

# Imaginative Perception Tokens Enhance Spatial Reasoning in Multimodal Language Models

## Supplementary Material

### 780 7. Training Details and Hyperparameters

#### 781 7.1. Training Setup

782 We fine-tune BAGEL-7B-MoT [12] using PyTorch FSDP  
783 (Fully Sharded Data Parallel) with bf16 mixed precision  
784 on 8 NVIDIA A100 80 GB GPUs. Table 7 summarizes the  
785 key hyperparameters.

786 **System prompts.** All training modes share one of two sys-  
787 tem prompts prepended to every input:

- 788 • **Thinking prompt** (used for IPT, Text CoT, and label-  
789 only):

790 *Let's think step by step to answer the ques-*  
791 *tion. For text-based thinking, enclose the process*  
792 *within <think> </think>. For visual think-*  
793 *ing, enclose the content within <image\_start>*  
794 *</image\_end>. Finally conclude with the final*  
795 *answer wrapped in <answer></answer> tags.*

- 796 • **Answer-only prompt** (used for the answer-only portion  
797 of mixed training):

798 *Answer the question directly. Wrap your answer*  
799 *in <answer></answer> tags. Do not think or*  
800 *generate any images.*

801 **Training modes.** Table 8 summarizes the five training con-  
802 figurations evaluated in this work. Each mode differs in the  
803 output format the model is trained to produce, and corre-  
804 spondingly in which loss terms are active. In IPT mode,  
805 both  $\mathcal{L}_{fm}$  and  $\mathcal{L}_{lm}$  are active; in Text CoT and label-only,  
806 only  $\mathcal{L}_{lm}$  is used.

807 **Mixed training.** Mixed training combines 50% IPT ex-  
808 amples (with the thinking prompt and visual generation  
809 targets) and 50% answer-only examples (with the answer-  
810 only prompt and direct answers). The two data subsets  
811 are mixed at the dataloader level: both dataset names and  
812 sample counts are specified in the training configuration,  
813 and the dataloader interleaves batches from both sources.  
814 The model learns to switch between generating imaginative  
815 perception tokens and producing direct answers based on  
816 which system prompt is provided, enabling a single check-  
817 point to operate in either mode at inference time.

818 **Text CoT generation.** Text chain-of-thought training tar-  
819 gets are generated by GPT-5.1 using ground-truth scene  
820 metadata from the simulator. For each training example,  
821 GPT-5.1 receives the input image, the ground-truth answer,  
822 and a task-specific instruction, and produces a step-by-step  
823 textual reasoning trace (100–300 words). Below we de-  
824 scribe the task-specific prompts.

**Path tracing CoT prompt.** The system prompt instructs  
the model to act as a spatial reasoning AI that solves indoor  
navigation questions by analyzing top-down views step by  
step. The task instruction is:

*You are navigating along a numbered path through an*  
*indoor scene. The top-down view shows the path with*  
*waypoints and midpoints. You need to determine what*  
*object is visible from a specific side (left or right) at a*  
*midpoint.*

*Reason step by step: (1) Identify the path direction*  
*(which waypoint to which). (2) Determine your ori-*  
*entation at the midpoint. (3) Figure out what “left”*  
*or “right” means given that orientation. (4) Analyze*  
*the top-down layout to identify objects on that side.*  
*(5) Compare against the answer choices and eliminate*  
*wrong ones.*

Figure 3 shows three Text CoT training examples for path  
tracing in the EgoDir setting, illustrating how the generated  
reasoning traces analyze the top-down layout, determine the  
agent’s orientation at the midpoint, and systematically elim-  
inate incorrect answer choices.

**Perspective Taking CoT prompt.** For Perspective Tak-  
ing, the model receives the input image together with priv-  
ileged hidden information, the ground-truth answer and a  
second image showing the target view after the camera mo-  
tion, to ensure correctness. The generated chain-of-thought  
must be written as if only the original image and question  
were available, without referencing any hidden information.  
The prompt is:

*You generate student-facing chain-of-thought explana-*  
*tions for visual navigation/spatial reasoning from a*  
*single image. You may receive hidden privileged in-*  
*formation (metadata, ground-truth answer, and a sec-*  
*ond image) to ensure correctness, but your explanation*  
*must be written as if you only saw the original image*  
*and question. Do not mention, quote, paraphrase, or*  
*allude to any hidden information.*

*Student-visible input (ONLY what the explanation may*  
*reference): (1) Image A (original scene with a red “X”*  
*on the floor). (2) Question: {Q}*

*Hidden privileged information (FOR CORRECTNESS*  
*ONLY — MUST NOT APPEAR IN THE EXPLANA-*  
*TION): (1) Correct answer: {A} (2) Image B (ground-*  
*truth final view after moving and turning)*

Table 7. Training hyperparameters.

Parameter	Value
Optimizer	AdamW ( $\beta_1=0.9, \beta_2=0.95, \epsilon=10^{-15}$ )
Learning rate	$1 \times 10^{-5}$ (constant after warmup)
Warmup steps	2,000
Gradient clipping	max norm = 1.0
EMA decay	0.9999
Max tokens per batch	32,768
Max tokens per sample	24,576
Input image resolution	1024 $\times$ 1024 (PET) / 512 $\times$ 512 (PT, MVC)
IPT latent resolution	64 $\times$ 64 (Latent-64)
Flow-matching loss weight ( $\lambda_{fm}$ )	1.0
Language modeling loss weight ( $\lambda_{lm}$ )	1.0
Frozen modules	VAE encoder & decoder
Fine-tuned modules	LLM (all layers), ViT encoder, connector

Table 8. Training modes and their output formats. [IMG] denotes the generated intermediate image tokens.

Mode	Prompt	Training target
Label-only	Think	<answer>A</answer>
Text CoT	Think	<think>reasoning</think> <answer>A</answer>
IPT	Think	<think>task prompt</think> <image_start>[IMG]<image_end> <answer>A</answer>
Mixed	50% Think / 50% Ans-only	50% IPT + 50% answer-only

869 *Task: Write a chain-of-thought that a student could*  
 870 *produce from Image A and the question alone (not more*  
 871 *than 10 sentences).*

872 *Hard constraints: Do NOT mention Image B, meta-*  
 873 *data, coordinates, angles, yaw/pitch/roll, agent states,*  
 874 *movement info, IDs, or any numbers from the hidden*  
 875 *block. The explanation must be plausible using only*  
 876 *Image A: talk only about visible objects, their relative*  
 877 *positions, and the effect of turning left/right. Keep it*  
 878 *short (3–6 steps). No extra commentary. End with*  
 879 *the correct final answer: My answer is [A]/[B]. Do*  
 880 *not answer based on your own perspective; stick to the*  
 881 *ground-truth answer provided.*

882 *Output format:*  
 883 *Step 1: Identify the object in the question and describe*  
 884 *its relative location in the original view.*  
 885 *Step 2: Describe the walk-to-X action (translation*  
 886 *without changing facing direction).*  
 887 *Step 3: Describe the turn action (update egocentric*  
 888 *left/right).*  
 889 *Step 4: Conclude whether the object is on the left or*  
 890 *right relative to the agent.*

891 An example of a generated thought:

892 *Step 1: The oven is on the right wall of the kitchen, a*  
 893 *bit ahead of the current viewpoint, and the red X is on*  
 894 *the floor between the viewer and that oven wall. Step 2:*  
 895 *Walking straight to the X without turning keeps me fac-*  
 896 *ing the same direction, so I move forward toward the*  
 897 *area in front of the oven and end up nearer that wall.*  
 898 *Step 3: From the X, turning right 90 degrees means I*  
 899 *would now face the wall with the oven and cabinets.*  
 900 *Step 4: Since I walked closer to that wall before turn-*  
 901 *ing, I end up standing nearer to the oven than I was at*  
 902 *the start. My answer is [A].*

903 **Multiview Counting CoT prompt.** For Multiview  
 904 Counting, the model receives per-frame metadata listing  
 905 which object instances (by ID) are visible in each frame,  
 906 and is asked to produce a frame-by-frame reasoning trace  
 907 without directly referencing the hidden instance IDs. The  
 908 prompt is:

909 *You are given an object-counting question based on*  
 910 *multiple images from the same scene. You will be pro-*  
 911 *vided with: the target object, the correct total count,*  
 912 *and frame-level information listing which object in-*  
 913 *stances are visible in each frame (this information is*  
 914 *hidden and should not be referenced directly).*

## Text CoT Training Examples (Path Tracing, EgoDir)

**Question:** You are navigating from waypoint 1 to 2, passing through M1, and viewing the space from a top-down perspective as shown in <image\_1>. At point 1, I'm facing forward and see: <image\_2>. At point 2, facing forward: <image\_3>. Which object can you see on your **left** side at waypoint M1?

