*
674 *arXiv:2505.21500*, 2025.
- 675 [22] Fangyu Liu, Guy Emerson, and Nigel Collier. Visual spatial
676 reasoning. *arXiv preprint arXiv:2205.00363*, 2022.
- 677 [23] Wufei Ma, Haoyu Chen, Guofeng Zhang, Yu-Cheng Chou,
678 Jieneng Chen, Celso M. de Melo, and Alan Yuille. 3DSR-
679 Bench: A comprehensive 3D spatial reasoning benchmark.
680 In *ICCV*, 2025.
- 681 [24] Manolis Savva*, Abhishek Kadian*, Oleksandr
682 Maksymets*, Yili Zhao, Erik Wijmans, Bhavana Jain,
683 Julian Straub, Jia Liu, Vladlen Koltun, Jitendra Malik, Devi
684 Parikh, and Dhruv Batra. Habitat: A Platform for Embodied
685 AI Research. In *Proceedings of the IEEE/CVF International*
686 *Conference on Computer Vision (ICCV)*, 2019.
- 687 [25] OpenAI. Thinking with images. OpenAI Blog, 2025.
- 688 [26] Xavi Puig, Eric Undersander, Andrew Szot, Mikael Dallaire
689 Cote, Ruslan Partsey, Jimmy Yang, Ruta Desai, Alexan-
690 der William Clegg, Michal Hlavac, Tiffany Min, Theo
691 Gervet, Vladimír Vondruš, Vincent-Pierre Berges, John
692 Turner, Oleksandr Maksymets, Zsolt Kira, Mrinal Kalakr-
693 ishnan, Jitendra Malik, Devendra Singh Chaplot, Unnat Jain,
694 Dhruv Batra, Akshara Rai, and Roozbeh Mottaghi. Habitat
695 3.0: A co-habitat for humans, avatars and robots, 2023.
- 696 [27] Arijit Ray, Ahmed Abdelkader, Chengzhi Mao, Bryan A.
697 Plummer, Kate Saenko, Ranjay Krishna, Leonidas Guibas,
698 and Wen-Sheng Chu. Mull-tokens: Modality-agnostic latent
699 thinking. *arXiv preprint arXiv:2512.10941*, 2025.
- 700 [28] Arijit Ray, Jiafei Duan, Ellis Brown, Reuben Tan, Dina
701 Bashkirova, Rose Hendrix, Kiana Ehsani, Aniruddha Kem-
702 bhavi, Bryan A. Plummer, Ranjay Krishna, Kuo-Hao Zeng,
703 and Kate Saenko. Sat: Dynamic spatial aptitude training for
704 multimodal language models, 2025.
- 705 [29] Andrew Szot, Alex Clegg, Eric Undersander, Erik Wijmans,
706 Yili Zhao, John Turner, Noah Maestre, Mustafa Mukadam,
707 Devendra Chaplot, Oleksandr Maksymets, Aaron Gokaslan,
708 Vladimir Vondrus, Sameer Dharur, Franziska Meier, Woj-
709 ciech Galuba, Angel Chang, Zsolt Kira, Vladlen Koltun, Ji-
710 tendra Malik, Manolis Savva, and Dhruv Batra. Habitat 2.0:
711 Training home assistants to rearrange their habitat. In *Ad-*
712 *vances in Neural Information Processing Systems (NeurIPS)*,
713 2021.
- 714 [30] Michael Tschannen, Alexey Gritsenko, Xiao Wang, Muham-
715 mad Ferjad Naeem, Ibrahim Alabdulmohsin, Nikhil
716 Parthasarathy, Talfan Evans, Lucas Beyer, Ye Xia, Basil
717 Mustafa, et al. Siglip 2: Multilingual vision-language en-
718 coders with improved semantic understanding, localization,
719 and dense features. *arXiv preprint arXiv:2502.14786*, 2025.
