

OMHBench: Benchmarking Balanced and Grounded Omni-Modal Multi-Hop Reasoning

Anonymous ACL submission

Abstract

Multimodal Large Language Models (MLLMs) have increasingly supported omni-modal processing across text, vision, and speech. However, existing evaluation frameworks for such models suffer from critical limitations, including modality shortcuts and biased reasoning paths. To address these challenges, we propose **OMHBench**, a novel benchmark designed to rigorously evaluate omni-modal multi-hop reasoning. It consists of 6,144 questions with balanced reasoning paths that are jointly grounded across all three modalities. Extensive evaluation of 13 state-of-the-art models reveals that (1) a large performance gap exists between proprietary and open-source MLLMs and (2) even proprietary models exhibit high sensitivity to reasoning path variations, resulting in **asymmetric omni-modal grounding**. Notably, models struggle when processing the speech modality, underscoring the need for balanced, multi-hop evaluation of omni-modal intelligence.

1 Introduction

Human perception and understanding are inherently complex, often requiring the integration of textual, visual, and auditory information. Accordingly, the ability to process such heterogeneous inputs in tandem is fundamental to achieving human-level AI (Baltrušaitis et al., 2019). Relatedly, recent Multimodal Large Language Models (MLLMs) have evolved from initial bi-modal variants (e.g., text-vision and text-audio) to more comprehensive ones that jointly process text, vision, and audio, often referred to as **omni-modal** (Microsoft et al., 2025; Xu et al., 2025b; Gemini Team et al., 2025).¹

The development of these models has also driven the emergence of new evaluation schemes, which fall into two main directions: **Omni-Modal Understanding (OMU)** (Li et al., 2025b; Hong

¹**Omni-modal** refers to the text-vision-audio setting, while **multi-modal** is a general term for more than one modality.

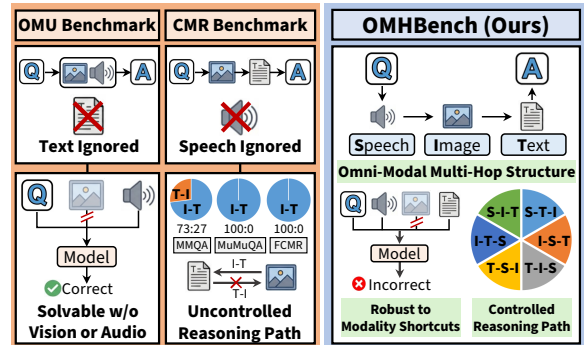


Figure 1: Omni-Modal Understanding (OMU) benchmarks lack textual context and suffer from modality shortcuts, while Cross-Modal Multi-Hop Reasoning (CMR) datasets exclude speech and exhibit imbalanced reasoning paths. OMHBench addresses these issues.

et al., 2025; Zhou et al., 2025; Chen et al., 2025; Nguyen et al., 2025), which emphasizes measuring a model’s ability to collectively handle text, vision, and audio; and **Cross-Modal Multi-Hop Reasoning (CMR)** (Talmor et al., 2021; Reddy et al., 2022; Kim et al., 2025; Foroutan et al., 2025; Jang et al., 2025), which focuses on its capability to perform multi-hop reasoning by composing information across modalities, typically in bi-modal settings. The key distinction between the two lies in whether the speech modality is incorporated and whether multi-hop reasoning is explicitly required.

In this work, we pose two crucial research questions regarding the current evaluation paradigms (see Figure 1): (1) If an OMU benchmark can be solved without leveraging all three modalities, can it truly be said to evaluate omni-modal understanding? (2) If a CMR benchmark is dominated by a single reasoning path, resulting in a heavily skewed composition distribution, can its results reliably reflect a model’s reasoning ability? We demonstrate through experiments that both suspected pitfalls are present in practice and substantially undermine the integrity of current omni-modal evaluations.

Benchmark	Text		Vision	Speech	CMR	Path Balance
	(Q)	(C)				
OmniBench	✓	✗	✓	✓	✗	—
WorldSense	✓	✗	✓	✓	✗	—
Daily-Omni	✓	✗	✓	✓	✗	—
OmniVideoBench	✓	✗	✓	✓	✗	—
UNO-Bench	✓	✗	✓	✓	✗	—
AV-SpeakerBench	✓	✗	✓	✓	✗	—
MMQA	✓	✓	✓	✗	✓	✗
MuMuQA	✓	✓	✓	✗	✓	✗
FCMR	✓	✓	✓	✗	✓	✗
ICT-QA	✓	✓	✓	✗	✓	✗
WikiMixQA	✓	✓	✓	✗	✓	✗
OMHBench(Ours)	✓	✓	✓	✓	✓	✓

Table 1: Comparison of OMU and CMR benchmarks by modality coverage and reasoning path balance. Text (Q) indicates question/option text, whereas Text (C) denotes separate contextual text. OMHBench uniquely supports all modalities with balanced multi-hop reasoning paths.

To systematically address these limitations, we propose **O(mnimodal)M(ulti)H(op)Bench**. By design, this novel benchmark departs from prior OMU and CMR datasets, as compared in Table 1. It requires omni-modal multi-hop reasoning over text, image, and speech, with each modality explicitly used at least once, eliminating shortcuts that allow models to solve tasks without access to a specific modality. Moreover, the ground-truth reasoning paths used as solutions are controlled with respect to modality order, enabling clearer identification of MLLMs’ strengths and weaknesses.²

Using OMHBench, we extensively test 13 proprietary and open-source MLLMs, uncovering several underexplored properties of omni-modal multi-hop reasoning. We find that (1) entity-attribute-based multi-hop structure effectively mitigates modality shortcut issues; (2) model performance rankings can vary considerably depending on the category of reasoning paths; (3) models exhibit strong asymmetry in omni-modal grounding, especially when transferring semantics into the speech modality.

