

# CONVERSATIONAL ORIENTATION REASONING: EGOCENTRIC-TO-ALLOCENTRIC NAVIGATION WITH MULTIMODAL CHAIN-OF-THOUGHT

**Anonymous authors**

Paper under double-blind review

## ABSTRACT

Conversational agents must translate egocentric utterances (e.g., “on my right”) into allocentric orientations (N/E/S/W). This challenge is particularly critical in indoor or complex facilities where GPS signals are weak and detailed maps are unavailable. While chain-of-thought (CoT) prompting has advanced reasoning in language and vision tasks, its application to multimodal spatial orientation remains underexplored. We introduce **Conversational Orientation Reasoning (COR)**, a new benchmark designed for Traditional Chinese conversational navigation projected from real-world environments, addressing egocentric-to-allocentric reasoning in non-English and ASR-transcribed scenarios. We propose a multimodal chain-of-thought (MCoT) framework, which integrates ASR-transcribed speech with landmark coordinates through a structured three-step reasoning process: (1) extracting spatial relations, (2) mapping coordinates to absolute directions, and (3) inferring user orientation. A curriculum learning strategy progressively builds these capabilities on Taiwan-LLM-13B-v2.0-Chat, a mid-sized model representative of resource-constrained settings. Experiments show that MCoT achieves 100% orientation accuracy on clean transcripts and 98.1% with ASR transcripts, substantially outperforming unimodal and non-structured baselines. Moreover, MCoT demonstrates robustness under noisy conversational conditions, including ASR recognition errors and multilingual code-switching. The model also maintains high accuracy in cross-domain evaluation and resilience to linguistic variation, domain shift, and referential ambiguity. These findings highlight the potential of structured MCoT spatial reasoning as a path toward interpretable and resource-efficient embodied navigation.

## 1 INTRODUCTION

Humans naturally describe spatial environments in *egocentric* (agent-centric) terms (e.g., “The exit is on my right”), whereas navigation systems typically operate on *allocentric* (world-centric) orientations such as north, south, east, and west. Conversational navigation has emerged as a promising paradigm that enables users to specify goals through dialogue, offering a natural and human-centered means of guidance in unfamiliar environments Sundar et al. (2024); Sheshadri & Hara (2024); Kaniwa et al. (2024); Liu et al. (2024); Levi & Kadar (2025). However, the crucial problem of grounding egocentric language into allocentric orientation remains underexplored. Current approaches typically assume access to GPS, detailed maps, or fixed global frames de Vries et al. (2018); Chen et al. (2020), and have concentrated primarily on English-based scenarios. Recent progress has also relied heavily on large-scale models Ghosh et al. (2024); Tang et al. (2023), which show strong reasoning abilities but demand substantial computational resources, hindering deployment in resource-constrained settings such as mobile navigation and edge devices.

Research in embodied AI and MCoT has advanced vision-language navigation and action planning Mu et al. (2023); Sun et al. (2024); Liu et al. (2025); Shen et al. (2025); Pareek et al. (2024), but orientation reasoning from natural language has been largely overlooked. These approaches typically assume that the agent’s orientation is already known or operate on high-level action spaces rather than inferring fundamental spatial relationships de Vries et al. (2018); Chen et al. (2020). Meanwhile, large audio-language models (LALMs) Zhang et al. (2023); Xie & Wu (2024); Fu et al.

(2025); Défossez et al. (2024) have advanced speech understanding and dialogue Tang et al. (2024); Gong et al. (2023); Ghosh et al. (2024); Kong et al. (2024), yet their reasoning abilities remain limited to perception-level tasks such as transcription or summarization Yu Huang et al. (2024); Yang et al. (2024); Wang et al. (2025); Shi et al. (2025). While recent efforts like Audio-CoT Ma et al. (2025) show promise for enhanced speech-based reasoning, the challenge of transforming egocentric spatial descriptions into allocentric orientation inference remains unaddressed.

To address this gap, we introduce **Conversational Orientation Reasoning (COR)**, a new benchmark for egocentric-to-allocentric orientation reasoning in Traditional Chinese conversational navigation. COR is derived from real-world urban transportation environments in Taiwan, projected into structured grid representations. Unlike prior studies that rely on vision or raw audio, COR combines ASR-transcribed egocentric language with structured landmark coordinates, evaluating both clean text and ASR transcripts to simulate realistic recognition errors. COR addresses the lack of non-English benchmarks in multimodal spatial reasoning, particularly under noisy ASR conditions.

Our study is guided by three research questions:

- **RQ1 (Effectiveness):** How effective is multimodal CoT prompting for orientation reasoning compared to unimodal and unstructured baselines?
- **RQ2 (Component analysis):** What are the contributions of ASR preprocessing, multimodal fusion, and structured CoT steps?
- **RQ3 (Robustness and generalization):** How robust is the approach to linguistic variation, and how well does it generalize across different spatial domains?

Our contributions are as follows:

1. **Task and benchmark.** We introduce the COR benchmark for egocentric-to-allocentric orientation reasoning, combining ASR-transcribed speech with landmark coordinates.
2. **Framework.** We develop a multimodal CoT framework with a structured three-step reasoning process that integrates noisy transcripts with spatial signals for orientation inference.
3. **Evaluation.** We provide extensive experiments in Traditional Chinese, demonstrating effectiveness, component contributions, and robustness validation across linguistic variation, cross-domain generalization, and referential ambiguity beyond English-centric research.

## 2 RELATED WORK

**Navigation and Spatial Orientation.** Natural language navigation tasks have long driven progress in embodied AI. The Room-to-Room (R2R) benchmark Anderson et al. (2018) first established visually grounded navigation instructions in real indoor environments, while Cooperative Vision-and-Dialog Navigation (CVDN) Thomason et al. (2019) extended this to conversational settings where agents interpret multi-turn human dialogues. Talk the Walk de Vries et al. (2018) grounds tourist utterances with masked attention but assumes orientation-agnostic actions. Touchdown Chen et al. (2020) extends navigation to urban environments with graph-based orientation, while SpatialRGPT Cheng et al. (2024) leverages 3D scene graphs for spatial grounding. Speaker-Follower models Fried et al. (2018) improve instruction following by embedding panoramic orientation, and Ego4D Grauman et al. (2022) provides large-scale egocentric video datasets for studying orientation and attention. Despite these advances, existing frameworks assume known agent orientations or operate on high-level action spaces such as Left, Right, Up, and Down rather than addressing the fundamental egocentric-to-allocentric transformation challenge de Vries et al. (2018); Fried et al. (2018). In contrast, our work tackles the egocentric-to-allocentric challenge, using MCoT to infer absolute directions from noisy ASR transcripts and landmark coordinates.

**Chain-of-Thought and Multimodal Reasoning.** CoT prompting has proven effective for enhancing reasoning in large language models (LLMs). Increasing the number of reasoning steps in demonstrations improves accuracy across diverse benchmarks Jin et al. (2024), and subsequent studies show that performance critically depends on the logical consistency and contextual relevance of rationales Wang et al. (2023); Prystawski et al. (2023); Tang et al. (2023). Recent extensions to MCoT enable embodied agents to reason jointly over multiple modalities. For instance, EmbodiedGPT Mu

et al. (2023) and E-CoT Lin et al. (2024) segment tasks into subgoals, while Emma-X Sun et al. (2024) incorporates grounded planning and predictive movement. SpatialCoT Liu et al. (2025) focuses on coordinate alignment, and MCoCoNav Shen et al. (2025) integrates semantic maps for multi-robot collaboration. Together, these advances highlight the growing importance of MCoT for embodied reasoning. However, none of these approaches tackles the transformation from egocentric natural language descriptions to allocentric orientation inference under noisy ASR conditions.

**Large Audio-Language Models.** Large audio-language models (LALMs) extend this line of research by incorporating speech as an additional input. They have advanced transcription, classification, and interactive dialogue Zhang et al. (2023); Xie & Wu (2024); Fu et al. (2025); Défossez et al. (2024), showing strong performance on perception-level tasks such as speech recognition and summarization Tang et al. (2024); Gong et al. (2023); Ghosh et al. (2024); Kong et al. (2024); yu Huang et al. (2024); Yang et al. (2024); Wang et al. (2025); Shi et al. (2025). Audio-CoT Ma et al. (2025) suggests that CoT can aid speech-derived reasoning, but the improvements remain modest and fail to generalize to complex reasoning tasks. While existing LALMs excel at perception, they lack mechanisms for structured multi-step reasoning, limiting their effectiveness in noisy or ambiguous conditions. Our work addresses this by treating ASR transcripts as a noisy textual modality, integrating them with landmark coordinates, and applying MCoT to enable orientation inference.

### 3 METHOD

#### 3.1 OVERVIEW

We propose a MCoT framework for egocentric-to-allocentric orientation reasoning in conversational navigation. To simulate realistic speech-driven conditions, we synthesize speech from clean descriptions and transcribe it using automatic speech recognition (ASR), thereby introducing natural transcription errors. This controlled approach allows us to systematically evaluate performance under varying levels of ASR noise while maintaining reproducible experimental conditions. The resulting transcript  $A'$  is combined with structured spatial coordinates  $T$  that describe the user’s position and nearby landmarks. The model then generates both an interpretable reasoning trace and the final allocentric orientation prediction  $D^* \in \{\text{north, east, south, west}\}$ .

As illustrated in Figure 1, the framework consists of three modules: (1) *Speech synthesis and transcription*, where original descriptions are converted into transcripts via text-to-speech (TTS) and ASR; (2) *Multimodal input preparation*, which encodes transcripts and spatial coordinates into a unified representation; and (3) *Orientation reasoning*, where the model applies structured CoT inference to derive the final allocentric prediction. Unlike prior LALMs that emphasize transcription quality, our framework focuses on spatial reasoning under noisy transcripts.