A: the countertop ✓  B: the house plant  C: the dining table  D: the dresser

**Generated Text CoT:**

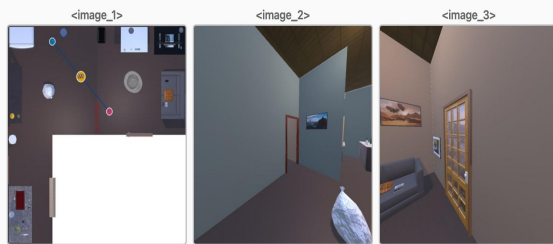
From the top-down map, the path goes from waypoint 1 (to the right) toward waypoint 2 (down-left), with M1 located slightly left of 1 along that diagonal segment. So at M1, "forward" is oriented down-left (southwest on the map). Facing southwest, the "left side" corresponds to the southeast side of the path (down-right relative to the direction of travel).

Looking at the top-down layout around M1: a long white counter run is positioned just to the right (east/southeast) of M1, aligned roughly horizontally near waypoint 1. The house plant is up-left (northwest) of M1, which would be on the right side when facing southwest. The dining table is far up-right, not adjacent to M1's left side. The dresser/cabinet block is also up-left of M1, again on the right side for this heading.

Therefore, the object visible on the left side at waypoint M1 is **the countertop (A)**.

**Question:** You are navigating from waypoint 1 to 2, passing through M1, and viewing the path from a top-down perspective as shown in <image\_1>. At point 1, I'm facing forward and see: <image\_2>. At point 2, facing forward: <image\_3>. Which object can you see on your **left** side at waypoint M1?

A: a book ✓  B: a shelving unit  C: a garbage bag  D: a stool

**Generated Text CoT:**

From the top-down map, the path goes diagonally from waypoint 1 (upper-left) to waypoint 2 (lower-right), with M1 in between. So at M1, "forward" means continuing along that same diagonal toward waypoint 2 (southeast). Facing southeast, your left side points toward the northeast (up-right on the map).

Using the egocentric views to anchor objects: at waypoint 2 we see a couch area with a small rectangular object on the couch that looks like a book. On the top-down view, that couch is on the right side of the room, above and to the right of waypoint 2, i.e., on the northeast side relative to M1.

Now check the choices against the left (northeast) side at M1: the garbage bag is positioned left-of-center (more west), not on the northeast side; the shelving unit appears in the far top-right but is farther and not the prominent item; the stool is near the center-right but not as clearly on the left-of-forward line at M1. The book on/near the couch lies on M1's left side when moving toward 2.

Therefore, the object visible on the left at M1 is **the book (A)**.

**Question:** You are following the path 1→2, looking at a top-down view at waypoint M1 as shown in <image\_1>. This is what you see looking forward from point 1: <image\_2>. This is what you see looking forward from point 2: <image\_3>. Which object can you see on your **left** side at waypoint M1?

A: the floor lamp  B: the armchair ✓  C: the sink  D: the bed

**Generated Text CoT:**

From the top-down map, the path goes from waypoint 1 on the right side of the blue-floored bedroom area toward waypoint 2 in the brown-floored bathroom area on the left. So at the midpoint M1, our forward travel direction is generally leftward (west) toward the bathroom doorway.

To interpret "left side" at M1, take that forward direction (toward 2); left of a westward heading is south (down on the map). Looking at the layout around M1, the large floral armchair sits in the lower (southern) part of the blue bedroom area, just below the path near waypoint 1/M1. The bed is far up at the top of the bedroom, so it would be on the right side (north) relative to the westward heading, not the left. The sink is inside the bathroom area to the left of the doorway, ahead/left but separated by the doorway and not on the immediate left side at M1. A floor lamp is not positioned on the south side near M1.

Therefore, the object visible on the left at M1 is **the armchair (B)**.

Figure 3. **Text CoT training examples for path tracing (EgoDir setting).** Each example shows the input images (top-down map and egocentric views at endpoints), the question with answer choices, and the GPT-5.1-generated reasoning trace. The reasoning follows the structured prompt: identifying path direction, determining orientation at  $M_1$ , interpreting left/right relative to that orientation, and eliminating distractors.

915 *IMPORTANT: An empty list for a frame means that no*  
916 *objects of the target type are visible in that frame.*

917 *Write a brief, frame-by-frame explanation describing*  
918 *what is visible in each frame. Do not mention object*

*IDs or refer to them explicitly. When explaining each*  
*frame, do not count objects that were already visible in*  
*previous frames. After the frame-by-frame explanation,*  
*conclude with: "The total number is X."*

919  
920  
921  
922

923	<i>Object: {O}</i>	969
924	<i>Correct total count: {answer}</i>	970
925	<i>Frames: {frames_text}</i>	971
926	An example of a generated thought:	972
927	<i>Frame 1: No bowls are visible. Frame 2: A bowl ap-</i>	973
928	<i>pears and is visible for the first time. Frame 3: No new</i>	974
929	<i>bowls appear; the same bowl from before may still be</i>	
930	<i>present. The total number is 1.</i>	
931	<b>7.2. Evaluation Setup</b>	
932	<b>Inference parameters.</b> Table 9 lists the inference hyperpa-	
933	rameters used across all evaluations.	
934	<b>Evaluation modes.</b> At inference, each model variant	
935	is evaluated in the mode that matches its training con-	
936	figuration. IPT models are evaluated in two settings:	
937	(1) <i>imagination mode</i> , where the model generates an in-	
938	termediate image before answering, and (2) <i>answer-only</i>	
939	<i>mode</i> , where the model produces only a text answer with-	
940	out generating any image. For models trained with vi-	
941	ual generation (IPT, Mixed), the VAE weights are al-	
942	ways loaded ( <code>visual_gen=True</code> ), and input images	
943	are encoded through both the ViT and VAE pathways	
944	( <code>vae_input=True</code> ) to match the training-time encoding.	
945	Without setting <code>vae_input=True</code> , a train-eval mismatch	
946	would occur: during training, input images pass through	
947	both VAE and ViT, but the default evaluation behavior sends	
948	inputs through ViT only. <sup>1</sup>	
949	<b>Answer extraction.</b> We extract the predicted answer let-	
950	ter from model outputs using a cascading rule-based proce-	
951	cedure: (1) parse <code>&lt;answer&gt;X&lt;/answer&gt;</code> tags; (2) extract	
952	from <code>\boxed{X}</code> format; (3) match patterns such as “the	
953	answer is X”; (4) detect bold letter formatting ( <code>**X**</code> );	
954	(5) fall back to the last single letter in the response. All	
955	benchmarks use the same unified scoring function that com-	
956	pares the extracted letter against the ground-truth answer.	
957	<b>8. Data Curation Details</b>	
958	<b>8.1. Path Tracing</b>	
959	We generate path tracing data from two sources: AI2-	
960	THOR (synthetic) and Matterport3D (real-world).	
961	<b>8.1.1. AI2-THOR</b>	
962	<b>Scene selection.</b> We use 120 standard iTHOR [20] scenes	
963	spanning four room types (kitchens, living rooms, bed-	
964	rooms, bathrooms), with 30 scenes per type, split into train	
965	(20), val (5), and test (5) per type. The training set ad-	
966	ditionally incorporates procedurally generated houses from	
967	ProcTHOR-10k [10], which include a fifth room type (hall-	
968	ways, offices, dining rooms).	
	<b>Path sampling.</b> For each scene, we sample feasible two-	969
	waypoint paths on the navigation mesh, balanced across	970
	room types and three distance bins: short (1–2 m), medium	971
	(2–4 m), and long ( $\geq 4$ m). Grid-based path sampling uses	972
	a spacing of 0.5 m with a minimum waypoint separation of	973
	1.0 m.	974
	<b>Camera configuration.</b> All views are rendered at $1024 \times$	975
	$1024$ resolution. To increase viewpoint diversity, we ran-	976
	domize the camera height (sampled from seven values be-	977
	tween 1.4 and 1.8 m), field of view (75–120), and pitch (–5	978
	to 5).	979
	<b>Rendering.</b> At each sample we render top-down views,	980
	egocentric forward views at both endpoints, and a sweep	981
	of candidate sideviews at the midpoint $M_1$ . The sideview	982
	sweep covers 7 yaw angles $\times$ 7 horizontal offsets $\times$ 3 pitch	983
	values, yielding 147 candidate views per midpoint. We se-	984
	lect the sideview that best exposes the queried object using	985
	simulator segmentation masks, requiring a minimum object	986
	coverage of 0.15% of the image area and a maximum view	987
	angle of 90 relative to the path direction.	988
	<b>Question generation.</b> Questions are generated from tem-	989
	plates (“Which object can you see on your {side} at way-	990
	point $M_1$ ?”) with four choices: the correct answer drawn	991
	from verified visible objects and three distractors drawn	992
	from the opposite side or a global object pool. Each base	993
	MCQ is expanded into eight input variants by combining	994
	different image types (top-down path, top-down with arrow,	995
	top-down with midpoint marker, dollhouse view) and ego-	996
	centric cue availability (with/without endpoint views).	997
	<b>TIFA filtering.</b> We apply TIFA-style filtering [16] to en-	998
	sure question quality. Each candidate is decomposed into	999
	binary visibility queries and verified by GPT-4.1 with three-	1000
	round majority voting, using early exit after round 2 when	1001
	unanimous. Samples are dropped if (1) the correct answer	1002
	is not visible in the sideview, (2) a distractor is also visible	1003
	in the sideview, or (3) the model answers incorrectly even	1004
	when provided the sideview. We further remove samples	1005
	where GPT-4.1 answers correctly from the top-down view	1006
	and egocentric endpoint views alone, ensuring the bench-	1007
	mark requires genuine spatial imagination.	1008
	<b>Debiasing.</b> Answer choices are reshuffled with per-	1009
	sample deterministic seeds to remove positional bias. We	1010
	additionally verify that per-object and per-room answer dis-	1011
	tributions remain approximately uniform.	1012
	<b>Statistics.</b> The synthetic training set contains 11,204 ex-	1013
	amples.	1014