In sum, this work (1) exposes fundamental limitations of existing OMU and CMR benchmarks; (2) presents OMHBench, an omni-modal benchmark with controlled and balanced multi-hop reasoning paths; and (3) reveals systematic weaknesses in omni-modal grounding, particularly in speech.

²The OMHBench dataset will be publicly released.

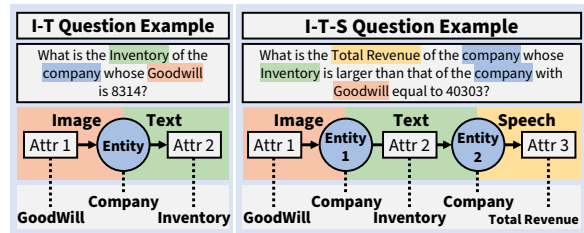


Figure 2: Illustration of the proposed task formulation with two example questions and their reasoning paths: I-T (left) and I-T-S (right). Attributes (e.g., Goodwill) are accessible only through specific modalities, while entities (e.g., Company) are shared across modalities.

2 Related Work

Multimodal Large Language Models (MLLMs)

MLLMs aim to extend the capabilities of LLMs beyond text by incorporating additional modalities such as vision and audio. Early efforts primarily focused on bi-modal settings, most notably text–vision (Alayrac et al., 2022; Liu et al., 2023a) or text–audio (Kong et al., 2024; Ghosh et al., 2024) integration, by attaching modality-specific encoders to off-the-shelf LLMs. Consequently, these models remain limited in their ability to jointly reason over more than two modalities, motivating recent efforts toward omni-modal architectures (Xu et al., 2025a; Yao et al., 2024; Microsoft et al., 2025; Ye et al., 2025) that natively support text, vision, and audio within a unified framework.

Omni-Modal Understanding (OMU) With the advent of omni-modal models, a few benchmarks have been proposed to test their capabilities. OmniBench (Li et al., 2025b) is an initial attempt to evaluate models under tri-modal inputs. WorldSense (Hong et al., 2025), Daily-Omni (Zhou et al., 2025), OmniVideoBench (Li et al., 2025a), and AV-SpeakerBench (Nguyen et al., 2025) focus on vision-audio understanding. UNO-Bench (Chen et al., 2025) further explores the relationship between uni-modal and omni-modal performance.

However, existing datasets emphasize visual and auditory signals, relegating text to questions or options, and often permit shortcuts that enable strong performance even without using all modalities.

Cross-Modal Multi-Hop Reasoning (CMR)

CMR evaluates a model’s ability to perform multi-hop reasoning by interleaving textual and visual evidence. Early benchmarks such as MMQA (Talmor et al., 2021) and MuMuQA (Reddy et al., 2022) require models to integrate information from text,

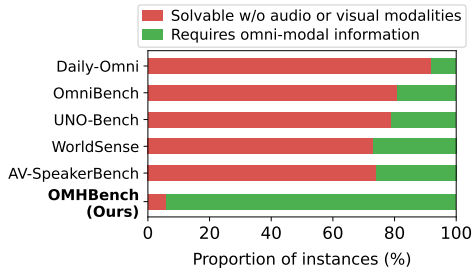


Figure 3: Fraction of instances in the OMU benchmarks that remain solvable *without* visual or auditory input. Across all datasets, nearly 70~80% of instances are susceptible, indicating that shortcuts allowing models to answer without using certain modalities are prevalent.

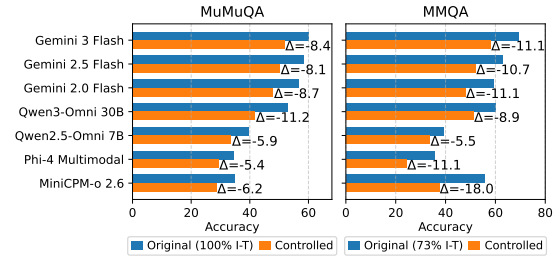


Figure 4: Performance comparison on the *original* MuMuQA and MMQA datasets vs. their *controlled* variants. When revised to ensure balanced reasoning paths, accuracy drops by up to 18%, raising concerns about the validity of previous evaluations on the original datasets.

tables, and images, while FCMR (Kim et al., 2025) extends this paradigm to the financial domain. ICTQA (Jang et al., 2025) and WikiMixQA (Foroutan et al., 2025) also explore multi-hop reasoning over structured sources, e.g., tables and charts.

Nevertheless, these benchmarks primarily focus on text-vision modalities and do not support audio-based reasoning. Moreover, the absence of explicit control over reasoning paths in these datasets often leads to heavily skewed reasoning pattern distributions, compromising the reliability of evaluation.

3 Preliminaries

3.1 Task Formulation: Omni-Modal Multi-Hop Reasoning

Here, we define the scope and formal definition of **omni-modal multi-hop reasoning** considered in this work. Multi-hop reasoning inherently operates over *entities* and their *attributes* as articulated in the context. From an omni-modal—more specifically, cross-modal—perspective, we consider scenarios in which entities are shared across modalities while attributes remain modality-specific, as illustrated in Figure 2. Answering such questions requires consulting modalities in a particular order, determined by the availability of the referenced attributes. We refer to this modality order as a **reasoning path**.

For instance, the question “*What is the inventory of the company whose goodwill is 8314?*” in Figure 2 necessitates integrating two attributes of the same entity, where *goodwill* is available in the image modality and *inventory* in text. In this case, the reasoning path is I-T. This formulation naturally generalizes to longer reasoning chains (e.g., I-T-S) through additional hops that identify subsequent entities. A reasoning problem is considered *omni-modal multi-hop* when the reasoning chain involves

all three modalities: image, text, and speech.