#### 3.2 TASK

Let  $A$  denote the clean egocentric description,  $A'$  the ASR transcript,  $T$  the spatial coordinate input, and  $D^*$  the ground-truth allocentric orientation. We formalize orientation reasoning on a discrete grid derived from the Gongguan MRT area, a busy transportation hub in Taiwan (Figure 2).

**Environment.** The real-world area is projected into a  $10 \times 10$  grid  $\mathcal{G}$ , with user position  $u = (x_u, y_u) \in \mathcal{G}$  and landmarks  $\mathcal{L} = \{\ell_i\}$  with coordinates  $p(\ell_i) \in \mathcal{G}$ .

**Egocentric description.** The user describes their context with an egocentric relation  $q \in \{\text{FRONT, BACK, LEFT, RIGHT}\}$  and a reference landmark  $\ell_r$ . After TTS and ASR processing, we obtain transcript  $A'$ . For example: “*I am at Exit 2, and restaurant 5 is on my right*”  $\rightarrow D^* = \text{north}$ , which implies that the user is facing north.

**Mapping rule.** We first compute the relative vector between the landmark and the user,

$$\Delta = p(\ell_r) - u = (\Delta_x, \Delta_y).$$

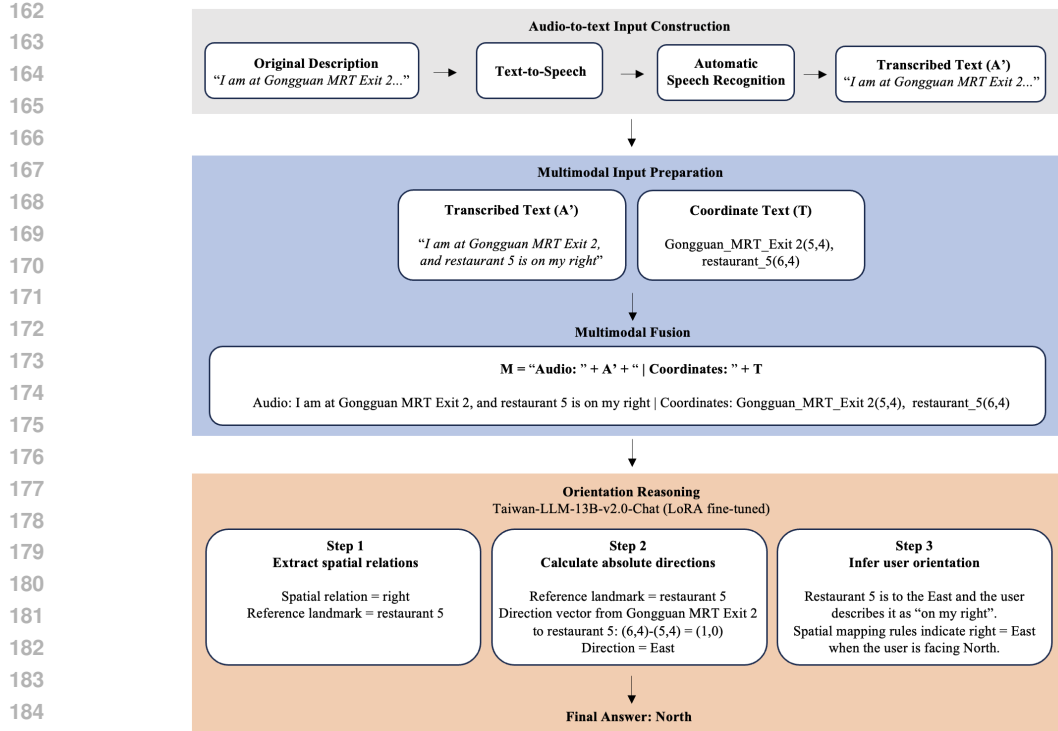


Figure 1: Pipeline of our MCot framework. It consists of three modules: (1) speech synthesis and transcription, (2) multimodal input preparation and fusion, and (3) orientation reasoning.

The absolute direction of  $\ell_r$  is then defined as

$$\text{AbsDir}(\Delta) = \begin{cases} \text{E} & \text{if } |\Delta_x| > |\Delta_y| \wedge \Delta_x > 0, \\ \text{W} & \text{if } |\Delta_x| > |\Delta_y| \wedge \Delta_x < 0, \\ \text{N} & \text{if } |\Delta_y| > |\Delta_x| \wedge \Delta_y > 0, \\ \text{S} & \text{if } |\Delta_y| > |\Delta_x| \wedge \Delta_y < 0, \end{cases}$$

and the user’s orientation  $D^*$  is derived by rotating  $d_{\text{abs}} = \text{AbsDir}(\Delta)$  according to relation  $q$ :

$$D^* = \begin{cases} d_{\text{abs}} & q = \text{FRONT}, \\ \text{Rot}(d_{\text{abs}}, 180^\circ) & q = \text{BACK}, \\ \text{Rot}(d_{\text{abs}}, +90^\circ) & q = \text{LEFT}, \\ \text{Rot}(d_{\text{abs}}, -90^\circ) & q = \text{RIGHT}. \end{cases}$$

Since the environment is represented as a discrete  $10 \times 10$  grid with only four cardinal neighbors  $(0, \pm 1)$ ,  $(\pm 1, 0)$ , i.e., a Manhattan grid, landmarks are always axis-aligned and diagonal cases do not occur. All ground-truth orientations  $D^*$  are automatically derived from these mapping rules. Each instance includes a step-by-step reasoning trace used to supervise CoT generation. To ensure quality, all automatically generated instances are verified by human annotators for correctness.

Table 1 lists the full set of relative-to-absolute mapping rules.

### 3.3 TRAINING

In preliminary experiments, directly training the model end-to-end led to unstable learning and poor generalization. We therefore adopt curriculum learning with stage-wise fine-tuning on Taiwan-LLM-13B-v2.0-Chat to incrementally build orientation reasoning capabilities. All training stages use clean text, while noisy ASR transcripts are introduced only during evaluation. This design choice allows the model to first master fundamental spatial reasoning without transcription noise, then tests performance when ASR errors are introduced.

	Facing	Front	Back	Right	Left
North	N	S	E	W	
East	E	W	S	N	
South	S	N	W	E	
West	W	E	N	S	

Table 1: Relative-to-absolute direction mapping rules.

(0, 9) Hospital	(1, 9) Beverage Shop 2	(2, 9) Band Practice Room	(3, 9) Photo Studio	(4, 9) Bookstore	(5, 9) Underground Exit	(6, 9) Fu Garden	(7, 9) NTU Main Gate	(8, 9) University Museum	(9, 9) Art Museum	↑ N
(0, 8) Convenience Store 3	(1, 8) Optical Shop 2	(2, 8) Sports Store 2	(3, 8) Office Building	(4, 8) Gongguan MRT Exit 4	(5, 8) Gongguan MRT Exit 3	(6, 8) Dormitory 2	(7, 8) Small Plaza 3	(8, 8) Cooperative Store	(9, 8) Fu Bell	
(0, 7) Foundation	(1, 7) Shoe Store	(2, 7) Accessory Shop	(3, 7) Law Firm	(4, 7) Pharmacy	(5, 7) Bike Rental Station 7	(6, 7) Post Office 2	(7, 7) Dormitory 3	(8, 7) Theater	(9, 7) Student Affairs Office	
(0, 6) High School	(1, 6) Restaurant 6	(2, 6) Convenience Store 2	(3, 6) Coffee Shop 2	(4, 6) Gym	(5, 6) Bus Stop 2	(6, 6) Dormitory 1	(7, 6) Dormitory 5	(8, 6) Coffee Shop 3	(9, 6) Bike Rental Station 6	
(0, 5) Coffee Shop 1	(1, 5) Temple	(2, 5) Beverage Shop 1	(3, 5) Music Studio	(4, 5) Bike Rental Station 3	(5, 5) Bike Rental Station 4	(6, 5) Convenience Store 1	(7, 5) Bike Rental Station 5	(8, 5) Bank 3	(9, 5) Science Center	
(0, 4) Hotel 2	(1, 4) Mobile Phone Shop	(2, 4) Sports Store 1	(3, 4) Restaurant 4	(4, 4) Gongguan MRT Exit 1	(5, 4) Gongguan MRT Exit 2	(6, 4) Restaurant 5	(7, 4) Dormitory 6	(8, 4) Li-Xian Building	(9, 4) Animal Museum	
(0, 3) Printing Shop	(1, 3) Flea Market	(2, 3) Restaurant 3	(3, 3) Hair Salon	(4, 3) Bus Stop 1	(5, 3) Elementary School	(6, 3) Kindergarten	(7, 3) Security Office	(8, 3) Academic Building E	(9, 3) Bank 2	
(0, 2) Shopping Mall	(1, 2) Parking Lot 3	(2, 2) Post Office 1	(3, 2) Fast Food Restaurant	(4, 2) Market	(5, 2) Bike Rental Station 2	(6, 2) Restaurant 2	(7, 2) Small Plaza 2	(8, 2) Student Activity Center 2	(9, 2) Academic Building D	
(0, 1) Water Park	(1, 1) Bar 2	(2, 1) Toy Store	(3, 1) Stationery Store	(4, 1) Bank 1	(5, 1) Hotel 1	(6, 1) Student Activity Center 1	(7, 1) Academic Building A	(8, 1) Academic Building B	(9, 1) Academic Building C	
(0, 0) Park	(1, 0) Government Office	(2, 0) Optical Shop 1	(3, 0) Bar 1	(4, 0) Restaurant 1	(5, 0) Convention Center	(6, 0) Parking Lot 1	(7, 0) Parking Lot 2	(8, 0) Bike Rental Station 1	(9, 0) Small Plaza 1	

Figure 2: Task environment. Gongguan MRT area projected into a  $10 \times 10$  grid map for testing.