<sup>1</sup>For Path-Tracing, we found that setting `vae_input=False` actually improves generalization to real environments.

Table 9. Inference hyperparameters for evaluation.

Parameter	Value
Text temperature	0.3
Sampling	do_sample = True
Max thinking tokens	4,096
Diffusion timesteps	50
Timestep shift	3.0
Text CFG scale	4.0
Image CFG scale	2.0
CFG interval	[0.0, 1.0]
Generated image resolution	1024 × 1024
Max generation rounds	1
GPU allocation	2 × A100 80 GB (model parallelism)

### 1015 8.1.2. Real-World Data (Matterport3D)

1016 To evaluate cross-domain transfer, we construct a real-  
1017 world test set from Matterport3D [6] top-down views.

1018 **Image collection.** We collect top-down screenshots from  
1019 Matterport 3D indoor tours, capturing per-floor views with  
1020 UI elements removed and dark borders cropped.

1021 **Auto-annotation.** We annotate walking paths on each im-  
1022 age using a two-pass GPT pipeline. In the first pass, GPT  
1023 proposes  $N$  candidate walking paths, each defined by three  
1024 waypoints (start “1”, midpoint “M1”, end “2”) placed on  
1025 open floor areas. In the second pass, the proposed way-  
1026 points are drawn on the image and GPT is asked to (a) ver-  
1027 ify that  $M_1$  lies on walkable floor, adjusting its position if  
1028 necessary, and (b) identify 2–5 visible furniture items or ob-  
1029 jects on each side (left/right) of the path at  $M_1$ .

1030 **Post-processing.** Several geometric filters are applied to  
1031 ensure path quality. Waypoints are clamped to image  
1032 bounds, and paths shorter than 30% of the shorter image  
1033 dimension are rejected.  $M_1$  is snapped onto the line seg-  
1034 ment between waypoints 1 and 2 and constrained to the  
1035 [0.2, 0.8] interval to avoid proximity to endpoints. Paths  
1036 that traverse dark or background regions (more than 20%  
1037 dark pixels along the path) are discarded.

1038 **TIFA filtering.** Because no side-view images exist for  
1039 real environments, TIFA verification operates on the top-  
1040 down image only. We verify three properties: (1) all three  
1041 waypoints lie on walkable floor (not on furniture, walls, or  
1042 background); (2) each annotated object is visible in the im-  
1043 age; and (3) each object is on the correct side of the path.  
1044 Paths with fewer than 2 verified objects on either side are  
1045 dropped. Majority voting across up to 3 rounds is used,  
1046 with early exit when the first two rounds agree.

**Human review.** After automated filtering, all surviving  
1047 annotations undergo human review, where annotators can  
1048 approve, delete, or edit individual paths and their associated  
1049 object lists. 1050

**Question generation.** Each verified path yields eight  
1051 question variants (forward/reverse × left/right × arrow/no-  
1052 arrow). When the walking direction is reversed ( $2 \rightarrow 1$ ),  
1053 left and right swap because the agent faces the opposite di-  
1054 rection. Distractors are drawn from opposite-side objects  
1055 first, then supplemented from a global pool of 50 com-  
1056 mon indoor objects. Answer choices are shuffled with per-  
1057 sample deterministic seeds to eliminate positional bias. Be-  
1058 cause the real environments lack egocentric viewpoints, we  
1059 evaluate on the Path and PathArr settings only. The real-  
1060 world benchmark contains 332 human-verified questions. 1061

## 8.2. Perspective Taking 1062

We generate perspective taking data from three sources:  
1063 AI2-THOR (synthetic), Habitat (photorealistic scans), and  
1064 Visual Spatial Tuning (real-world images). All sources  
1065 share the same task structure—given a first-person view  
1066 with a marked target position, answer a spatial question  
1067 about the scene from the new viewpoint—but differ in vi-  
1068 sual domain and 3D engine. 1069

### 8.2.1. AI2-THOR 1070

**Scene selection.** We use procedurally generated indoor  
1071 scenes from ProcTHOR [10], which provides diverse house  
1072 layouts with varied room configurations. For each scene, we  
1073 sample multiple camera positions from the navigable area  
1074 and generate perspective-taking examples at each position. 1075

**Target position placement.** The target position (marked  
1076 with a red “X” on the input image) is determined by ray-  
1077 casting from the camera into the scene. To ensure physically  
1078 plausible movement targets, we filter ray hits to *ground-only*  
1079 surfaces: floors, carpets, rugs, and tiles. Hits on furniture,  
1080