3.2 Shortcuts in OMU Benchmarks

To verify whether OMU benchmarks demand the exhaustive use of omni-modal input, we investigate five cases—OmniBench, WorldSense, Daily-Omni, UNO-Bench, and AV-SpeakerBench—by measuring the proportion of instances that remain solvable without visual or auditory input. This experiment is conducted using Gemini 3 Flash (Google, 2025b).

Figure 3 reports that nearly 70–80% of instances in OMU benchmarks can be answered without access to certain modalities. In other words, existing OMU benchmarks fail to genuinely assess the utilization of all three modalities due to insufficient structural constraints on cross-modal dependence.

3.3 Biases in CMR Benchmarks

CMR datasets, e.g., MuMuQA and MMQA, are characterized by skewed distributions in their representation of reasoning paths. As they consider only visual and textual inputs, there exist two possible orders: I-T and T-I. Nonetheless, these benchmarks are imbalanced: MuMuQA and FCMR contain only I-T instances, while MMQA is skewed toward I-T, with roughly twice as many I-T as T-I.³

To examine the effect of these biases, we conduct a controlled experiment by reversing the reasoning path direction (e.g., from I-T to T-I) for each question without changing its semantics. This yields variants with uniform reasoning path distributions while preserving the originals’ unique properties. Specifically, we create balanced versions of MuMuQA, MMQA and FCMR, each with an equal split (50:50) between I-T and T-I instances.⁴

³ICTQA (Jang et al., 2025) and WikiMixQA (Foroutan et al., 2025) are not publicly available, but their dataset construction does not consider reasoning path distributions.

⁴The detailed process is described in Appendix A.

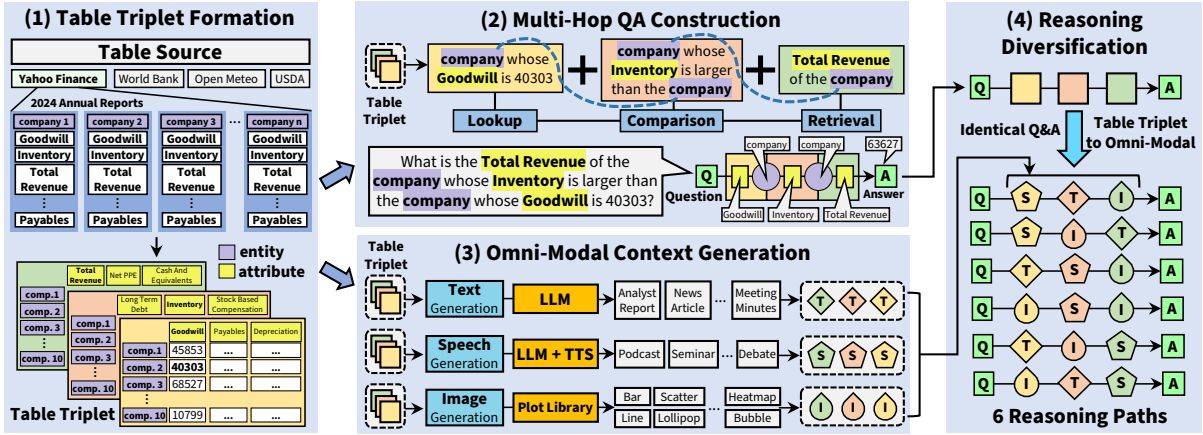


Figure 5: Overview of the OMHBench pipeline. (1) **Table Triplet Formation** constructs table triplets that share the same entities yet having separate attributes. (2) **Multi-Hop QA Construction** yields multi-hop QA pairs by utilizing content from table triplets. (3) **Omni-Modal Context Generation** converts each table into text, image, and speech modalities. (4) **Reasoning Diversification** guarantees multiple reasoning paths via modality permutation.

Figure 4 shows that performance exhibits a pronounced and consistent decline under controlled settings, in some cases exceeding a 10% drop in accuracy (see Figure 11 in the Appendix for FCMR). This indicates that balanced reasoning path distributions are essential for accurate MLLM evaluation; however, simple path reversal remains a partial solution limited to bi-modal settings, necessitating the development of a new benchmark.

4 Proposed Benchmark: OMHBench

Motivated by the largely isolated research streams on OMU and CMR and their respective limitations, we propose **OMHBench**, a new benchmark that bridges the two paradigms while resolving the issues discussed so far. It enables controlled evaluation of omni-modal multi-hop reasoning with three key desiderata: (1) it prevents shortcut-prone evaluation by enforcing multi-hop reasoning, (2) it jointly incorporates textual, visual, and speech modalities, and (3) it explicitly controls reasoning paths to support reliable and unbiased assessment.

As shown in Figure 5, OMHBench is constructed through four stages: (1) Table Triplet Formation, (2) Multi-Hop QA Construction, (3) Omni-Modal Context Generation, and (4) Reasoning Diversification. We refer readers to Appendix B for details.

4.1 Table Triplet Formation

OMHBench covers four domains—finance, economics, climate, and nutrition—where reasoning over text, images, and speech naturally occurs, with instances evenly distributed across domains. We

employ real-world table data from Yahoo Finance, World Bank, Open-Meteo, and USDA.

Given an original tabular source, we construct table triplets of three smaller tables that share the same set of *entities* but contain distinct *attributes*. Each table has a size of 10 entities \times 3 attributes and includes both relevant and distractor entities and attributes, requiring models to retrieve correct clues under information overload.

The intuition is that practical information is often organized in tabular form, yet its content can be realized across modalities—images for visual comparison, text for detailed description, and speech for public announcements. The table triplets serve as the core intermediate representation, which are used to construct question-answer pairs (§4.2) and corresponding omni-modal contexts (§4.3).