**Stage-wise fine-tuning.** Let  $f_{\theta^{(i)}}$  denote the model at stage  $i$ .

**(S1) Relation extraction.** At the first stage, the model identifies the egocentric relation and reference landmark from clean description  $A$ :

$$r_1 = f_{\theta^{(0)}}(A) \rightarrow (q, \ell_r). \quad (1)$$

**(S2) Coordinate mapping.** At the second stage, the model converts user and landmark coordinates into an absolute direction:

$$r_2 = f_{\theta^{(1)}}(u, p(\ell_r)) \rightarrow d_{\text{abs}} \in \{\text{N, E, S, W}\}. \quad (2)$$

**(S3) Orientation reasoning.** At the third stage, the model predicts the final orientation given the absolute direction and egocentric relation:

$$r_3 = f_{\theta^{(2)}}(d_{\text{abs}}, q) \rightarrow D^*. \quad (3)$$

**(S4) End-to-end integration.** Finally, transcripts and coordinates are serialized into a single textual multimodal input:

$$M = \text{``Audio: } ' + A' + \text{``} \mid \text{Coordinates: } ' + T \quad (4)$$

where “+” denotes string concatenation and  $T$  serializes  $(u, \{\ell_i, p(\ell_i)\})$  into tokens. The model then produces a full reasoning trace:

$$S = f_{\theta^{(3)}}(M), \quad (5)$$

which includes  $(r_1, r_2, r_3, D^*)$ .

Standard Prompting	Multimodal Chain-of-Thought Prompting
<b>Model Input</b> Audio: "I'm at Gongguan MRT Exit 2, and restaurant 5 is on my right." Coordinates: Gongguan_MRT_Exit_2(5,4), restaurant_5(6,4) Query: Which direction is the user facing?	<b>Model Input</b> Audio: "I'm at Gongguan MRT Exit 2, and restaurant 5 is on my right." Coordinates: Gongguan_MRT_Exit_2(5,4), restaurant_5(6,4) Query: Which direction is the user facing? Please reason step by step:
<b>Model Output</b> Final answer: East X	<b>Model Output</b> <b>Step 1: Extract spatial relations</b> Spatial relation = right Reference landmark = restaurant 5 <b>Step 2: Calculate absolute directions</b> Reference landmark = restaurant 5 Direction vector from Gongguan MRT Exit 2 to restaurant 5: (6,4)-(5,4) = (1,0) Direction = East <b>Step 3: Infer user orientation</b> Restaurant 5 is to the East and the user describes it as "on my right". Spatial mapping rules indicate right = East when the user is facing North. Therefore, the user is facing North. Final answer: North ✓

Figure 3: Comparison of standard prompting and MCoT. Standard prompting fails under ambiguous egocentric descriptions, whereas MCoT uses structured steps for better accuracy and interpretability.

**Objective.** Training supervision is applied to intermediate reasoning steps and final predictions:

$$\mathcal{L} = - \sum_{t=1}^{|S|} \log P(s_t | s_{<t}, M). \tag{6}$$

### 3.4 MULTIMODAL CHAIN-OF-THOUGHT

CoT prompting enhances reasoning by decomposing tasks into interpretable steps Jin et al. (2024); Wu et al. (2023). Extending this to multimodal settings improves interpretability and stability Xie et al. (2025); Cui et al. (2024). Building on these insights, our MCoT framework handles ASR noise, provides interpretable rationales, and localizes errors to specific reasoning stages. This makes MCoT well suited for conversational navigation in indoor or GPS-limited environments.

**Three-step reasoning.** Our MCoT decomposes orientation reasoning into three steps: (1) Relation extraction: identify relation  $q$  and landmark  $\ell_r$  from input transcript, (2) Coordinate mapping: compute  $d_{\text{abs}}$  from  $(u, p(\ell_r))$ , and (3) Orientation reasoning: infer  $D^*$  given  $(d_{\text{abs}}, q)$ .

**Comparison to standard prompting.** As shown in Figure 3, standard prompting directly maps input to output, often failing under ambiguous egocentric language. In contrast, MCoT introduces intermediate steps that align spatial relations with mapping rules, yielding more accurate and interpretable predictions. While our experiments focus on a  $10 \times 10$  grid for tractability, the framework can be extended to larger or continuous maps by adapting the coordinate mapping rules.

## 4 EXPERIMENTAL RESULTS

### 4.1 DATASET

We evaluate on the COR benchmark, which we construct for egocentric-to-alloentric orientation reasoning in conversational navigation. COR contains 4,600 instances, each consisting of an egocentric utterance in Traditional Chinese, structured landmark coordinates, and the alloentric orientation label. Each instance includes a step-by-step reasoning trace for supervising CoT generation. To simulate speech-driven conditions, we use TTS and ASR to generate noisy transcripts. All instances are automatically produced from grid-based mapping rules and verified by human annotators.

**Data splits.** The dataset is divided into 3,216 training, 688 validation, and 696 test examples, with training using clean text and ASR transcripts only during evaluation. To avoid distributional bias, the training set is balanced to cover all combinations of spatial relations. Each single orientation

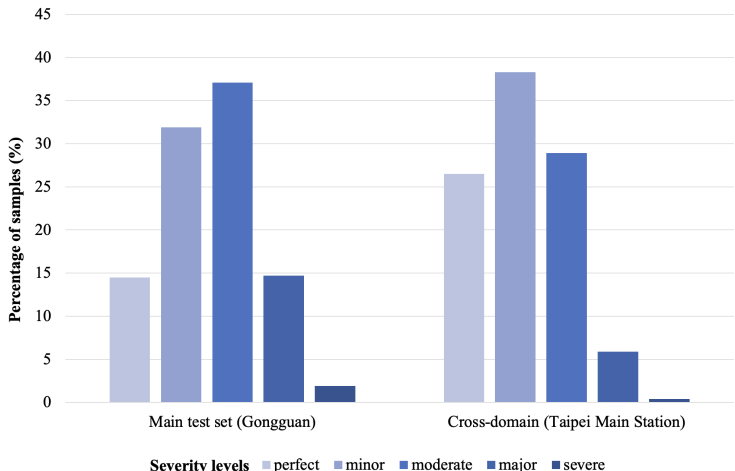


Figure 4: ASR error severity distribution in the two evaluation sets.

Table 2: Effectiveness on the egocentric spatial orientation task. Our MCoT with curriculum learning outperforms unimodal and non-structured baselines under both clean and ASR conditions.

ID	Method	Setting	Acc. (%)	Format err. (%)	Reasoning
B1	Zero-shot	Clean	25.0	5.2	–
B2	Few-shot (no CoT)	Clean	25.9	4.7	–
B3	Few-shot (with CoT)	Clean	21.1	39.2	0.534
B4	Fine-tuned (no CoT)	Clean	12.8	50.4	–
B5a	<b>MCoT + curriculum (ours)</b>	<b>Clean</b>	<b>100.0</b>	<b>0.0</b>	<b>1.000</b>
B5b	<b>MCoT + curriculum (ours)</b>	<b>ASR</b>	<b>98.1</b>	<b>0.0</b>	<b>1.000</b>

Reasoning quality reported only for methods with step-by-step reasoning; “–” means not applicable.

(front/back/left/right) contains 320 utterances. Double-orientation combinations (e.g., front+left, back+right, etc.), triple- and quadruple-orientation combinations, each contain about 280 utterances. This balanced design ensures fair coverage for training and evaluation. A subset of 400 examples introduces controlled linguistic variations (e.g., synonym substitutions, word-order changes), distributed across the splits. Test sets include multilingual elements common in Taiwan, with English landmark names appearing in 4.7% of main test samples and 46.5% of cross-domain samples, as well as occasional simplified Chinese variants from ASR outputs.

**Evaluation subsets.** Beyond the main test set (696 examples from Gongguan area), we build two additional evaluation sets for RQ3: (1) Cross-domain, 540 examples projected into a 10 × 10 grid from unseen Taipei Station area (Figure 5); and (2) Referential ambiguity, 200 cases with ambiguous references, disfluent or incomplete utterances, and semantically underspecified mentions.

**ASR error profile.** Figure 4 shows the ASR error severity distribution across the main and cross-domain test sets. Exact counts are provided in Appendix 5 (Table 7).

**Example.** “I am at Gongguan MRT Exit 2, and restaurant 5 is on my right” → Label: North.

#### 4.2 EXPERIMENTS ON ORIENTATION REASONING (RQ1: METHOD EFFECTIVENESS)

We report *orientation accuracy* as the primary metric, the proportion of predictions matching the ground-truth orientation. For step-by-step methods, we also measure *reasoning quality*, the match rate of intermediate steps (0–1, higher is better), and *format error rate*, the proportion of outputs violating the expected schema. Table 2 summarizes results across baselines and our method.