1081	tabletops, or other elevated surfaces are rejected. The hit	
1082	point must lie within $\pm 0.2$ m of the ground-level height. Up	
1083	to 40 raycasting attempts are made per sample; if no valid	
1084	ground hit is found, the sample is skipped.	
1085	<b>Viewpoint transformation.</b> The agent is teleported to the	
1086	target position and rotated 90 in a randomly chosen direc-	
1087	tion (left or right with equal probability), simulating a re-	
1088	alistic movement-and-turn action. A ground-truth novel-	
1089	viewpoint image is rendered at this new pose to serve as	
1090	the imaginative perception target.	
1091	<b>Object filtering.</b> To ensure well-defined questions, target	
1092	objects must satisfy several criteria:	
1093	• Visible in <i>both</i> the original and new viewpoint (dual-view	
1094	visibility).	
1095	• Occupy at least 0.4% of the image area.	
1096	• Lie within 5 m of the camera.	
1097	• Fall at least 150 px from the image edge (to avoid partially	
1098	visible objects).	
1099	• For relative position questions: the object must be unam-	
1100	biguously on one side of the image center, with a 150 px	
1101	margin from the center line.	
1102	Additionally, we enforce a <i>left-right eligibility</i> constraint:	
1103	an object is only used for relative position questions if it	
1104	appears on at most one side of the image (not straddling the	
1105	center), ensuring that the left/right answer is unambiguous.	
1106	<b>Question generation.</b> We generate two types of questions	
1107	with 10 template variants each, varying in person perspec-	
1108	tive (first, second, third person) and formality level:	
1109	• <b>Distance change:</b> “After moving to ‘X’ and turning	
1110	{direction} 90, will the {object} get closer or further?”	
1111	Requires a minimum distance change of $\pm 0.5$ m between	
1112	the old and new camera positions.	
1113	• <b>Relative position:</b> “After moving to ‘X’ and turning	
1114	{direction} 90, will the {object} be on your left or right?”	
1115	Left/right is determined by the object’s 2D position in the	
1116	new-viewpoint image.	
1117	<b>Sub-categories.</b> As described in Sec. 3.1, the six balanced	
1118	sub-categories arise from the combination of question type	
1119	and answer: two for distance change ( <i>closer</i> , <i>further</i> ) and	
1120	four for relative position ( <i>left</i> $\rightarrow$ <i>left</i> , <i>left</i> $\rightarrow$ <i>right</i> , <i>right</i> $\rightarrow$ <i>left</i> ,	
1121	<i>right</i> $\rightarrow$ <i>right</i> ), where the notation indicates the object’s lat-	
1122	eral position before and after the viewpoint transformation.	
1123	The training set is balanced across all six sub-categories.	
1124	<b>Image annotation.</b> Each input image is annotated in two	
1125	versions: (1) with only a red “X” marking the target po-	
1126	sition, and (2) with the “X” plus a blue directional arrow	
1127	indicating the agent’s facing direction after rotation. The	
1128	arrow version provides an additional spatial cue at evalua-	
1129	tion time.	
1130	<b>Statistics.</b> The AI2-THOR training set contains 20,531	
1131	examples across 98 scenes, with an average of $\sim 210$ ques-	
1132	tions per scene.	
	<b>8.2.2. Habitat</b>	1133
	<b>Scene source.</b> We use photorealistic 3D scans from HM3D	1134
	(Habitat-Matterport 3D) [6] with semantic annotations.	1135
	Only single-floor scenes are selected (floor-level Y-variance	1136
	$< 2.0$ m) to avoid cross-level ambiguities.	1137
	<b>Camera configuration.</b> Images are rendered at $1024 \times$	1138
	$1024$ resolution with a horizontal field of view of 90. The	1139
	sensor height is set to 1.25 m above the navigable floor sur-	1140
	face, matching a standing human eye level.	1141
	<b>Target position and viewpoint.</b> Camera A (original view-	1142
	point) is placed at a random navigable point with a random	1143
	yaw. The target position (Camera B) is determined by se-	1144
	lecting a visible object as an anchor: Camera B is placed	1145
	at the object’s XZ coordinates, offset slightly along the fac-	1146
	ing direction to avoid clipping into geometry, and snapped	1147
	to the nearest navigable point on the mesh (within a 1.0 m	1148
	snap radius). The ground-truth novel-viewpoint image is	1149
	rendered at Camera B’s position and orientation.	1150
	<b>Object filtering.</b> We apply a strict whitelist of $\sim 70$ main-	1151
	stream furniture categories (seating, tables, beds, storage,	1152
	appliances, bathroom fixtures) and exclude structural ele-	1153
	ments (walls, floors, doors). Objects must occupy at least	1154
	0.8% of the image area in the original frame and at least	1155
	0.5% in the imagined frame. To avoid ambiguous refer-	1156
	ences, only objects whose category is <i>unique</i> in the frame	1157
	are used ( <i>e.g.</i> , if two chairs are visible, neither is selected as	1158
	a question target). An edge margin of 200 px is applied.	1159
	<b>Left/right determination.</b> Object laterality is determined	1160
	by the mean x-coordinate of the object’s semantic segmen-	1161
	tation mask in the rendered image, with a 180 px margin	1162
	from the image center (512 px). Objects falling in the cen-	1163
	ter zone ( $332 < x < 692$ ) are excluded as ambiguous.	1164
	<b>Statistics.</b> The Habitat training set contains 19,998 exam-	1165
	ples balanced across the six sub-categories.	1166
	<b>8.2.3. Real-World Data (VST)</b>	1167
	To bridge the synthetic-to-real domain gap, the <i>mixed</i> train-	1168
	ing variant incorporates 15,000 real-world examples drawn	1169
	from the camera motion subset of the Visual Spatial Tuning	1170
	(VST) dataset [37]. Each example contains a pair of multi-	1171
	view images captured from different viewpoints in real in-	1172
	door scenes, along with a question about the camera motion	1173
	between them and a corresponding answer.	1174
	<b>Filtering uncertain answers.</b> We first filter out examples	1175
	whose answers are uncertain or underspecified using GPT-	1176
	5.1, prompting it as a binary classifier. Any example whose	1177
	answer contains phrases such as “cannot be determined,”	1178
	“unknown,” or “insufficient information” is removed. The	1179
	filtering prompt is:	1180

1181	<i>You are a binary classifier.</i>			
1182	<i>Proposed Answer:</i>			
1183	<code>&lt;start_answer&gt;{A}&lt;end_answer&gt;</code>			
1184	<i>If the answer contains ANY of the following</i>			
1185	<i>phrases or meanings, output “NO”:</i>			
1186	• <i>cannot be determined</i>			1230
1187	• <i>insufficient information</i>			1231
1188	• <i>not enough information</i>			1232
1189	• <i>unknown</i>			
1190	• <i>unclear</i>			
1191	• <i>cannot tell</i>			
1192	• <i>impossible to determine</i>			
1193	• <i>indeterminate</i>			
1194	<i>Only output “YES” if the answer clearly states a</i>			
1195	<i>specific, determined choice (e.g., a direction, lo-</i>			
1196	<i>cation, label, or concrete option).</i>			
1197	<i>Output exactly one token: YES or NO.</i>			
1198	<b>Rewriting into generation prompts.</b> After filtering, we			
1199	use GPT-5.1 to rewrite each question, answer pair into a			
1200	generation prompt describing the camera motion <i>from the</i>			
1201	<i>first image to the second.</i> This requires careful handling			
1202	of reference frame direction: if the original question asks			
1203	where the first camera is relative to the second image, the			
1204	motion direction must be inverted before constructing the			
1205	prompt. The rewriting prompt is:			
1206	<i>You are given a question about camera motion be-</i>			
1207	<i>tween two images and its correct answer.</i>			
1208	<b>IMPORTANT INVERSION RULE:</b>			
1209	• <i>If the question asks “where is the FIRST cam-</i>			
1210	<i>era relative to the SECOND image” (or uses</i>			
1211	<i>the second image as reference), then the answer</i>			
1212	<i>describes motion FROM second TO first.</i>			
1213	• <i>In this case, you MUST invert the direction to</i>			
1214	<i>get motion FROM first TO second.</i>			
1215	• <i>If the question asks “where is the SECOND</i>			
1216	<i>camera relative to the FIRST image” (or uses</i>			
1217	<i>the first image as reference), NO inversion is</i>			
1218	<i>needed.</i>			
1219	<b>STEPS:</b>			
1220	1. <i>Determine which image is the reference point</i>			
1221	<i>in the question.</i>			
1222	2. <i>If the reference is the second image, invert the</i>			
1223	<i>direction in the answer.</i>			
1224	3. <i>Using the final motion FROM first TO second,</i>			
1225	<i>create a generation prompt.</i>			
1226	<b>INVERSION EXAMPLES:</b>			
1227	• <i>“right” → “left”</i>			
1228	• <i>“left” → “right”</i>			
1229	• <i>“front” → “back”</i>			
		• <i>“back” → “front”</i>		1230
		• <i>“front left” → “back right”</i>		1231
		• <i>“back right” → “front left”</i>		1232
		<i>Your output should start with “generate” and</i>		1233
		<i>describe viewing the scene from the new camera</i>		1234
		<i>position after applying the motion from the first</i>		1235
		<i>image.</i>		1236
		<i>Question:</i> <code>&lt;start_question&gt;{Q}&lt;end_question&gt;</code>		1237
		<i>Answer:</i> <code>&lt;start_answer&gt;{A}&lt;end_answer&gt;</code>		1238
		<i>First, identify the reference image. Then apply</i>		1239
		<i>inversion if needed. Then generate your output.</i>		1240
		The resulting generation prompts condition the model on		1241
		the first view and the inferred camera motion, with the sec-		1242
		ond view serving as the imaginative perception target. Be-		1243
		cause no programmatic 3D annotation is available for these		1244
		real scenes, this data serves as a domain bridge rather than		1245
		a source of the full six-sub-category question format.		1246
		<b>Mixed training composition.</b> The mixed PET training		1247
		variant combines AI2-THOR (20,531), Habitat (19,998),		1248
		and VST (15,000) examples, totaling 55,529 samples.		1249
		<b>8.3. Multiview Counting</b>		1250
		We construct multiview counting data from both synthetic		1251
		and real-image sources. Our main training set is generated		1252
		from ProcTHOR/AI2-THOR environments, which provide		1253
		full 3D supervision for both egocentric observations and		1254
		top-down bird’s-eye-view (BEV) targets. To complement		1255
		this synthetic source, we additionally curate two real-image		1256
		multiview counting sets from MessyTable and ScanNet++,		1257
		which expose the model to real visual appearance and par-		1258
		tial observability under natural image statistics.		1259
		<b>8.3.1. ProcTHOR / AI2-THOR</b>		1260
		We generate the main multiview counting training set from		1261
		AI2-THOR environments, using two trajectory types that		1262
		capture complementary modes of partial observability.		1263
		<b>Trajectory types.</b>		1264
		• <b>Rotation:</b> The agent remains at a fixed position and ro-		1265
		tates in 90° increments through four cardinal directions		1266
		(0°, 90°, 180°, 270°), producing four frames that together		1267
		cover a 360° panorama of the surrounding area.		1268
		• <b>Multi-camera:</b> The agent traverses a square path, cap-		1269
		turing one frame at each of four corners. This setup sim-		1270
		ulates a multi-camera rig where viewpoints are spatially		1271
		distributed around the scene.		1272
		Both trajectory types produce exactly four input frames per		1273
		sample.		1274
		<b>Bird’s-eye view (BEV) generation.</b> The ground-truth in-		1275
		termediate image is a top-down BEV map rendered from an		1276