4.2 Multi-Hop QA Construction

From each table triplet, we formulate a question that requires three-hop reasoning to answer. In detail, we define eight reasoning operations—Lookup, Ranking, Comparison, Range, Proximity, Retrieval, Mean, and Summation—sample three of them, and sequentially apply each to the tables to derive questions.⁵ The first two reasoning steps focus on entity-level reasoning, applying operations like Lookup, Comparison, and Range to select, compare, or filter entities and pass either a single entity or several entities forward. The final step produces the answer

⁵In practice, not all combinations of reasoning operations are valid; their applicability depends on the entities and attributes involved. We therefore discard infeasible combinations and retain only those that apply to each case.

258	by either retrieving an attribute of a selected entity	(§4.3) to create multiple QA variants. They pre-	308
259	or aggregating attributes over a filtered entity set	serve the same question and answer, but differ in	309
260	using operations such as Mean or Summation.	how contextual evidence is organized. By permut-	310
261	Note that the proposed procedure is determinis-	ing modality assignments over the three tables, we	311
262	tic and rule-based, enabling scalable question gen-	obtain $3! = 6$ possible reasoning paths. Note that	312
263	eration given a sufficient number of tabular data	across these variants, the informational content re-	313
264	sources, without relying on costly external tools	remains unchanged; only the modality sequence re-	314
265	such as generative AI. Consequently, the QA con-	quired for inference (i.e., the reasoning path) varies.	315
266	struction process is efficient and fully automated.		
267	Finally, we partition the QA pairs into two cat-	4.5 Quality Control	316
268	egories based on their underlying reasoning oper-	Lastly, we apply four quality control methodologies	317
269	ation structures. OMHBench-Connect includes	to ensure dataset reliability and fair evaluation.	318
270	cases where intermediate results remain single enti-		
271	ties throughout the reasoning process, following a	Entity Anonymization Given prior findings that	319
272	fixed sequence of <i>Lookup-Comparison-Retrieval</i>	CMR datasets may allow shortcuts via parametric	320
273	operations. By contrast, OMHBench-Reasoning	knowledge (Kim et al., 2025), we anonymize entity	321
274	covers cases where intermediate results expand into	names with alphabetical codes (e.g., B, X), forcing	322
275	sets of entities, requiring aggregation operations.	models to rely solely on the provided context.	323
276			
277	4.3 Omni-Modal Context Generation	Consistency Checking We perform QA-based	324
278	At this stage, the three tabular sources from §4.1 are	consistency checks (Fabbri et al., 2021) by deriv-	325
279	transformed into contextual representations across	ing factoid questions from the original tables (§4.1)	326
280	three modalities—text, image, and speech. We	and validating answers using the converted context	327
281	explain this part using the financial domain as an	(§4.3). We also apply a test where an LLM recon-	328
282	example; the same process can be applied to others.	structs the original tables from the converted modal-	329
283	For the text modality, we define scenarios such	ities (§4.3) and compares them with the originals.	330
284	as analyst reports, news articles, and meeting mi-	Both checks achieve 100% consistency, confirming	331
285	utes, following prior work in financial text mining	the absence of factual loss or distortion.	332
286	(Kumar and Ravi, 2016; Pejić Bach et al., 2019;		
287	Gupta et al., 2020). Task-specific prompts are then	Question Rephrasing To enhance linguistic di-	333
288	used to guide LLMs in generating natural language	versity, we paraphrase questions using multiple	334
289	descriptions grounded in the underlying tables. For	LLMs: GPT-5.1, Grok-4, and Claude-Sonnet-4.5.	335
290	the image modality, we generate visualizations us-	Paraphrasing quality is evaluated with the Lexical	336
291	ing plotting libraries, following common practices	Deviation (LD) metric (Liu et al., 2022), where our	337
292	in financial data visualization (Ko et al., 2016; Ud-	dataset achieves higher LD scores than the widely	338
293	din et al., 2024; Christensen et al., 2024). For	used PAWS dataset (Zhang et al., 2019) (0.32 vs.	339
294	the speech modality, we adopt the taxonomy of	0.13), indicating greater lexical diversity.	340
295	financial speech scenarios proposed by Cao et al.		
296	(2025). The generated scripts are synthesized into	TTS Validation We evaluate TTS fidelity using	341
297	speech using Kokoro-82M TTS (Hexgrad, 2025).	ASR-based error rates (WER and CER) and speech	342
298	These are constructed as a multi-speaker dialogue	quality metrics (STOI and SI-SDR), following Ku-	343
299	in which each speaker describes different attributes	mar et al. (2023). The results demonstrate high	344
300	of the same entity, encouraging models to lever-	transcription accuracy and audio quality (WER:	345
301	age both semantic content and acoustic cues. To	0.03, CER: 0.02, STOI: 99.2, SI-SDR: 21.0).	346
302	enhance linguistic and stylistic diversity, we use		
303	three LLMs—GPT-5.1 (OpenAI, 2025), Grok-4	Final Dataset Statistics OMHBench comprises	347
304	(xAI, 2025), and Claude-Sonnet-4.5 (Anthropic,	6,144 instances evenly distributed across six rea-	348
305	2025)—for text and speech script generation.	soning paths. The benchmark is divided into two	349
306		subsets—OMHBench-Connect and OMHBench-	350
307		Reasoning—each with 3,072 instances, based on	351
		the required reasoning operations. Comprehensive	352
		dataset statistics, including diversity control across	353
		the three modalities, are reported in Table 4.	354