378  
379  
380  
381  
382  
383  
384  
385  
386  
387  
388  
389  
390  
391  
392  
393  
394  
395  
396  
397  
398  
399  
400  
401  
402  
403  
404  
405  
406  
407  
408  
409  
410  
411  
412  
413  
414  
415  
416  
417  
418  
419  
420  
421  
422  
423  
424  
425  
426  
427  
428  
429  
430  
431

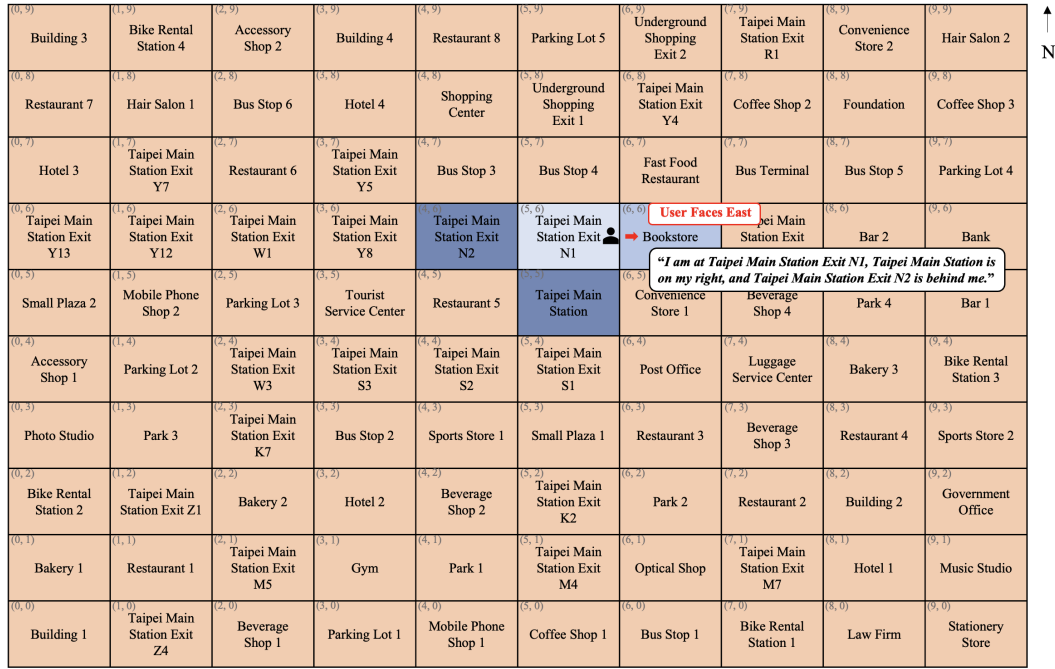


Figure 5: Cross-domain evaluation environment. Taipei Station area projected into a  $10 \times 10$  grid.

Direction Understanding Error	Relation Extraction Error	ASR Misrecognition Error
<b>Model Input</b> Audio: "I am at the cooperative store, small plaza 3 is in front of me, and theater is on my left." Coordinates: cooperative_store(8,8), small_plaza_3(7,8), theater(8,7) Query: Which direction is the user facing?	<b>Model Input</b> Audio: "I am at the market, bike rental 2 is in front of me, and bank 1 is on my right." Coordinates: market(4,2), bike_rental_2(5,2), bank_1(4,1) Query: Which direction is the user facing? Please reason step by step.	<b>Model Input</b> Audio: "I am at drink shop 2, and January is behind me." Coordinates: drink_shop_2(1,9), hospital(0,9) Query: Which direction is the user facing? Please reason step by step.
<b>Model Output</b> <b>Step 1: Extract spatial relations</b> Spatial relation 1 = front, Reference landmark = small plaza 3 Spatial relation 2 = left, Reference landmark = theater <b>Step 2: Calculate absolute directions</b> Reference landmark 1 = small plaza 3, Direction vector: (7,8) - (8,8) = (-1,0), Direction = West Reference landmark 2 = theater, Direction vector: (8,7) - (8,8) = (0,-1), Direction = South <b>Step 3: Infer user orientation</b> Small plaza 3 is to the West and the user describes it as "in front of me". Spatial mapping rules indicate front = West when the user is facing North. Theater is to the South and the user describes it as "on my left". Spatial mapping rules indicate left = South when the user is facing North. Therefore, the user is facing North. Final answer: North X	<b>Model Output</b> <b>Step 1: Extract spatial relations</b> Spatial relation 1 = front, Reference landmark = market Spatial relation 2 = right, Reference landmark = bank 1 <b>Step 2: Calculate absolute directions</b> Reference landmark 1 = market, Direction vector from bike rental 2 to market: (4,2) - (5,2) = (-1,0), Direction = West Reference landmark 2 = bank 1, Direction vector from market to bank 1: (4,1) - (4,2) = (0,-1), Direction = South <b>Step 3: Infer user orientation</b> The market is to the West and the user describes it as "in front of me". Spatial mapping rules indicate front = West when the user is facing West. Bank 1 is to the South and the user describes it as "on my right". Spatial mapping rules indicate right = South when the user is facing East. Therefore, the user is facing West. Final answer: West X	<b>Model Output</b> <b>Step 1: Extract spatial relations</b> Spatial relation = back, Reference landmark = January <b>Step 2: Calculate absolute directions</b> Reference landmark = January, Direction vector from drink shop 2 to January: (1,9) - (0,9) = (1,0), Direction = East <b>Step 3: Infer user orientation</b> January is to the East and the user describes it as "behind me". Spatial mapping rules indicate back = East when the user is facing West. Therefore, the user is facing West. Final answer: West X

Figure 6: Representative error cases falling into three categories: direction understanding errors, relation extraction errors, and ASR misrecognition errors.

**Findings and error analysis.** Our curriculum-trained MCoT improves accuracy by 74.1 percentage points over the strongest baseline (B2), achieving 100.0% on clean inputs and 98.1% with ASR transcripts. In the ASR setting, the model makes 13 residual errors (13/696, 1.9%). These errors mostly involve direction understanding (9 cases), with relation extraction mistakes (2 cases) and ASR misrecognition errors (3 cases; categories may overlap). Representative examples are shown in Figure 6. Full reasoning traces for all residual cases are provided in Appendix A.1.

#### 4.3 ABLATION STUDIES (RQ2: COMPONENT ANALYSIS)

We conduct ablation studies to assess the contributions of ASR transcripts, spatial coordinates, and structured CoT reasoning. Results are reported in Table 3.

**Findings.** Adding coordinates to ASR transcripts (A2→A3) improves accuracy by 10.2 percentage points and reduces format errors. Introducing structured CoT reasoning on top of multimodal

Table 3: Ablation study results. Spatial coordinates reduce ASR errors, while structured CoT provides the largest gains. The complete system achieves highest accuracy with no format errors.

ID	Configuration	Accuracy (%)	Format error (%)
A1	Clean text only (no coords)	25.0	0.7
A2	ASR text only (no coords)	16.2	35.8
A3	ASR text + coordinates (no CoT)	26.4	3.0
A4a	<b>Complete (clean text + coords + CoT)</b>	<b>100.0</b>	<b>0.0</b>
A4b	<b>Complete (ASR text + coords + CoT)</b>	<b>98.1</b>	<b>0.0</b>

Table 4: Robustness and generalization results. The model maintains high accuracy across linguistic variations, unseen domains, and referential ambiguity.

ID	Setting	Accuracy	Format error	Reasoning
R1	Linguistic variation	100%	0.0%	1.000
R2	Cross-domain (Taipei Station grid)	94.6% (511/540)	0.0%	1.000
R3	Referential ambiguity	99.5% (199/200)	0.0%	1.000

inputs (A3→A4b) contributes over 70 additional points. The complete system (A4a, A4b) eliminates all format errors, with 100% accuracy on clean inputs and 98.1% under ASR transcripts.

#### 4.4 ROBUSTNESS ANALYSIS (RQ3: ROBUSTNESS AND GENERALIZATION)

We evaluate robustness to linguistic variation, cross-domain generalization, and referential ambiguity. Results are reported in Table 4. Detailed experimental configurations for each robustness test are provided in Appendix A.4.

**Findings and discussion.** Across R1–R3, the model sustains perfect reasoning quality and zero format errors. These results suggest that the structured three-step MCoT effectively handles linguistic variation, transfers to unseen domains (Taipei Station), and manages referential ambiguity. The 29 residual errors in R2 are mostly direction understanding errors (22/29), followed by ASR misrecognition errors (6/29) and relation extraction errors (1/29).

## 5 CONCLUSION AND LIMITATIONS

**Conclusion.** This paper introduced a MCoT framework for egocentric-to-alloentric orientation reasoning in conversational navigation. The three-step reasoning process achieves 100% accuracy on clean text and 98.1% on ASR transcripts in Traditional Chinese. Experiments show that (1) curriculum-trained MCoT substantially outperforms unimodal and non-structured baselines, (2) spatial coordinates and structured reasoning help the model remain accurate under noisy transcripts, and (3) the approach demonstrates generalization across linguistic variation, unseen domains, and referential ambiguity. Overall, the results indicate that structured reasoning can enhance both accuracy and interpretability in conversational navigation.

**Limitations and future work.** Despite these contributions, this study has limitations. Evaluation is conducted within a 10×10 grid environment, and the framework is developed primarily for Traditional Chinese using synthesized speech data with grid-based rules, human verification, and controlled variation to approximate realistic conversational navigation. Future work should extend to larger continuous environments, incorporate real-time multilingual speech recognition under diverse noise, and explore modalities such as vision or motion cues for more realistic settings.