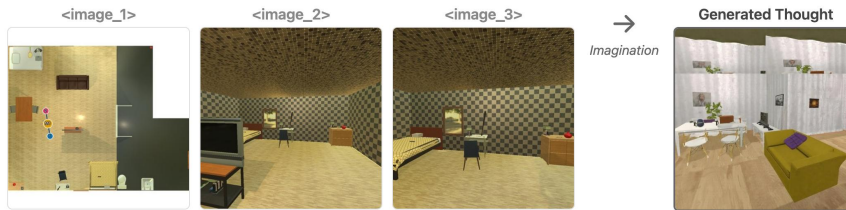
1277	overhead camera in the 3D scene. To ensure the BEV only	10% from the singleton bucket. Input images are selected	1323
1278	covers the <i>explored</i> area (the region visible from the input	from the eight surrounding non-top cameras, with priority	1324
1279	frames), we crop the map with trajectory-aware padding:	given to views in which the target category is absent. When	1325
1280	5 m around the agent position for rotation trajectories and	too few such views exist, they are supplemented with addi-	1326
1281	4 m around the traversed path for multi-camera trajectories.	tional non-adjacent views to maintain viewpoint diversity.	1327
1282	Object counts in the cropped BEV are validated against the	This makes the final count require aggregation across mul-	1328
1283	segmentation maps to ensure consistency.	multiple views rather than inspection of a single image.	1329
1284	<b>Object filtering and category balancing.</b> Structural ele-	<b>Top-view supervision and question generation.</b> Each	1330
1285	ments (walls, floors, ceilings, doorways) are excluded from	sample also stores the canonical top-view image as the rea-	1331
1286	counting targets. Target objects must be visible in both the	soning target. In the exported JSONL format, image inputs	1332
1287	first-person frames and the cropped top-down segmentation	use centered crops derived from the union of all annotated	1333
1288	map (with a minimum coverage of 0.1%). Because initial	object boxes in each selected camera view, which reduces	1334
1289	generation heavily favors count= 1 questions (~82%), we	empty borders while preserving the visible object layout.	1335
1290	apply iterative rebalancing: high-frequency categories are	Questions are instantiated from a diverse pool of natural-	1336
1291	capped at 9.9% of the dataset, and count= 1 samples are	language counting templates such as “How many {object}	1337
1292	downsampled. This produces a more uniform distribution	are in this scene?”	1338
1293	across object categories and count values.	<b>8.3.3. ScanNet++</b>	1339
1294	<b>Question format and distractor generation.</b> Questions	We further construct a real indoor multiview counting	1340
1295	follow the template: “How many {category}(s) are in this	set from ScanNet++ [41], using iPhone image trajectories	1341
1296	area?” with four answer choices (A–D). Distractors are	paired with labeled 3D scene reconstructions. Compared	1342
1297	sampled from $\pm 1$ and $\pm 2$ of the correct count, producing	with MessyTable, this set covers larger indoor spaces, more	1343
1298	plausible alternatives. Negative counts are removed, and	varied viewpoints, and more realistic household layouts.	1344
1299	the four options are shuffled with a per-sample determinis-	<b>Top-down map and candidate view generation.</b> For	1345
1300	tic seed to eliminate positional bias.	each scene, we first generate a top-down map from the	1346
1301	<b>Statistics.</b> The synthetic training set contains 17,079 ex-	labeled 3D reconstruction. Because raw point-cloud ren-	1347
1302	amples generated from ProcTHOR [10] scenes, covering	derings are visually sparse and unrealistic, we further use	1348
1303	both trajectory types.	Qwen-Image-Edit [32] to transform the rendered top-down	1349
1304	<b>8.3.2. MessyTable</b>	visualization into a more realistic top-down image while	1350
1305	To expose the model to real tabletop imagery with severe	preserving the scene layout. We then build a candidate ego-	1351
1306	clutter and occlusion, we construct an additional multiview	centric view pool by combining a small set of canonical	1352
1307	counting set from MessyTable [4]. Each scene contains	iPhone views with additional randomly sampled frames, re-	1353
1308	multiple camera views of the same tabletop arrangement to-	quiring each extra frame to differ from the canonical views	1354
1309	gether with instance-level annotations.	by at least a minimum yaw angle.	1355
1310	<b>Scene-level counting targets.</b> For each scene, we aggre-	<b>Visibility estimation and target selection.</b> We estimate	1356
1311	gate annotations across all cameras and de-duplicate in-	which object instances are visible in each candidate frame	1357
1312	stance IDs across views, so that the ground-truth answer	by projecting the labeled 3D scene into the camera views	1358
1313	corresponds to the number of unique physical objects rather	using mesh ray-casting, and use the semantic annotations	1359
1314	than the sum of per-view detections. Counting targets are	to obtain scene-level category counts. Top-down maps are	1360
1315	defined at the subclass level and mapped to readable cate-	filtered by automatic quality rules to remove blurry, blank,	1361
1316	gory names.	or low-texture renderings. Candidate counting targets are	1362
1317	<b>Target and view sampling.</b> For each scene, we sample	restricted to non-structural object categories with bounded	1363
1318	one target category from the categories present in the scene.	scene-level counts and sufficient visible support in the	1364
1319	To reduce the dominance of trivial singleton cases, sam-	candidate views. To avoid metadata leakage, target selection	1365
1320	pling is biased toward categories with count $\geq 2$ : when both	is performed in a blind setting: the model is shown the top-	1366
1321	singleton and multi-instance categories are available, 90%	down image and candidate labels, but not their annotated	1367
1322	of samples are drawn from the multi-instance bucket and	counts, and we keep only categories whose visually pre-	1368
		dicted count matches the ground truth. When multiple valid	1369
		categories remain, we preferentially sample categories with	1370
		counts greater than one.	1371



## Path Tracing: Inference with Imaginative Perception Tokens

**Question:** You are moving along the path 2→1 as shown in the top-down view <image\_1>. Looking ahead from point 2, you see: <image\_2>. Looking ahead from point 1, you see: <image\_3>. Which object can you see on your **left** side at waypoint M1?

A: the sink basin    **B: the garbage can ✓**    C: the towel    D: the cloth    **✓ Correct**



"My left view at the midpoint M1:"

**Question:** You are navigating from waypoint 2 to waypoint 1 and passing through M1, as shown in the top-down view <image\_1>. Here's your view looking forward from point 2: <image\_2>. And from point 1: <image\_3>. Which object can you see on your **left** side at waypoint M1?