Model	Accuracy by Reasoning Path (%)						Avg. (Acc.)	PBS
	S-I-T	S-T-I	I-S-T	T-S-I	I-T-S	T-I-S		
Proprietary Models								
Gemini 3 Flash	97.5	98.4	75.4	75.0	60.2	63.5	78.3	32.2
Gemini 2.5 Pro	94.5	96.9	66.4	71.1	55.5	50.8	72.5	25.0
Gemini 2.5 Flash	82.0	85.9	50.8	54.7	26.6	21.9	53.6	4.7
Gemini 2.5 Flash-lite	49.2	60.9	38.3	35.2	5.5	4.7	32.3	0.0
Gemini 2.0 Flash	28.9	33.6	26.6	29.7	4.7	6.2	21.6	0.0
Gemini 2.0 Flash-lite	35.9	32.8	21.1	11.7	2.3	2.3	17.7	0.0
Open-Source Models								
Qwen3-Omni 30B	75.8	77.0	46.7	49.6	16.0	16.0	46.8	2.3
Phi-4 Multimodal	26.6	23.6	21.5	18.4	0.6	0.0	15.1	0.0
Qwen2.5-Omni 7B	22.7	20.9	19.3	20.5	2.0	1.8	14.5	0.0
Qwen2.5-Omni 3B	12.7	17.6	15.6	14.6	1.2	2.0	10.6	0.0
OmniVinci	14.8	8.6	14.8	7.0	0.8	0.6	7.8	0.0
MiniCPM-o 2.6	8.0	10.9	7.4	8.4	1.2	0.2	6.0	0.0
Omni-AutoThink	7.6	6.6	8.0	6.1	0.6	0.0	4.8	0.0

Table 2: Accuracies and Path Balance Scores (PBSs) across six reasoning paths in **OMHBench-Connect**. Avg denotes macro-averaged accuracy. PBSs (§5.2) measure robustness to reasoning path variations.

5 Experiments

5.1 Experimental Setup

We evaluate 13 MLLMs in total: both proprietary models—Gemini series (Google, 2025a; Comanici et al., 2025; Google, 2025b)—and open-source ones—Qwen3-Omni 30B (Xu et al., 2025c), Phi-4-Multimodal (Microsoft et al., 2025), Qwen2.5-Omni (Xu et al., 2025a), OmniVinci (Ye et al., 2025), MiniCPM-o 2.6 (Yao et al., 2024), and Omni-AutoThink (Yang et al., 2025).⁶ For models that support explicit reasoning modes, we enable this capability by setting a thinking budget of 8,192 tokens. All models are prompted using zero-shot chain-of-thought with a brief instruction and no fixed reasoning format. Model outputs are parsed into discrete answers and scored as correct or incorrect. We report *exact match* accuracies across all six reasoning paths, along with the macro-average. As Tan et al. (2024) shows that input modality order can affect model behavior—which we also observe (see Appendix C)—we randomize the arrangement of omni-modal contexts to prevent such biases.

5.2 Path Balance Score (PBS)

Beyond accuracy, we propose the **Path Balance Score (PBS)** as a novel metric to measure model robustness to variations in reasoning paths. In OMHBench, each question is instantiated across all permutations of the available modalities; with N modalities, this yields $N!$ reasoning paths sharing

⁶As of 2026-01-01, the Gemini series is the only proprietary model family that supports native omni-modal reasoning.

Model	Accuracy by Reasoning Path (%)						Avg. (Acc.)	PBS
	S-I-T	S-T-I	I-S-T	T-S-I	I-T-S	T-I-S		
Proprietary Models								
Gemini 3 Flash	55.9	58.8	49.8	49.6	40.0	42.6	49.4	8.6
Gemini 2.5 Pro	53.9	51.6	52.3	47.7	41.4	46.1	48.8	10.9
Gemini 2.5 Flash	32.0	30.5	17.2	24.2	10.9	10.9	21.0	0.0
Gemini 2.5 Flash-lite	18.8	21.1	15.6	8.6	0.0	0.0	10.7	0.0
Gemini 2.0 Flash	4.7	11.7	4.7	6.2	0.8	0.0	4.7	0.0
Gemini 2.0 Flash-lite	3.9	5.5	3.9	2.3	0.8	0.0	2.7	0.0
Open-Source Models								
Qwen3-Omni 30B	27.3	28.5	14.1	14.6	2.7	2.7	15.0	0.0
Phi-4 Multimodal	0.6	0.4	0.2	0.0	0.2	0.2	0.3	0.0
Qwen2.5-Omni 7B	0.4	1.0	1.0	0.6	0.2	1.2	0.7	0.0
Qwen2.5-Omni 3B	0.8	0.6	0.2	0.0	0.4	0.2	0.4	0.0
OmniVinci	0.6	0.2	0.2	0.4	0.0	0.0	0.2	0.0
MiniCPM-o 2.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Omni-AutoThink	0.4	0.2	0.4	0.2	0.0	0.0	0.2	0.0

Table 3: Accuracies and Path Balance Scores (PBSs) across six reasoning paths in **OMHBench-Reasoning**. Avg and PBSs are defined as in Table 2.

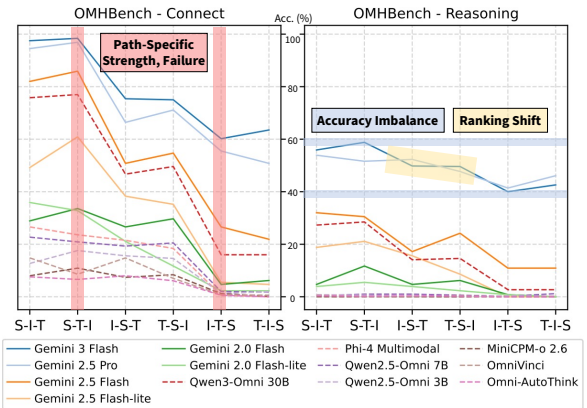


Figure 6: Visualization of core trends from Table 2 and Table 3, highlighting **path-specific strengths and failures**, **accuracy gaps**, and **ranking changes by reasoning paths**.

the same question. PBS evaluates whether a model can consistently answer all such paths.