## REFERENCES

Peter Anderson, Qi Wu, Damien Teney, Jake Bruce, Mark Johnson, Niko Sünderhauf, Ian Reid, Stephen Gould, and Anton van den Hengel. Vision-and-language navigation: Interpreting

- 486 visually-grounded navigation instructions in real environments, 2018. URL <https://arxiv.org/abs/1711.07280>.
- 487
- 488
- 489 Howard Chen, Alane Suhr, Dipendra Misra, Noah Snaveley, and Yoav Artzi. Touchdown: Natural  
490 language navigation and spatial reasoning in visual street environments, 2020. URL <https://arxiv.org/abs/1811.12354>.
- 491
- 492 An-Chieh Cheng, Hongxu Yin, Yang Fu, Qiushan Guo, Ruihan Yang, Jan Kautz, Xiaolong Wang,  
493 and Sifei Liu. Spatialrgpt: Grounded spatial reasoning in vision language models, 2024. URL  
494 <https://arxiv.org/abs/2406.01584>.
- 495
- 496 Yingqian Cui, Pengfei He, Xianfeng Tang, Qi He, Chen Luo, Jiliang Tang, and Yue Xing. A the-  
497 oretical understanding of chain-of-thought: Coherent reasoning and error-aware demonstration,  
498 2024. URL <https://arxiv.org/abs/2410.16540>.
- 499
- 500 Harm de Vries, Kurt Shuster, Dhruv Batra, Devi Parikh, Jason Weston, and Douwe Kiela. Talk the  
501 walk: Navigating new york city through grounded dialogue, 2018. URL <https://arxiv.org/abs/1807.03367>.
- 502
- 503 Alexandre Défossez, Laurent Mazaré, Manu Orsini, Amélie Royer, Patrick Pérez, Hervé Jégou,  
504 Edouard Grave, and Neil Zeghidour. Moshi: a speech-text foundation model for real-time dia-  
505 logue, 2024. URL <https://arxiv.org/abs/2410.00037>.
- 506
- 507 Daniel Fried, Ronghang Hu, Volkan Cirik, Anna Rohrbach, Jacob Andreas, Louis-Philippe  
508 Morency, Taylor Berg-Kirkpatrick, Kate Saenko, Dan Klein, and Trevor Darrell. Speaker-follower  
509 models for vision-and-language navigation, 2018. URL <https://arxiv.org/abs/1806.02724>.
- 510
- 511 Chaoyou Fu, Haojia Lin, Xiong Wang, Yi-Fan Zhang, Yunhang Shen, Xiaoyu Liu, Haoyu Cao,  
512 Zuwei Long, Heting Gao, Ke Li, Long Ma, Xiawu Zheng, Rongrong Ji, Xing Sun, Caifeng Shan,  
513 and Ran He. Vita-1.5: Towards gpt-4o level real-time vision and speech interaction, 2025. URL  
514 <https://arxiv.org/abs/2501.01957>.
- 515
- 516 Sreyan Ghosh, Sonal Kumar, Ashish Seth, Chandra Kiran Reddy Evuru, Utkarsh Tyagi, S Sakshi,  
517 Oriol Nieto, Ramani Duraiswami, and Dinesh Manocha. Gama: A large audio-language model  
518 with advanced audio understanding and complex reasoning abilities, 2024. URL <https://arxiv.org/abs/2406.11768>.
- 519
- 520 Yuan Gong, Alexander H. Liu, Hongyin Luo, Leonid Karlinsky, and James Glass. Joint audio and  
521 speech understanding, 2023. URL <https://arxiv.org/abs/2309.14405>.
- 522
- 523 Kristen Grauman, Andrew Westbury, Eugene Byrne, Zachary Chavis, Antonino Furnari, Rohit Gird-  
524 har, Jackson Hamburger, Hao Jiang, Miao Liu, Xingyu Liu, Miguel Martin, Tushar Nagarajan,  
525 Ilija Radosavovic, Santhosh Kumar Ramakrishnan, Fiona Ryan, Jayant Sharma, Michael Wray,  
526 Mengmeng Xu, Eric Zhongcong Xu, Chen Zhao, Siddhant Bansal, Dhruv Batra, Vincent Car-  
527 tillier, Sean Crane, Tien Do, Morrie Doulaty, Akshay Erapalli, Christoph Feichtenhofer, Adriano  
528 Fragomeni, Qichen Fu, Abrham Gebreselasie, Cristina Gonzalez, James Hillis, Xuhua Huang,  
529 Yifei Huang, Wenqi Jia, Weslie Khoo, Jachym Kolar, Satwik Kottur, Anurag Kumar, Federico  
530 Landini, Chao Li, Yanghao Li, Zhenqiang Li, Karttikeya Mangalam, Raghava Modhugu, Jonathan  
531 Munro, Tullie Murrell, Takumi Nishiyasu, Will Price, Paola Ruiz Puentes, Meray Ramazanov,  
532 Leda Sari, Kiran Somasundaram, Audrey Southerland, Yusuke Sugano, Ruijie Tao, Minh Vo,  
533 Yuchen Wang, Xindi Wu, Takuma Yagi, Ziwei Zhao, Yunyi Zhu, Pablo Arbelaez, David Cran-  
534 dall, Dima Damen, Giovanni Maria Farinella, Christian Fuegen, Bernard Ghanem, Vamsi Krishna  
535 Ithapu, C. V. Jawahar, Hanbyul Joo, Kris Kitani, Haizhou Li, Richard Newcombe, Aude Oliva,  
536 Hyun Soo Park, James M. Rehg, Yoichi Sato, Jianbo Shi, Mike Zheng Shou, Antonio Torralba,  
537 Lorenzo Torresani, Mingfei Yan, and Jitendra Malik. Ego4d: Around the world in 3,000 hours of  
egocentric video, 2022. URL <https://arxiv.org/abs/2110.07058>.
- 538
- 539 Mingyu Jin, Qinkai Yu, Dong Shu, Haiyan Zhao, Wenye Hua, Yanda Meng, Yongfeng Zhang,  
and Mengnan Du. The impact of reasoning step length on large language models, 2024. URL  
<https://arxiv.org/abs/2401.04925>.

- 540 Yuka Kaniwa, Masaki Kuribayashi, Seita Kayukawa, Daisuke Sato, Hironobu Takagi, Chieko  
541 Asakawa, and Shigeo Morishima. Chitchatguide: Conversational interaction using large lan-  
542 guage models for assisting people with visual impairments to explore a shopping mall. *Proc.*  
543 *ACM Hum.-Comput. Interact.*, 8(MHCI), September 2024. doi: 10.1145/3676492. URL  
544 <https://doi.org/10.1145/3676492>.
- 545  
546 Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping, Rafael Valle, and Bryan Catanzaro. Audio  
547 flamingo: A novel audio language model with few-shot learning and dialogue abilities, 2024.  
548 URL <https://arxiv.org/abs/2402.01831>.
- 549  
550 Elad Levi and Ilan Kadar. Intelligent: A multi-agent framework for evaluating conversational ai  
551 systems, 2025. URL <https://arxiv.org/abs/2501.11067>.
- 552  
553 Ming-Yi Lin, Ou-Wen Lee, and Chih-Ying Lu. Embodied ai with large language models: A survey  
554 and new hri framework. In *2024 International Conference on Advanced Robotics and Mecha-*  
555 *tronics (ICARM)*, pp. 978–983, 2024. doi: 10.1109/ICARM62033.2024.10715872.
- 556  
557 Shuijing Liu, Aamir Hasan, Kaiwen Hong, Runxuan Wang, Peixin Chang, Zachary Mizrachi, Justin  
558 Lin, D. Livingston McPherson, Wendy A. Rogers, and Katherine Driggs-Campbell. Dragon: A  
559 dialogue-based robot for assistive navigation with visual language grounding. *IEEE Robotics and*  
560 *Automation Letters*, 9(4):3712–3719, 2024. doi: 10.1109/LRA.2024.3362591.
- 561  
562 Yuecheng Liu, Dafeng Chi, Shiguang Wu, Zhanguang Zhang, Yaochen Hu, Lingfeng Zhang,  
563 Yingxue Zhang, Shuang Wu, Tongtong Cao, Guowei Huang, Helong Huang, Guangjian Tian,  
564 Weichao Qiu, Xingyue Quan, Jianye Hao, and Yuzheng Zhuang. Spatialcot: Advancing spatial  
565 reasoning through coordinate alignment and chain-of-thought for embodied task planning, 2025.  
566 URL <https://arxiv.org/abs/2501.10074>.
- 567  
568 Ziyang Ma, Zhuo Chen, Yuping Wang, Eng Siong Chng, and Xie Chen. Audio-cot: Exploring chain-  
569 of-thought reasoning in large audio language model, 2025. URL <https://arxiv.org/abs/2501.07246>.
- 570  
571 Yao Mu, Qinglong Zhang, Mengkang Hu, Wenhai Wang, Mingyu Ding, Jun Jin, Bin Wang, Jifeng  
572 Dai, Yu Qiao, and Ping Luo. Embodiedgpt: Vision-language pre-training via embodied chain of  
573 thought, 2023. URL <https://arxiv.org/abs/2305.15021>.
- 574  
575 Ritawari Pareek, Divyansh Chauhan, Sonal Tuteja, and Kapil Madan. Enhancing campus navigation:  
576 A conversational ai agent for location assistance, 2024. URL <https://doi.org/10.1145/3675888.3676146>.
- 577  
578 Ben Prystawski, Michael Y. Li, and Noah D. Goodman. Why think step by step? reasoning emerges  
579 from the locality of experience, 2023. URL <https://arxiv.org/abs/2304.03843>.
- 580  
581 Zhixuan Shen, Haonan Luo, Kexun Chen, Fengmao Lv, and Tianrui Li. Enhancing multi-robot  
582 semantic navigation through multimodal chain-of-thought score collaboration, 2025. URL  
583 <https://arxiv.org/abs/2412.18292>.
- 584  
585 Smitha Sheshadri and Kotaro Hara. Conversational localization: Indoor human localization through  
586 intelligent conversation. *Proc. ACM Interact. Mob. Wearable Ubiquitous Technol.*, 7(4), January  
587 2024. doi: 10.1145/3631404. URL <https://doi.org/10.1145/3631404>.
- 588  
589 Jiatong Shi, Hye jin Shim, Jinchuan Tian, Siddhant Arora, Haibin Wu, Darius Petermann, Jia Qi Yip,  
590 You Zhang, Yuxun Tang, Wangyou Zhang, Dareen Safar Alharthi, Yichen Huang, Koichi Saito,  
591 Jionghao Han, Yiwen Zhao, Chris Donahue, and Shinji Watanabe. Versa: A versatile evaluation  
592 toolkit for speech, audio, and music, 2025. URL <https://arxiv.org/abs/2412.17667>.
- 593  
594 Qi Sun, Pengfei Hong, Tej Deep Pala, Vernon Toh, U-Xuan Tan, Deepanway Ghosal, and Soujanya  
595 Poria. Emma-x: An embodied multimodal action model with grounded chain of thought and  
596 look-ahead spatial reasoning, 2024. URL <https://arxiv.org/abs/2412.11974>.