- 720 [31] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
721 Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny
722 Zhou. Chain-of-thought prompting elicits reasoning in large
723 language models. *arXiv preprint arXiv:2201.11903*, 2022.
- 724 [32] Chenfei Wu, Jiahao Li, Jingren Zhou, Junyang Lin, Kaiyuan
725 Gao, Kun Yan, Sheng ming Yin, Shuai Bai, Xiao Xu, Yilei
726 Chen, Yuxiang Chen, Zecheng Tang, Zekai Zhang, Zhengyi
727 Wang, An Yang, Bowen Yu, Chen Cheng, Dayiheng Liu, De-
728 qing Li, Hang Zhang, Hao Meng, Hu Wei, Jingyuan Ni, Kai
729 Chen, Kuan Cao, Liang Peng, Lin Qu, Minggang Wu, Peng
Wang, Shuting Yu, Tingkun Wen, Wensen Feng, Xiaoxiao
Xu, Yi Wang, Yichang Zhang, Yongqiang Zhu, Yujia Wu,
Yuxuan Cai, and Zenan Liu. Qwen-image technical report,
2025.
- [33] Jinheng Xie, Zhenheng Yang, and Mike Zheng Shou. Show-
o2: Improved native unified multimodal models. *arXiv*
preprint arXiv:2506.15564, 2025.
- [34] Jihan Yang, Shusheng Yang, Anjali W. Gupta, Rilyn Han,
Li Fei-Fei, and Saining Xie. Thinking in space: How mul-
timodal large language models see, remember, and recall
spaces. *arXiv preprint arXiv:2412.14171*, 2025.
- [35] Kaiyu Yang, Olga Russakovsky, and Jia Deng. SpatialSense:
An adversarially crowdsourced benchmark for spatial rela-
tion recognition. In *ICCV*, 2019.
- [36] Rui Yang, Ziyu Zhu, Yanwei Li, Jingjia Huang, Shen Yan,
Siyuan Zhou, Zhe Liu, Xiangtai Li, Shuangye Li, Wen-
qian Wang, et al. Visual spatial tuning. *arXiv preprint*
arXiv:2511.05491, 2025.
- [37] Sihan Yang, Runsen Xu, Yiman Xie, et al. MMSI-Bench:
A benchmark for multi-image spatial intelligence. *arXiv*
preprint arXiv:2505.23764, 2025.
- [38] Zeyuan Yang, Xueyang Yu, Delin Chen, Maohao Shen, and
Chuang Gan. Machine mental imagery: Empower multi-
modal reasoning with latent visual tokens. *arXiv preprint*
arXiv:2506.17218, 2025.
- [39] Chun-Hsiao Yeh, Chenyu Wang, Shengbang Tong, Ta-Ying
Cheng, Ruoyu Wang, Tianzhe Chu, Yuexiang Zhai, Yubei
Chen, Shenghua Gao, and Yi Ma. Seeing from another
perspective: Evaluating multi-view understanding in mlms.
arXiv preprint arXiv:2504.15280, 2025.
- [40] Chandan Yeshwanth, Yueh-Cheng Liu, Matthias Nießner,
and Angela Dai. Scannet++: A high-fidelity dataset of 3d in-
door scenes. In *Proceedings of the IEEE/CVF International*
Conference on Computer Vision, pages 12–22, 2023.
- [41] Baiqiao Yin, Qineng Wang, Pingyue Zhang, et al. Spa-
tial mental modeling from limited views. *arXiv preprint*
arXiv:2506.21458, 2025.
- [42] Pingyue Zhang, Zihan Huang, Yue Wang, Jieyu Zhang,
Letian Xue, Zihan Wang, Qineng Wang, Keshigeyan Chan-
drasegaran, Ruohan Zhang, Yejin Choi, Ranjay Krishna, Ji-
ajun Wu, Li Fei-Fei, and Manling Li. Theory of space: Can
foundation models construct spatial beliefs through active
exploration? In *International Conference on Learning Rep-*
resentations (ICLR), 2026.