Formally, let the dataset contain $|D|$ instances, forming $|G| = |D|/N!$ groups. For the i -th group, let $a_{i,j} \in \{0, 1\}$ denote whether the model correctly answers the j -th path. PBS is defined as:

$$\text{PBS} = \frac{1}{|G|} \sum_{i=1}^{|G|} \mathbb{I} \left(\sum_{j=1}^{N!} a_{i,j} = N! \right),$$

where $\mathbb{I}(\cdot)$ is the indicator function. Intuitively, PBS counts a group as correct only if the model can answer the question under *all* different reasoning paths, reflecting its robustness to path variations.

5.3 Main Results

Table 2 and Table 3 report the performance of LLMs on OMHBench-Connect and -Reasoning, respectively.⁷ Figure 6 provides a visual summary

⁷For models allowing reasoning mode, only thinking variants are reported; full results are shown in Tables 6 and 7.

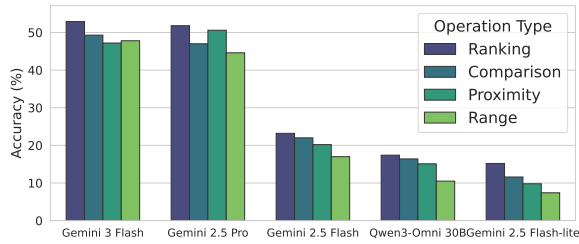


Figure 7: Accuracy by operation type on OMHBench-Reasoning, computed over instances whose reasoning chains include the corresponding operation, showing decreasing performance from Ranking to Range.

of these results and highlights several key trends.

Overall Trends For both OMHBench-Connect and OMHBench-Reasoning, proprietary models consistently outperform open-source models, with Qwen3-Omni 30B standing out as the strongest open-source model. Importantly, we observe substantial performance variations across reasoning paths, with Qwen3-Omni 30B achieving 77% on S-T-I versus 16% on I-T-S in OMHBench-Connect. OMHBench-Reasoning is more challenging than OMHBench-Connect, due to multi-entity intermediate states and numerical operations: even the best model reaches only 49.4% accuracy, while most open-source models perform near zero.

Difficulty by Reasoning Path We discover that MLLMs may answer the same question correctly or incorrectly depending on the composition of the omni-modal context. This behavior is clearly reflected by the low PBS scores observed for most models, indicating limited robustness to reasoning-path variations. In particular, *accuracy is strongly influenced by the position of the speech modality*: paths where speech appears earlier generally achieve higher performance, whereas those where speech appears later are substantially more challenging. We further analyze this effect in §6.2.

In addition, Figure 6 shows that model rankings can shift across different paths. In OMHBench-Reasoning, Gemini 2.5 Pro outperforms Gemini 3 Flash on I-S-T, but underperforms on T-S-I. This suggests that single-path evaluation (e.g., 100% I-T in MuMuQA) fails to characterize model behavior, underscoring the limitation of existing benchmarks.

Difficulty by Reasoning Operation Instances in OMHBench-Reasoning are crafted to require diverse operations (e.g., Range, Ranking) for their solution. To examine difficulty by operation type, we group instances accordingly and average accuracies

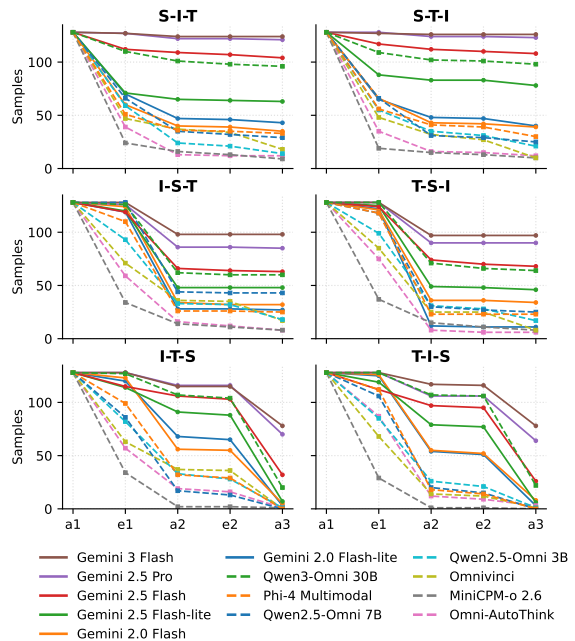


Figure 8: Step-by-step failure analysis for six reasoning paths. Each subplot reports the number of samples (out of 128 per reasoning path, 768 in total) successfully completing each reasoning stage— a_1 , e_1 , a_2 , e_2 , and a_3 —with success depending on all prior steps being correct. This view reveals key bottlenecks in OMHBench.

over the top five models. Figure 7 highlights that performance degrades from Ranking to Comparison, Proximity, and Range, implying that MLLMs handle ordinal or pairwise comparisons well, while operations involving numerical neighborhoods or interval constraints remain challenging.

6 Analysis

6.1 Modality Shortcut Validation

We verify that the modality shortcut issue observed in prior OMU benchmarks is no longer present in OMHBench. Following the protocol in §3.2, we measure the proportion of instances that remain solvable when one modality is removed. Figure 3 confirms that OMHBench exhibits almost no shortcut-prone cases, reflecting the effectiveness of its explicit multi-hop design. The few remaining solvable cases mainly stem from lookup-type questions, where correct answers can sometimes be obtained by chance through keyword retrieval.

6.2 Step-by-Step Failure Analysis

We conduct a step-by-step failure analysis to identify key bottlenecks in OMHBench, with results

	Connect						Reasoning					
	S-I-T	S-T-I	I-S-T	T-S-I	I-T-S	T-I-S	S-I-T	S-T-I	I-S-T	T-S-I	I-T-S	T-I-S
Economics	98.4	99.2	81.2	85.9	66.4	70.3	62.5	68.8	57.0	60.9	47.7	51.6
Finance	99.2	100.0	77.3	75.0	64.1	67.2	56.2	61.7	53.9	56.2	46.1	50.8
Climate	96.9	96.9	73.4	73.4	60.2	59.4	50.8	55.5	46.1	42.2	36.7	35.9
Nutrition	95.3	97.7	69.5	65.6	50.0	57.0	53.9	49.2	42.2	39.1	29.7	32.0

Figure 9: Domain-wise accuracies of Gemini 3 Flash on OMHBench-Connect and -Reasoning, with rows denoting domains and columns denoting reasoning paths. Performance gaps are amplified for challenging paths.

depicted in Figure 8.⁸ Leveraging the task’s step-wise structure, we categorize failures by the stage at which the model fails to identify the required entity (e) or attribute (a), using Gemini 3 Flash.