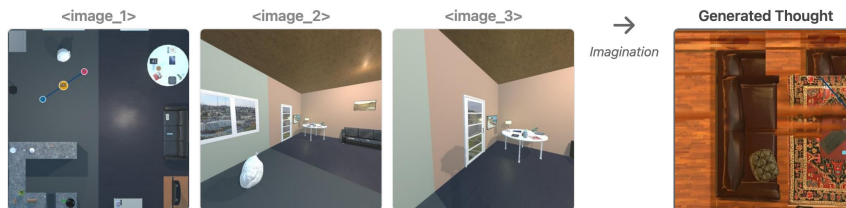
**A: the countertop ✓**    B: the chair    C: the dining table    D: the garbage bag    **✓ Correct**



"At midpoint M1, looking to the left, I see:"

**Question:** You are moving along the path 2→1, and you should refer to the top-down view at <image\_1>. Looking ahead from point 2, you see: <image\_2>. Looking ahead from point 1, you see: <image\_3>. Which object can you see on your **right** side at waypoint M1?

A: a sink basin    **B: a garbage can ✓**    **C: a pot ✗**    D: a pepper shaker    **✗ Incorrect**



"At midpoint M1, looking to the right, I see:"

Figure 4. **Path tracing with imaginative perception tokens (EgoDir setting).** The model receives a top-down map (<image\_1>) and egocentric views at the two endpoints (<image\_2>, <image\_3>), generates a visual thought (imagined sideview at  $M_1$ ), and predicts an answer. Although the generated thoughts exhibit spatial imprecision and artifacts, the model still arrives at the correct answer in the first two examples, suggesting that imagination training encourages internalized spatial reasoning rather than reliance on pixel-accurate intermediate outputs. The third row shows a failure case. Correct answers are highlighted in green; incorrect predictions in red.

1468 discrete imaginative perception tokens directly to the language model vocabulary of state-of-the-art vision-language models. Inspired by Aurora [2], we first trained a VQ-

VAE from scratch on the intermediate RGB images in our datasets, including novel viewpoint renders, top-down BEV maps, and sideview images. However, the reconstruction

1471  
1472  
1473

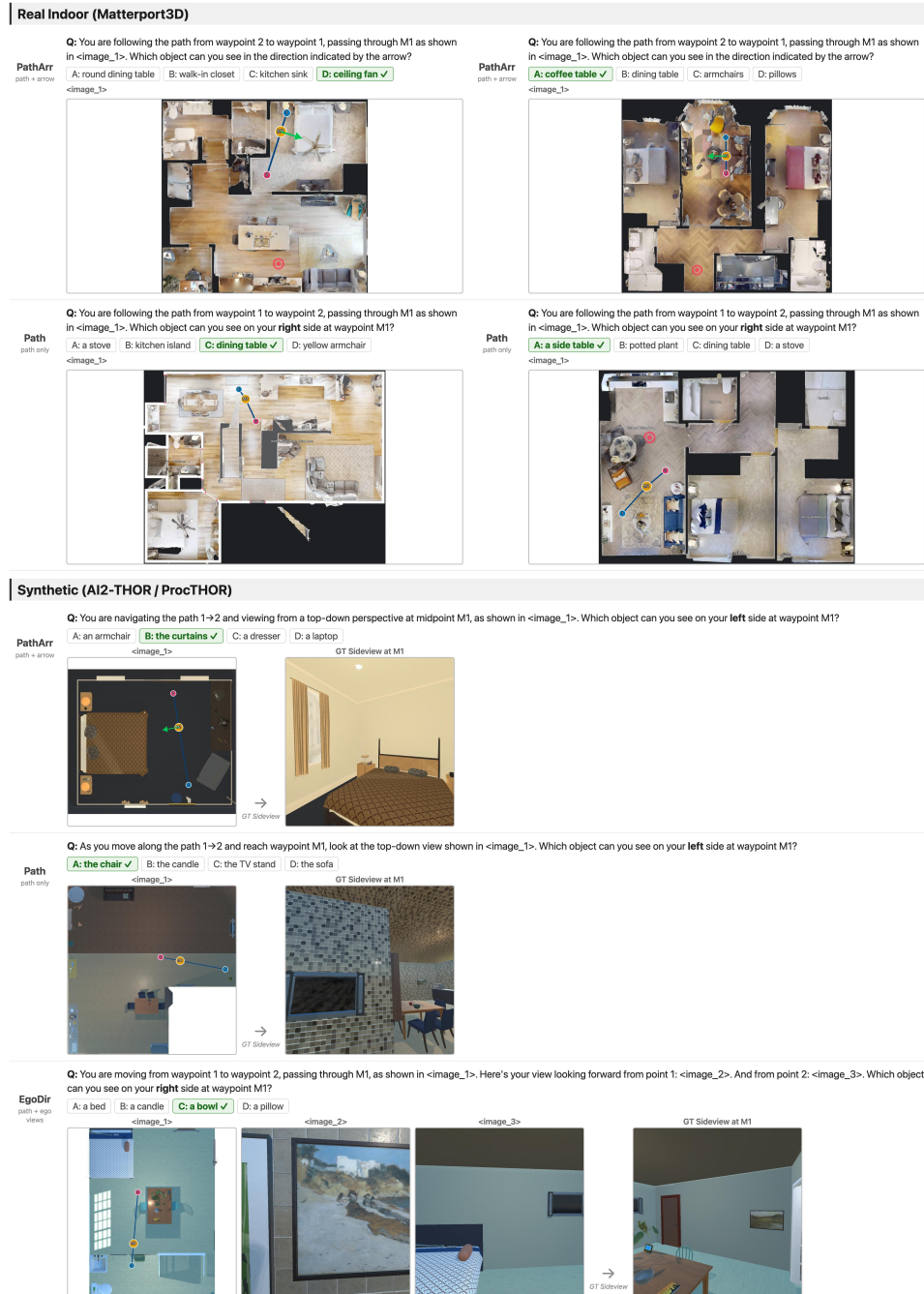


Figure 5. **Path tracing dataset examples.** Top: real-world examples from Matterport3D in the PathArr and Path settings. Bottom: synthetic examples from AI2-THOR/ProcTHOR in the PathArr, Path, and EgoDir settings, with ground-truth sideviews at midpoint  $M_1$  shown on the right. Each example shows the input image(s), question, and four answer choices with the correct answer highlighted.

1474 quality of these simple VQ-VAEs was insufficient for su-  
1475 pervising models on their intermediate token sequences as  
1476 image outputs.

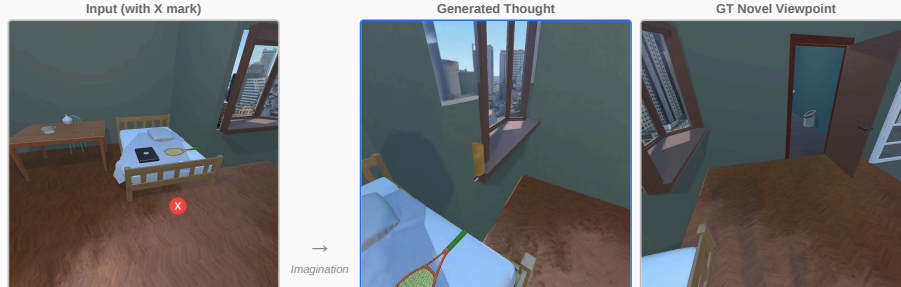
1477 We therefore switched to off-the-shelf pretrained VQ-  
1478 GANs [14] with varying configurations. These configura-

tions differ along two axes: codebook size (e.g., 1K, 8K, 1479  
and 16K entries) and spatial downsampling ratio ( $f = 8$  vs. 1480  
 $f = 16$ ). Each choice involves a tradeoff: a larger codebook 1481  
improves representational fidelity but inflates the model voc- 1482  
abulary, while a smaller downsampling ratio yields higher 1483

## Perspective Taking: Inference with Imaginative Perception Tokens

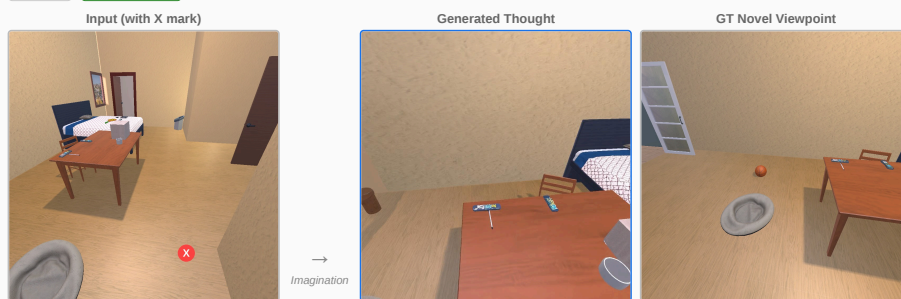
**Question:** If someone at this viewpoint moves to the 'X' on the ground while maintaining their orientation, then rotates right 90 degrees, will the window be closer or further to them?