- 594 Rahul Sundar, Shreyash Gadgil, Tankala Satya Sai, Sathi Sai Krishna Reddy, Gautam B, Ishita  
595 Mittal, Jyotsna Sree Guduguntla, and Shanmukesh Pujala. Innoguidegpt: Integrating conversa-  
596 tional interface and command interpretation for navigation robots. In *Proceedings of the Third*  
597 *International Conference on AI-ML Systems*, AIMLSystems '23, New York, NY, USA, 2024.  
598 Association for Computing Machinery. ISBN 9798400716492. doi: 10.1145/3639856.3639915.  
599 URL <https://doi.org/10.1145/3639856.3639915>.
- 600 Changli Tang, Wenyi Yu, Guangzhi Sun, Xianzhao Chen, Tian Tan, Wei Li, Lu Lu, Zejun Ma, and  
601 Chao Zhang. Salmonn: Towards generic hearing abilities for large language models, 2024. URL  
602 <https://arxiv.org/abs/2310.13289>.
- 603 Xiaojuan Tang, Zilong Zheng, Jiaqi Li, Fanxu Meng, Song-Chun Zhu, Yitao Liang, and Muhan  
604 Zhang. Large language models are in-context semantic reasoners rather than symbolic reasoners,  
605 2023. URL <https://arxiv.org/abs/2305.14825>.
- 607 Jesse Thomason, Michael Murray, Maya Cakmak, and Luke Zettlemoyer. Vision-and-dialog navi-  
608 gation, 2019. URL <https://arxiv.org/abs/1907.04957>.
- 609 Bin Wang, Xunlong Zou, Geyu Lin, Shuo Sun, Zhuohan Liu, Wenyu Zhang, Zhengyuan Liu, AiTi  
610 Aw, and Nancy F. Chen. Audiobench: A universal benchmark for audio large language models,  
611 2025. URL <https://arxiv.org/abs/2406.16020>.
- 613 Boshi Wang, Sewon Min, Xiang Deng, Jiaming Shen, You Wu, Luke Zettlemoyer, and Huan  
614 Sun. Towards understanding chain-of-thought prompting: An empirical study of what mat-  
615 ters. In Anna Rogers, Jordan Boyd-Graber, and Naoaki Okazaki (eds.), *Proceedings of the*  
616 *61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Pa-*  
617 *pers)*, Toronto, Canada, July 2023. Association for Computational Linguistics. URL <https://aclanthology.org/2023.acl-long.153/>.
- 619 Yifan Wu, Pengchuan Zhang, Wenhan Xiong, Barlas Oguz, James C. Gee, and Yixin Nie. The role  
620 of chain-of-thought in complex vision-language reasoning task, 2023. URL <https://arxiv.org/abs/2311.09193>.
- 622 Zhifei Xie and Changqiao Wu. Mini-omni: Language models can hear, talk while thinking in  
623 streaming, 2024. URL <https://arxiv.org/abs/2408.16725>.
- 625 Zhifei Xie, Mingbao Lin, Zihang Liu, Pengcheng Wu, Shuicheng Yan, and Chunyan Miao. Audio-  
626 reasoner: Improving reasoning capability in large audio language models, 2025. URL <https://arxiv.org/abs/2503.02318>.
- 628 Qian Yang, Jin Xu, Wenrui Liu, Yunfei Chu, Ziyue Jiang, Xiaohuan Zhou, Yichong Leng, Yuan-  
629 jun Lv, Zhou Zhao, Chang Zhou, and Jingren Zhou. Air-bench: Benchmarking large audio-  
630 language models via generative comprehension, 2024. URL <https://arxiv.org/abs/2402.07729>.
- 632 Chien yu Huang, Ke-Han Lu, Shih-Heng Wang, Chi-Yuan Hsiao, Chun-Yi Kuan, Haibin Wu,  
633 Siddhant Arora, Kai-Wei Chang, Jiatong Shi, Yifan Peng, Roshan Sharma, Shinji Watanabe,  
634 Bhiksha Ramakrishnan, Shady Shehata, and Hung yi Lee. Dynamic-superb: Towards a dy-  
635 namic, collaborative, and comprehensive instruction-tuning benchmark for speech, 2024. URL  
636 <https://arxiv.org/abs/2309.09510>.
- 637 Dong Zhang, Shimin Li, Xin Zhang, Jun Zhan, Pengyu Wang, Yaqian Zhou, and Xipeng Qiu.  
638 Speechgpt: Empowering large language models with intrinsic cross-modal conversational abil-  
639 ities, 2023. URL <https://arxiv.org/abs/2305.11000>.

641  
642  
643  
644  
645  
646  
647

## A APPENDIX

### A.1 ERROR TAXONOMY AND EXAMPLES

We categorize residual errors into three main types. Table 5 summarizes their definitions, and we provide one representative example for each type below.

Error type	Definition and example
Direction Understanding	The model correctly extracts landmarks and computes their coordinates, but misapplies the spatial mapping rules between relative orientation terms (front, back, left, right) and absolute directions (North, South, East, West).
Relation Extraction	The model fails at relation extraction, misidentifying or omitting a landmark. Even if coordinates are computed correctly, reasoning is built on faulty relations.
ASR Misrecognition	Errors caused by transcription mistakes in the ASR system, which corrupt the input before reasoning.

Table 5: Error taxonomy: three types of residual errors observed in our analysis.

Direction Understanding Error
<b>Model Input</b>
Audio: “我現在在台北車站S2出口，前面是餐廳5，左邊是台北車站S3出口” Coordinates: 台北車站S2出口(4,4), 餐廳5(4,5), 台北車站S3出口(3,4) 問題: 使用者面朝哪個方向? 請逐步推理:
<b>Model Output</b>
第一步: 提取空間關係 關係1: 空間關係=前面, 參考地標=餐廳5 關係2: 空間關係=左邊, 參考地標=台北車站S3出口 第二步: 計算絕對方向 參考地標1=餐廳5, 從台北車站S2出口到餐廳5方向向量: (4,5)-(4,4)=(0,1), 方向: 北方 參考地標2=台北車站S3出口, 從台北車站S2出口到台北車站S3出口方向向量: (3,4)-(4,4)=(-1,0), 方向: 西方 第三步: 推理朝向 餐廳5在北方且在使用者的左邊, 空間映射: 面朝南方時, 左邊=北方 台北車站S3出口在西方且在使用者的左邊, 空間映射: 面朝南方時, 左邊=西方 結論: 使用者面朝南方
Final answer: 南方 X

Figure 7: Direction understanding error with original Chinese model output.

### A.1.1 DIRECTION UNDERSTANDING ERROR

**Error Analysis.** The model correctly extracted relations and computed absolute directions but incorrectly applied spatial mapping rules in the final step, concluding the user faces South. Figure 7 shows the original output.

#### Input.

Audio: I am at Taipei Main Station Exit S2, restaurant 5 is in front of me, and Taipei Main Station Exit S3 is on my left | Coordinates: Taipei\_Main\_Station\_Exit\_S2(4,4), restaurant\_5(4,5), Taipei\_Main\_Station\_Exit\_S3(3,4)

#### Model Output.

##### 1. Step 1: Extract spatial relations

Spatial relation 1 = front, landmark = restaurant 5

Spatial relation 2 = left, landmark = Taipei Main Station Exit S3

##### 2. Step 2: Calculate absolute directions

Reference landmark 1 = restaurant 5, Direction vector from Taipei Main Station Exit S2 to restaurant 5: (4,5)-(4,4) = (0,1), Direction = North

Reference landmark 2 = Taipei Main Station Exit S3, Direction vector from Taipei Main Station Exit S2 to Taipei Main Station Exit S3: (3,4)-(4,4) = (-1,0), Direction = West

##### 3. Step 3: Infer user orientation

Restaurant 5 is to the North and the user describes it as “on my left”. Spatial mapping rules

702  
703  
704  
705  
706  
707  
708  
709  
710  
711  
712  
713  
714  
715  
716  
717  
718  
719  
720  
721  
722  
723  
724  
725  
726  
727  
728  
729  
730  
731  
732  
733  
734  
735  
736  
737  
738  
739  
740  
741  
742  
743  
744  
745  
746  
747  
748  
749  
750  
751  
752  
753  
754  
755

Relation Extraction Error	
<b>Model Input</b>	<p>Audio: “我現在在公車站2，前面是台北車站S3出口，左邊是台北車站推騎出口，右邊是運動用品店1”</p> <p>Coordinates: 公車站2(3,3), 台北車站S3出口(3,4), 台北車站K7出口(2,3), 運動用品店1(4,3)</p> <p>問題: 使用者面朝哪個方向? 請逐步推理:</p>
<b>Model Output</b>	<p>第一步: 提取空間關係</p> <p>關係1: 空間關係=左邊, 參考地標=台北車站S3出口</p> <p>關係2: 空間關係=右邊, 參考地標=運動用品店1</p> <p>第二步: 計算絕對方向</p> <p>參考地標1=台北車站S3出口, 從公車站2到台北車站S3出口方向向量: (3,4)-(3,3)=(0,1), 方向: 北方</p> <p>參考地標2=運動用品店1, 從公車站2到運動用品店1方向向量: (4,3)-(3,3)=(1,0), 方向: 東方</p> <p>第三步: 推理朝向</p> <p>台北車站S3出口在北方且在使用者的左邊, 空間映射: 面朝南方時, 左邊=北方</p> <p>運動用品店1在東方且在使用者的右邊, 空間映射: 面朝南方時, 右邊=東方</p> <p>結論: 使用者面朝南方</p> <p>Final answer: 南方 X</p>

Figure 8: Relation extraction error with original Chinese model output.

indicate left = North when the user is facing South.