Weaker models frequently fail at early reasoning stages, especially in identifying e_1 or a_2 , regardless of the reasoning path, reflecting shortcomings in single-modal entity detection and cross-modal grounding. In contrast, stronger models generally succeed in identifying e_1 , but exhibit divergent performance at the a_2 stage depending on the reasoning path. By decomposing each three-hop path into two cross-modal grounding steps, we find that transition between text and image (T-I and I-T), as well as from speech to other modalities (S-I and S-T), is relatively robust. However, *reasoning that moves to the speech modality (I-S and T-S) proves particularly challenging*. We define this as **asymmetric omni-modal grounding**, underscoring inconsistencies in processing across modality orders.

6.3 Domain-Specific Analysis

Figure 9 presents the domain-specific performance of Gemini 3 Flash. Performance varies across domains, with a maximum gap of 21.8% between the economics and nutrition domains under T-S-I in OMHBench-Reasoning. This implies that even the best model lacks uniform domain generalization, performing better on common domains (e.g., economics) than on technical ones (e.g., nutrition).

6.4 Case Study

To complement quantitative analyses, we present case studies in Figure 10 with three key findings. (1) Even without explicit guidance in the prompt (i.e., zero-shot CoT), models consistently attempt to follow the intended reasoning path, confirming that OMHBench requires structured multi-hop reasoning. (2) Models sometimes behave as if the

⁸We perform this analysis on OMHBench-Connect for a controlled study of the required reasoning operations.

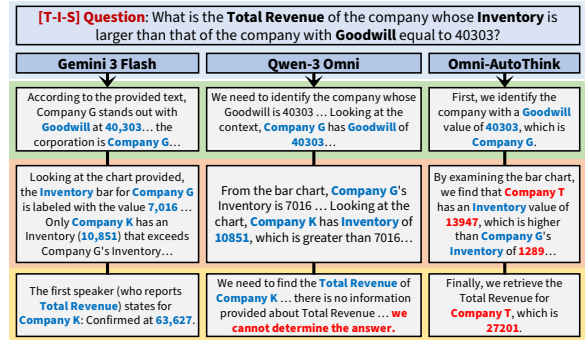


Figure 10: Three case studies show (1) adherence to the intended multi-hop reasoning path without explicit guidance, (2) neglect of the speech modality depending on its position, (3) error accumulation in weaker models.

speech modality were absent, reporting missing evidence despite speech information being provided; this behavior depends on the position of speech within the reasoning path (with similar input context lengths). (3) Weaker models exhibit error accumulation, where early mistakes propagate and result in cascading failures at later reasoning stages.

6.5 Prompting Alone is Not Enough

To examine whether the identified challenge can be alleviated by simply adopting advanced prompting techniques, we evaluate three multi-hop prompting strategies: Self-Ask (Press et al., 2023), Least-to-Most (Zhou et al., 2022), and Plan-and-Solve (Wang et al., 2023). As reported in Figure 13, none of these methods yields consistent improvements over the standard chain-of-thought baseline. This suggests that *asymmetric omni-modal grounding* is not primarily caused by insufficient prompt optimization, but instead reflects a fundamental limitation in transferring semantic representations across modalities, particularly into the speech modality.

7 Conclusion

We present OMHBench, a dataset for robust evaluation of omni-modal multi-hop reasoning. It addresses limitations of prior benchmarks by enforcing joint grounding across all three modalities and ensuring balanced proportions of distinct reasoning paths. Experiments on OMHBench provide new insights into how MLLMs perform multi-hop reasoning, revealing that they are sensitive to modality orders and struggle particularly with cross-modal grounding. In future work, we plan to explore training methods to improve the core multi-hop reasoning capabilities of omni-modal models.

530 Limitations

531 OMHBench adopts an entity-attribute based formu-
532 lation with fixed three-hop omni-modal reasoning
533 chains to enable controlled and balanced evalua-
534 tion of reasoning paths and to prevent modality
535 shortcuts. While this design facilitates precise anal-
536 ysis of modality interactions and fair comparison
537 across different reasoning paths, it primarily targets
538 reasoning scenarios that can be expressed through
539 explicit entity-attribute relations and fixed-depth
540 chains. Extending the benchmark to support more
541 diverse reasoning patterns is a promising direction
542 for future work.

543 References

544 Jean-Baptiste Alayrac, Jeff Donahue, Pauline Luc,
545 Antoine Miech, Iain Barr, Yana Hasson, Karel
546 Lenc, Arthur Mensch, Katherine Millican, Malcolm
547 Reynolds, and 1 others. 2022. Flamingo: a visual
548 language model for few-shot learning. *Advances in*
549 *neural information processing systems*, 35:23716–
550 23736.

551 Anthropic. 2025. Claude sonnet 4.5 sys-
552 tem card. [https://assets.anthropic.](https://assets.anthropic.com/m/12f214efcc2f457a/original/Claude-Sonnet-4-5-System-Card.pdf)
553 [com/m/12f214efcc2f457a/original/](https://assets.anthropic.com/m/12f214efcc2f457a/original/Claude-Sonnet-4-5-System-Card.pdf)
554 [Claude-Sonnet-4-5-System-Card.pdf](https://assets.anthropic.com/m/12f214efcc2f457a/original/Claude-Sonnet-4-5-System-Card.pdf). Ac-
555 cessed: 2025-12-30.