Distance — Closer

A: Further  B: Closer  ✓ Correct

**Question:** This image shows my current perspective. If I move straight to the 'X' marked point on the ground while keeping my orientation, then turn left for 90 degrees, will the chair be on my left or right?

Position — Left - Right

A: left  B: right  ✓ Correct

**Question:** After reaching the 'X' position on the ground while keeping my current orientation, then making a 90-degree turn left, on which side will the sofa be - left or right?

Position — Left - Left

A: right  X B: left  ✓ X Incorrect

Figure 6. **Perspective taking with imaginative perception tokens.** The model receives an input view with an “X” mark on the ground, imagines the scene from the target viewpoint, and predicts spatial relationships. Generated thoughts are compared against ground-truth novel viewpoints. Correct answers are highlighted in green; incorrect predictions in red.

1484 reconstruction quality at the cost of a longer token sequence  
1485 per image, increasing context length. Figure 10 illustrates  
1486 reconstruction quality across these settings.

1487 We selected Qwen2.5-VL [32] in two sizes (3B and 7B)  
1488 as our backbone for discrete token finetuning experiments.

We note that the training data at this stage was of lower  
1489 quality than our final datasets described in main text; these  
1490 experiments were intended solely to probe whether discrete  
1491 imaginative perception tokens can serve as a useful inter-  
1492 mediate, not to achieve peak performance. We focused on  
1493



Figure 7. **Perspective taking dataset examples.** Examples from AI2-THOR/ProctHOR across the four sub-categories. Each example shows the input view with “X” mark and the ground-truth novel viewpoint. The correct answer is highlighted in green.

1494 Path Tracing (PT) and Perspective Taking (PET).

1495 For each model we trained three variants: answer-only  
1496 finetuning, Text CoT, and image chain-of-thought with discrete  
1497 IPTs, where the VQGAN codebook tokens are appended to the model  
1498 vocabulary and the model first autoregressively generates the imaginative  
1499 perception token se-

quence before predicting the final answer plus a zeroshot  
1500 baseline. For IPT variants we tested two VQGAN settings  
1501 that keep sequence length manageable: CB 16K  $f=16$  and  
1502 CB 1K  $f=16$ . Results are shown in Table [II](#).  
1503

IPT consistently outperforms both answer-only finetuning  
1504 and Text CoT on Path Tracing, with CB 1K  $f=16$  yield-  
1505

1506 ing the best results for both model sizes (55.0 for 3B, 55.9  
1507 for 7B). On Perspective Taking, gains are modest and near  
1508 the zero-shot baseline, suggesting that lower data quality  
1509 and the representational limitations of discrete token recon-  
1510 struction are a bottleneck for this more visually demand-  
1511 ing task. A substantial gap remains relative to our final  
1512 BAGEL-based results, motivating the move to a unified  
1513 model.

1514 Figure 11 shows ground-truth imagination images along-  
1515 side the corresponding IPTs decoded from the Qwen2.5-  
1516 VL 3B model. The decoded outputs are visually degraded,  
1517 lacking the spatial structure and object detail present in the  
1518 ground truth, which helps explain the remaining perfor-  
1519 mance gap and further motivated our switch to continuous  
1520 latent representations.

1521 We further investigated alternative intermediate image  
1522 representations. Instead of training the model to generate  
1523 imaginative perception tokens for RGB thought images, we  
1524 replaced them with tokens for grayscale images and tokens  
1525 for pseudo depth maps obtained from the DepthAnything  
1526 model [36]. The intuition is that simplifying the generation  
1527 target, from full RGB to grayscale, reduces the difficulty of  
1528 the token prediction task and may improve spatial reasoning  
1529 downstream. Results in Table 12 show that switching from  
1530 RGB to grayscale does boost performance (55.0→59.6 on  
1531 PT, 50.0→55.5 on PET), while depth tokens perform com-  
1532 parably to RGB. Nevertheless, a substantial gap remains,  
1533 and Figure 12 shows that the decoded grayscale outputs  
1534 are still visually degraded, indicating that generation quality  
1535 rather than representation type is the primary bottleneck.

1536 Figure 12 shows ground-truth imagination images along-  
1537 side the corresponding IPTs decoded from the Qwen2.5-  
1538 VL 3B model. The decoded outputs are visually degraded,  
1539 lacking the spatial structure and object detail present in the  
1540 ground truth, which helps explain the remaining perfor-  
1541 mance gap and further motivated our switch to continuous  
1542 latent representations.

1543 These findings collectively motivated us to move away  
1544 from discrete token generation in non-unified VLMs and in-  
1545 stead adopt a unified model: BAGEL, that natively supports  
1546 interleaved image understanding and generation through  
1547 continuous latent representations.

Table 10. **Path tracing per-split results.** Accuracy (%) broken down by input setting. The main paper reports the average across these splits. For our models, accuracy reports the maximum between answer-only and free-generation inference. Best per group in **bold**.

Model	AI2-THOR			Different Env.	
	EgoDir	Path	PathArr	Real	Real+Arr
<i>VQA Models</i>					
GPT-5	<b>61.1</b>	<b>56.8</b>	<b>62.6</b>	<b>74.5</b>	87.3
GPT-5.2	22.1	40.2	36.3	57.5	68.4
Gemini 2.5 Flash	45.1	37.9	41.5	58.0	84.8
Gemini 3 Flash	48.7	39.1	39.2	70.1	<b>96.2</b>
InternVL3.5-8B	43.4	29.0	35.1	45.4	49.4
Qwen2.5-VL-7B	44.2	32.5	35.1	47.1	42.4
Qwen3-VL-8B	31.9	33.7	42.1	52.9	75.3
<i>Unified Models</i>					
Janus-Pro-7B	36.3	30.2	33.9	34.5	36.1
Chameleon 7B	5.3	23.7	19.9	23.0	25.9
<i>Ours (fine-tuned BAGEL)</i>					
Bagel (base)	36.3	26.0	27.5	39.7	45.6
Bagel (label-only)	<b>73.5</b>	61.5	62.0	46.6	62.7
+ Text CoT	53.1	47.9	48.0	<b>52.3</b>	51.3
+ IPT	61.1	43.2	42.7	46.6	<b>68.4</b>
+ Mixed Training	71.7	<b>65.1</b>	<b>63.2</b>	50.6	66.5

Table 11. **Discrete IPT experiments on Qwen2.5-VL.** Accuracy (%) on Path Tracing (PT) and Perspective Taking (PET). Training data at this stage was of lower quality than the final datasets. “-” denotes experiments not conducted.



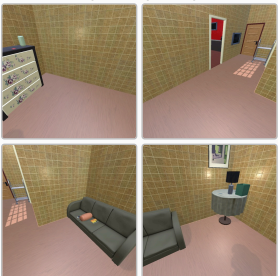
Model	Method	PT	PET
Qwen2.5-VL 3B	Zero-shot	33.0	50.0
	Answer-only	48.6	50.5
	Text CoT	43.0	-
	IPT (CB 16K, $f=16$ )	50.5	48.5
	IPT (CB 1K, $f=16$ )	<b>55.0</b>	50.0
Qwen2.5-VL 7B	Zero-shot	38.5	47.2
	Answer-only	37.6	-
	Text CoT	35.7	-
	IPT (CB 16K, $f=16$ )	55.0	-
	IPT (CB 1K, $f=16$ )	<b>55.9</b>	-

**Multiview Counting: Inference with Imaginative Perception Tokens**

**Question: How many sofa(s) are in this area?** Rotation — 4 cardinal dirs

A: 4 **B: 1 ✓** C: 2 D: 3 ✓ Correct

Input Views (4 frames) Generated Thought GT Top-Down Map





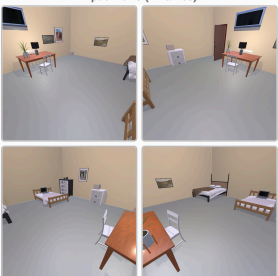
Imagination

🗨️ "The top-down view of this area is:"

**Question: How many bed(s) are in this area?** Multi-camera — square path

A: 2 ✓ B: 3 C: 1 D: 4 ✓ Correct

Input Views (4 frames) Generated Thought GT Top-Down Map




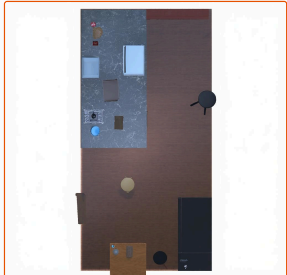

Imagination

🗨️ "The top-down view of this area is:"

**Question: How many stool(s) are in this area?** Rotation — 4 cardinal dirs

A: 2 ✗ **B: 3 ✓** C: 4 D: 1 ✗ Incorrect

Input Views (4 frames) Generated Thought GT Top-Down Map



Imagination

🗨️ "The top-down view of this area is:"

Figure 8. **Multiview counting with imaginative perception tokens.** The model receives four egocentric views, imagines a top-down BEV map, and counts target objects. Generated thoughts are compared against ground-truth top-down maps. Correct answers are highlighted in green; incorrect predictions in red.