Taipei Main Station Exit S3 is to the West and the user describes it as “on my left”. Spatial mapping rules indicate left = West when the user is facing South.

**Final Answer.** The user is facing South.

#### A.1.2 RELATION EXTRACTION ERROR

**Error Analysis.** The model incorrectly extracted spatial relations from the input, misidentifying which landmark corresponds to which direction. Figure 8 shows the original output.

##### Input.

Audio: I am at bus stop 2, Taipei Main Station Exit S3 is in front of me, Taipei Main Station Exit K7 is on my left, and sports store 1 is on my right | Coordinates: bus\_stop.2 (3, 3), Taipei\_Main\_Station\_Exit\_S3 (3, 4), Taipei\_Main\_Station\_Exit\_K7 (2, 3), sports\_store.1 (4, 3)

##### Model Output.

###### 1. Step 1: Extract spatial relations

Spatial relation 1 = left, landmark = Taipei Main Station Exit S3

Spatial relation 2 = right, landmark = sports store 1

###### 2. Step 2: Calculate absolute directions

Reference landmark 1 = Taipei Main Station Exit S3, Direction vector from bus stop 2 to Taipei Main Station Exit S3: (3,4)-(3,3) = (0,1), Direction = North

Reference landmark 2 = sports store 1, Direction vector from bus stop 2 to sports store 1: (4,3)-(3,3) = (1,0), Direction = East

###### 3. Step 3: Infer user orientation

Taipei Main Station Exit S3 is to the North and the user describes it as “on my left”. Spatial mapping rules indicate left = North when the user is facing South.

Sports store 1 is to the East and the user describes it as “on my right”. Spatial mapping rules indicate right = East when the user is facing South.

**Final Answer.** The user is facing South.

756  
757  
758  
759  
760  
761  
762  
763  
764  
765  
766  
767  
768  
769  
770  
771  
772  
773  
774  
775  
776  
777  
778  
779  
780  
781  
782  
783  
784  
785  
786  
787  
788  
789  
790  
791  
792  
793  
794  
795  
796  
797  
798  
799  
800  
801  
802  
803  
804  
805  
806  
807  
808  
809

ASR Misrecognition Error	
<b>Model Input</b>	Audio: “我現在在公園4，前面是引溜電4，後面是981，左邊是念包電3，右邊是982” Coordinates: 公園4(8,5), 飲料店4(7,5), 酒吧1(9,5), 麵包店3(8,4), 酒吧2(8,6) 問題: 使用者面朝哪個方向? 請逐步推理:
<b>Model Output</b>	<p>第一步: 提取空間關係</p> <p>關係1: 空間關係=前面, 參考地標=引溜電4 關係2: 空間關係=後面, 參考地標=引溜電4 關係3: 空間關係=左邊, 參考地標=念包電3 關係4: 空間關係=右邊, 參考地標=酒吧2</p> <p>第二步: 計算絕對方向</p> <p>參考地標1=引溜電4, 從公園4到引溜電4方向向量: (8,5)-(7,5)=(1,0), 方向: 東方 參考地標2=飲料店4, 從公園4到飲料店4方向向量: (7,5)-(8,5)=(-1,0), 方向: 西方 參考地標3=念包電3, 從公園4到念包電3方向向量: (8,4)-(8,5)=(0,-1), 方向: 南方 參考地標4=酒吧2, 從公園4到酒吧2方向向量: (8,6)-(8,5)=(0,1), 方向: 北方</p> <p>第三步: 推理朝向</p> <p>引溜電4在東方且在使用者的前面, 空間映射: 面朝東方時, 前面=東方 飲料店4在西方且在使用者的後面, 空間映射: 面朝東方時, 後面=西方 念包電3在南方且在使用者的左邊, 空間映射: 面朝東方時, 左邊=南方 酒吧2在北方且在使用者的右邊, 空間映射: 面朝東方時, 右邊=北方 結論: 使用者面朝東方</p> <p>Final answer: 東方 X</p>

Figure 9: ASR misrecognition error with original Chinese model output.

### A.1.3 ASR MISRECOGNITION ERROR

**Error Analysis.** ASR misrecognition caused the model to incorrectly map multiple spatial relations to the same landmark, leading to inconsistent reasoning. Figure 9 shows the original output.

#### Input.

Audio: I am at park 4, yin-liu-dian 4 is in front of me, 981 is behind me, nian-bao-dian 3 is on my left, and 982 is on my right | Coordinates: park\_4(8,5), drink.shop\_4(7,5), bar\_1(9,5), bakery\_3(8,4), bar\_2(8,6)

#### Model Output.

##### 1. Step 1: Extract spatial relations

Spatial relation 1 = front, reference landmark = yin-liu-dian 4  
Spatial relation 2 = back, reference landmark = yin-liu-dian 4  
Spatial relation 3 = left, reference landmark = nian-bao-dian 3  
Spatial relation 4 = right, reference landmark = bar 2

##### 2. Step 2: Calculate absolute directions

Reference landmark 1 = yin-liu-dian 4, Direction vector from park 4 to yin-liu-dian 4: (8,5)-(7,5) = (1,0), Direction = East  
Reference landmark 2 = drink shop 4, Direction vector from park 4 to drink shop 4: (7,5)-(8,5) = (-1,0), Direction = West  
Reference landmark 3 = nian-bao-dian 3, Direction vector from park 4 to nian-bao-dian 3: (8,4)-(8,5) = (0,-1), Direction = South  
Reference landmark 4 = bar 2, Direction vector from park 4 to bar 2: (8,6)-(8,5) = (0,1), Direction = North

##### 3. Step 3: Infer user orientation

Yin-liu-dian 4 is to the East and the user describes it as “in front of me”. Spatial mapping rules indicate front = East when the user is facing East.  
Drink shop 4 is to the West and the user describes it as “behind me”. Spatial mapping rules indicate back = West when the user is facing East.  
Nian-bao-dian 3 is to the South and the user describes it as “on my left”. Spatial mapping rules indicate left = South when the user is facing East.

Bar 2 is to the North and the user describes it as “on my right”. Spatial mapping rules indicate right = North when the user is facing East.

**Final Answer.** East

## A.2 IMPLEMENTATION DETAILS FOR THE EXPERIMENTS

**Model Architecture.** We use Taiwan-LLM-13B-v2.0-Chat as the base model with LoRA fine-tuning under 4-bit quantization (rank  $r = 16$ ,  $\alpha = 32$ , dropout rate 0.05).

**Training Parameters.** All experiments run for 5 epochs with batch size 1 and gradient accumulation steps of 32. The learning rate is set to  $5 \times 10^{-5}$ , with a maximum sequence length of 768 tokens. Training is conducted in bfloat16 mixed precision on NVIDIA A100 GPUs.

**Data Format.** Inputs follow the LLaMA chat template, combining ASR transcripts and spatial coordinates. The model is trained to output structured three-step reasoning traces aligned with the MCoT design.

**ASR Preprocessing.** Speech inputs were transcribed with Whisper-base (zh) using greedy decoding (temperature 0, no beam search) on 16 kHz mono audio. We report character error rate (CER) on the TTS→ASR loop as a proxy of noise, rather than Whisper’s internal probabilities. TTS was used solely to synthesize controlled variants for evaluation experiments.

Table 6: Training and ASR preprocessing details.

Model	Taiwan-LLM-13B-v2.0-Chat (LoRA, 4-bit)
Rank / $\alpha$ / Dropout	16 / 32 / 0.05
Optimizer	AdamW, lr = $5 \times 10^{-5}$
Batch size / Accumulation steps	1 / 32
Precision	bfloat16
Hardware	NVIDIA A100
ASR	Whisper-base (zh), greedy, 16 kHz mono
Metric	CER on TTS→ASR loop

## A.3 FULL PROMPTS

### A.3.1 TRAINING PROMPTS

#### Step 0: Spatial Relation Extraction

**Task.** Extract all spatial relations and reference landmarks from the following sentence.

**Input.** “I am at Gongguan MRT Exit 1, Restaurant 4 is in front of me, Gongguan MRT Exit 2 is behind me, Bus Stop 1 is on my left, and Bike Rental Station 3 is on my right.”

**Output.**

```
Spatial relation 1 = front, Reference landmark 1 =
Restaurant 4
Spatial relation 2 = behind, Reference landmark 2 = Gongguan
MRT Exit 2
Spatial relation 3 = left, Reference landmark 3 = Bus Stop 1
Spatial relation 4 = right, Reference landmark 4 = Bike
Rental Station 3
```

#### Step 1: Vector Direction Learning

**Input.**

```
From = (0,0), To = (1,0)
```

864 **Output.**

865

866 Direction vector:  $(1,0) - (0,0) = (1,0)$ , Direction = East

867

868 **Step 2: Orientation Inference**

869 **Task.** Infer user orientation based on the landmark absolute direction and spatial relation.

870

871 **Input.**

872 Landmark absolute direction = East

873

874 Spatial relation = behind

875

875 **Output.**

876

877 Analysis: The landmark is to the East, and the user  
878 describes it as ``behind me.`` Spatial mapping rules  
879 indicate behind = East when the user is facing West.  
Therefore, the user is facing West.

880

881 **Step 3: Complete Chain-of-Thought Reasoning**

882

882 **Input.**

883

884 Audio: I am at Academic Building A, Student Activity Center  
885 1 is in front of me, Academic Building B is behind me,  
886 Parking Lot 2 is on my left, and Small Plaza 2 is on my  
887 right | Coordinates: Academic Building A(7,1), Student  
888 Activity Center 1(6,1), Academic Building B(8,1), Parking  
889 Lot 2(7,0), Small Plaza 2(7,2)

889

890 **Output.**

891

892 1. **Step 1: Extract spatial relations**

893

894 Spatial relation 1 = front, Reference landmark 1 = Student Activity Center 1

895

896 Spatial relation 2 = behind, Reference landmark 2 = Academic Building B

897

898 Spatial relation 3 = left, Reference landmark 3 = Parking Lot 2

899

900 Spatial relation 4 = right, Reference landmark 4 = Small Plaza 2

901

902 2. **Step 2: Calculate absolute directions**

903

904 Reference landmark 1 = Student Activity Center 1, Direction vector:  $(6,1) - (7,1) = (-1,0)$ ,

905

906 Direction = West

907

908 Reference landmark 2 = Academic Building B, Direction vector:  $(8,1) - (7,1) = (1,0)$ ,

909

910 Direction = East

911

912 Reference landmark 3 = Parking Lot 2, Direction vector:  $(7,0) - (7,1) = (0,-1)$ , Direction =

913

914 South

915

916 Reference landmark 4 = Small Plaza 2, Direction vector:  $(7,2) - (7,1) = (0,1)$ , Direction =

917

North

918

919 3. **Step 3: Infer user orientation**

920

921 Student Activity Center 1 is to the West, and the user describes it as “in front of me.” Spatial  
922 mapping rules indicate front = West when the user is facing West.

923

924 Academic Building B is to the East, and the user describes it as “behind me.” Spatial  
925 mapping rules indicate behind = East when the user is facing West.

926

927 Parking Lot 2 is to the South, and the user describes it as “on my left.” Spatial mapping  
928 rules indicate left = South when the user is facing West.

929

930 Small Plaza 2 is to the North, and the user describes it as “on my right.” Spatial mapping  
931 rules indicate right = North when the user is facing West.

932

933 Therefore, the user is facing West.

934

935 A.3.2 BASELINE EVALUATION PROMPTS

936

937 **Training data format.** All training uses the LLaMA chat template with Traditional Chinese inputs  
and outputs, as Taiwan-LLM-13B-v2.0-Chat is optimized for Traditional Chinese. All training and

918 evaluation were conducted in Traditional Chinese. The prompts below are English translations for  
 919 clarity and reproducibility.

920 **B1: Zero-shot baseline**

921  
 922 Question: Audio: I am at Gongguan MRT Exit 3,  
 923 and Dormitory 2 is on my right | Coordinates:  
 924 Gongguan\_MRT\_Exit\_3(5,8), Dormitory\_2(6,8)  
 925 Which direction is the user facing? Please answer North,  
 926 South, East, or West.  
 927 Answer:

928 **B2: Few-shot prompting (no CoT)**

929  
 930 Instruction: Based on the audio description and coordinate  
 931 information, determine which direction the user is facing.

932  
 933 Example: Audio: I am at the gym, and the pharmacy is in  
 934 front of me | Coordinates: gym(4,6), pharmacy(4,7)  
 935 Answer: North

936  
 937 Example: Audio: I am at the park, and the water park is  
 938 behind me | Coordinates: park(0,0), water\_park(0,1)  
 939 Answer: South

940  
 941 Example: Audio: I am at the foundation, and the high  
 942 school is on my right | Coordinates: foundation(0,7),  
 943 high\_school(0,6)  
 944 Answer: East

945  
 946 Example: Audio: I am at the cooperative store,  
 947 and the theater is on my left | Coordinates:  
 948 cooperative\_store(8,8), theater(8,7)  
 949 Answer: West

950  
 951 Question: {user\_input}  
 952 Answer:

953 **B3: Few-shot prompting with CoT**

954  
 955 Instruction: Use three-step reasoning to determine the  
 956 user's facing direction given the audio description and  
 957 coordinates.

958 **Example 1**

959 Input: Audio: I am at the gym, and the pharmacy is in  
 960 front of me | Coordinates: gym(4,6), pharmacy(4,7)  
 961 Output:  
 962 Step 1: Extract spatial relations  
 963 Spatial relation = front  
 964 Reference landmark = pharmacy  
 965 Step 2: Calculate absolute directions  
 966 Direction vector from gym to pharmacy:  $(4,7) - (4,6) =$   
 967  $(0,1)$   
 968 Direction = North  
 969 Step 3: Infer user orientation  
 970 The pharmacy is to the North, and the user describes it as  
 971 ``in front of me.``  
 Spatial mapping rules indicate front = North when the user  
 is facing North.  
 Therefore, the user is facing North.

972 **Example 2**

973 Input: Audio: I am at the park, and the water park is

972           behind me | Coordinates: park(0,0), water\_park(0,1)  
 973           Output:  
 974           Step 1: Extract spatial relations  
 975           Spatial relation = behind  
 976           Reference landmark = water park  
 977           Step 2: Calculate absolute directions  
 978           Direction vector from park to water park:  $(0,1) - (0,0) =$   
 979            $(0,1)$   
 980           Direction = North  
 981           Step 3: Infer user orientation  
 982           The water park is to the North, and the user describes it as  
 983           ``behind me.``  
 984           Spatial mapping rules indicate behind = North when the user  
 985           is facing South.  
 986           Therefore, the user is facing South.  
 987  
 988           Now use the same three-step reasoning:  
 989           Input: {user\_input}  
 990           Output:

#### 989 **B4: Fine-tuned direct classification**

991           USER: {user\_input}  
 992           ASSISTANT:  
 993

### 994 A.4 SPATIAL ROBUSTNESS DETAILS

#### 996 A.4.1 LINGUISTIC VARIATION ROBUSTNESS (R1)

997 To evaluate robustness to natural linguistic variations, we constructed test sets in Traditional Chinese  
 998 with diverse expression patterns, while ensuring identical spatial semantics and orientation outputs.  
 999

##### 1000 **Variation types**

- 1001           • **Word order variations:** sentence inversion, argument permutation, and syntactic para-
- 1002           phrasing.
- 1003           • **Synonym substitutions:** spatial term substitution, position verb substitution, and landmark
- 1004           term substitution.
- 1005

#### 1006 A.4.2 REFERENTIAL AMBIGUITY ROBUSTNESS (R3)

1007 We test the model’s ability to handle ambiguous or underspecified references commonly encountered  
 1008 in natural conversational navigation.  
 1009

##### 1010 **Variation types**

- 1011           • **Referential ambiguity:** generic references (“this building”), and demonstrative pronouns
- 1012           (“that place”).
- 1013           • **Incomplete utterances:** disfluency (“I am at... um...”), uncertainty markers (“should be”),
- 1014           and hesitation patterns.
- 1015           • **Semantic underspecification:** vague location terms (“some building”) and imprecise ref-
- 1016           erences (“over there”).
- 1017
- 1018

##### 1019 **Example of R3 test cases**

1020           **Original:** Audio: I am at security office, and Dormitory  
 1021           6 is behind me | Coordinates: security\_office(7,3),  
 1022           dormitory\_6(7,4)  
 1023           **Referential ambiguity:** Audio: I am at this building, that  
 1024           dormitory is behind me | Coordinates: security\_office(7,3),  
 1025           dormitory\_6(7,4)

1026 **Incomplete utterance:** Audio: I am at... um... security  
 1027 office, Dormitory 6 should be behind | Coordinates:  
 1028 security\_office(3,5), dormitory\_6(3,4)

1029  
 1030 **Semantic underspecification:** Audio: I am at some place,  
 1031 that building over there is behind me | Coordinates:  
 1032 security\_office(3,5), dormitory\_6(3,4)

1033

## 1034 A.5 ASR ERROR SEVERITY STATISTICS

1035

1036 Table 7 reports the distribution of ASR error severity in both evaluation sets.

1037

1038

Table 7: ASR error severity distribution with exact counts.

1039

1040

1041

1042

1043

1044

1045

1046

1047

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

1061

1062

1063

1064

1065

1066

1067

1068

1069

1070

1071

1072

1073

1074

1075

1076

1077

1078

1079

Evaluation set	Severity	Count	Percent
Main test set (Gongguan)	perfect	101	14.5%
	minor	222	31.9%
	moderate	258	37.1%
	major	102	14.7%
	severe	13	1.9%
Cross-domain (Taipei Station)	perfect	143	26.5%
	minor	207	38.3%
	moderate	156	28.9%
	major	32	5.9%
	severe	2	0.4%