**Multiview Counting — AI2-THOR (ProcTHOR) — Rotation Trajectory**

Q: Count the number of coffee-table(s) visible in the explored area.

Rotation  A.4  B.2  C.3  B.3 ✓

Input Views (4 frames) GT Top-Down Map

---

Q: How many tostand(s) are in this area?

Rotation  A.3  B.2  C.4  B.3 ✓

Input Views (4 frames) GT Top-Down Map

---

**Multiview Counting — AI2-THOR (ProcTHOR) — Multi-Camera Trajectory**

Q: How many chair(s) are in this area?

Multi-cam  A.6 ✓  B.5  C.7  D.4

Input Views (8 frames) GT Top-Down Map

---

**Multiview Counting — MessyTable (Real Tabletop Photos)**

Q: Please compute the number of banana in this scene.

Real Photo  Answer: 3

Input Views (2 frames) GT Reasoning Image

---

Q: Can you count the vacuum-packed bagged snack in this scene?

Real Photo  Answer: 2

Input Views (4 frames) GT Reasoning Image

---

**Multiview Counting — ScanNet (Real Indoor Scenes)**

Q: What is the total number of office chairs in this scene?

Real Photo  Answer: 8

Input Views (4 frames) GT Reasoning Image

---

Q: What is the total count of bed pillows in this scene?

Real Photo  Answer: 3

Input Views (4 frames) GT Reasoning Image

Figure 9. **Multiview counting dataset examples.** Examples from AI2-THOR/ProcTHOR showing both rotation (top) and multi-camera (bottom) trajectories. Each example shows four input views in a  $2 \times 2$  grid and the ground-truth top-down map. The correct answer is highlighted in green.

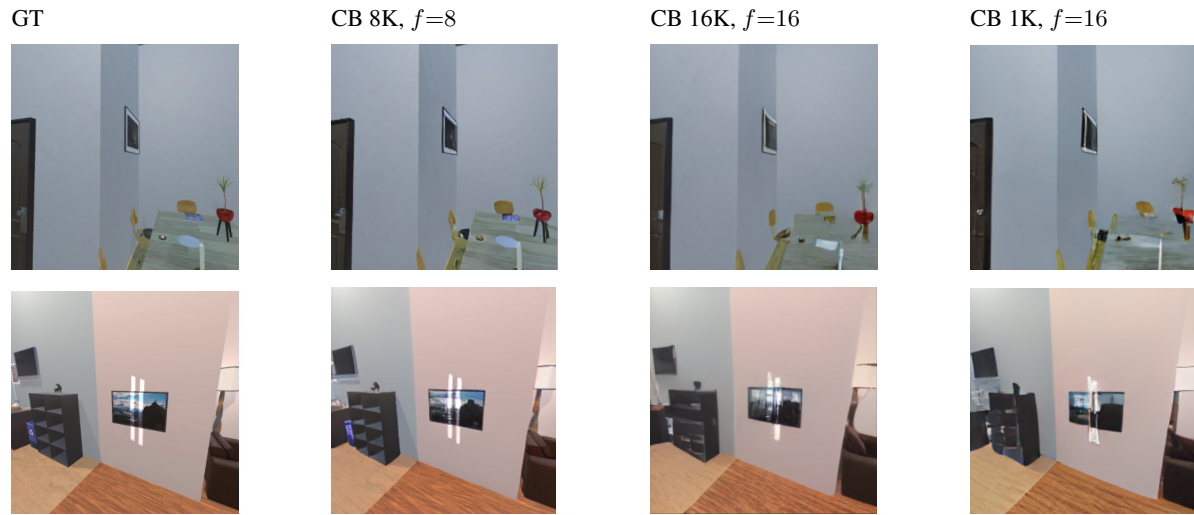


Figure 10. **VQGAN reconstruction quality across codebook and downsampling settings.** Larger codebooks and smaller  $f$  improve fidelity but increase vocabulary size and sequence length respectively.



Figure 11. **Ground-truth vs. decoded IPTs from Qwen2.5-VL 3B.** The model-generated imagination tokens decode into visually degraded images that fail to preserve the spatial structure of the ground truth, highlighting the limitations of discrete token generation in non-unified VLMs.

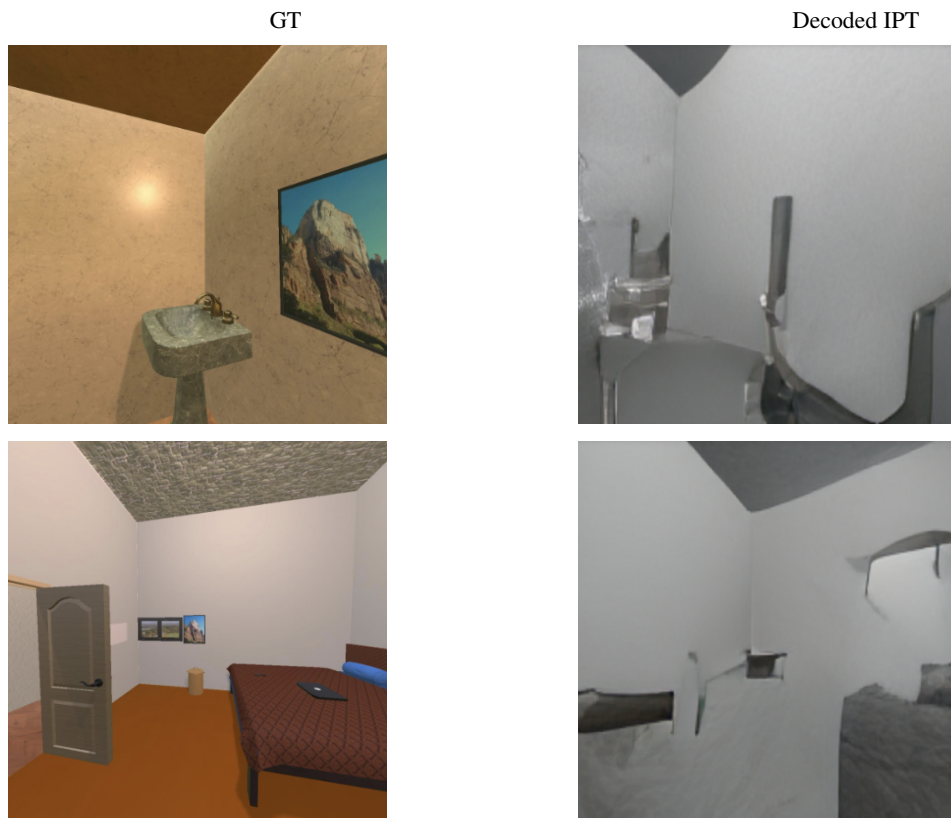


Figure 12. **Decoded IPTs from Qwen2.5-VL 3B for grayscale representations.** Despite the simpler generation target, grayscale decoded outputs remain visually degraded and fail to preserve spatial structure.

Table 12. **Effect of intermediate image representation on Qwen2.5-VL 3B.** Accuracy (%) on Path Tracing (PT) and Perspective Taking (PET).

<b>Method</b>	<b>PT</b>	<b>PET</b>
IPT RGB (CB 1K, $f=16$ )	55.0	50.0
IPT Grayscale (CB 1K, $f=16$ )	<b>59.6</b>	<b>55.5</b>
IPT Depth (Aurora VQVAE)	55.0	54.7