556 Tadas Baltrušaitis, Chaitanya Ahuja, and Louis-Philippe
557 Morency. 2019. Multimodal machine learning: A
558 survey and taxonomy. *IEEE Transactions on Pattern*
559 *Analysis and Machine Intelligence*, 41(2):423–443.

560 Yupeng Cao, Haohang Li, Yangyang Yu, Shashid-
561 har Reddy Javaji, Yueru He, Jimin Huang, Zining
562 Zhu, Qianqian Xie, Xiao-yang Liu, Koduvayur Sub-
563 balakshmi, and 1 others. 2025. Finaudio: A bench-
564 mark for audio large language models in financial
565 applications. *arXiv preprint arXiv:2503.20990*.

566 Chen Chen, ZeYang Hu, Fengjiao Chen, Liya Ma, Jiax-
567 ing Liu, Xiaoyu Li, Ziwen Wang, Xuezhi Cao, and
568 Xunliang Cai. 2025. Uno-bench: A unified bench-
569 mark for exploring the compositional law between
570 uni-modal and omni-modal in omni models. *arXiv*
571 *preprint arXiv:2510.18915*.

572 Xinyun Chen, Ryan A Chi, Xuezhi Wang, and Denny
573 Zhou. 2024. Premise order matters in reason-
574 ing with large language models. *arXiv preprint*
575 *arXiv:2402.08939*.

576 Theodore E Christensen, Karson E Fronk, Joshua A Lee,
577 and Karen K Nelson. 2024. Data visualization in
578 10-k filings. *Journal of Accounting and Economics*,
579 77(2-3):101631.

Gheorghe Comanici, Eric Bieber, Mike Schaeckermann,
Ice Pasupat, Noveen Sachdeva, Inderjit Dhillon, Mar-
cel Blistein, Ori Ram, Dan Zhang, Evan Rosen, and
1 others. 2025. Gemini 2.5: Pushing the frontier with
advanced reasoning, multimodality, long context, and
next generation agentic capabilities. *arXiv preprint*
arXiv:2507.06261.

Alexander R Fabbri, Chien-Sheng Wu, Wenhao Liu,
and Caiming Xiong. 2021. Qafacteval: Improved
qa-based factual consistency evaluation for summa-
rization. *arXiv preprint arXiv:2112.08542*.

Negar Foroutan, Angelika Romanou, Matin Ansari-
pour, Julian Martin Eisenschlos, Karl Aberer, and Rémi
Lebret. 2025. Wikimixqa: A multimodal benchmark
for question answering over tables and charts. *arXiv*
preprint arXiv:2506.15594.

Gemini Team, Rohan Anil, Sebastian Borgeaud, Jean-
Baptiste Alayrac, Jiahui Yu, Radu Soricut, and 1
others. 2025. Gemini: A family of highly capable
multimodal models. *Preprint*, arXiv:2312.11805.

Sreyan Ghosh, Sonal Kumar, Ashish Seth, Chandra Ki-
ran Reddy Evuru, Utkarsh Tyagi, S Sakshi, Oriol
Nieto, Ramani Duraiswami, and Dinesh Manocha.
2024. Gama: A large audio-language model with ad-
vanced audio understanding and complex reasoning
abilities. *arXiv preprint arXiv:2406.11768*.

Google. 2025a. Gemini 2.0 flash model card.
[https://modelcards.withgoogle.com/assets/](https://modelcards.withgoogle.com/assets/documents/gemini-2-flash.pdf)
documents/gemini-2-flash.pdf. Accessed:
2025-12-30.

Google. 2025b. Gemini 3 flash model card. [https://storage.googleapis.com/deepmind-media/](https://storage.googleapis.com/deepmind-media/Model-Cards/Gemini-3-Flash-Model-Card.pdf)
Model-Cards/Gemini-3-Flash-Model-Card.
pdf. Accessed: 2025-12-30.

Aaryan Gupta, Vinya Dengre, Hamza Abubakar
Kheruwala, and Manan Shah. 2020. Comprehen-
sive review of text-mining applications in finance.
Financial Innovation, 6(1):39.

Hexgrad. 2025. Kokoro-82m (revision d8b4fc7).

Jack Hong, Shilin Yan, Jiayin Cai, Xiaolong Jiang, Yao
Hu, and Weidi Xie. 2025. Worldsense: Evaluating
real-world omnimodal understanding for multimodal
llms. *arXiv preprint arXiv:2502.04326*.

Youngrok Jang, Hyesoo Kong, Gyeonghun Kim, Yejin
Lee, Jungkyu Choi, and Kyunghoon Bae. 2025. Ict-
qa: Question answering over multi-modal contexts
including image, chart, and text modalities. In *Pro-*
ceedings of the Computer Vision and Pattern Recog-
nition Conference, pages 138–148.

Seunghee Kim, Changhyeon Kim, and Taeuk Kim. 2025.
FCMR: robust evaluation of financial cross-modal
multi-hop reasoning. In *Proceedings of the 63rd*
Annual Meeting of the Association for Computational
Linguistics (Volume 1: Long Papers), pages 23352–
23380.

635	Sungahn Ko, Isaac Cho, Shehzad Afzal, Calvin Yau,	Le Thien Phuc Nguyen, Zhuoran Yu, Samuel Low Yu	689
636	Junghoon Chae, Abish Malik, Kaethe Beck, Yun	Hang, Subin An, Jeongik Lee, Yohan Ban, Se-	690
637	Jang, William Ribarsky, and David S Ebert. 2016.	ungEun Chung, Thanh-Huy Nguyen, JuWan Maeng,	691
638	A survey on visual analysis approaches for finan-	Soochahn Lee, and 1 others. 2025. See, hear, and un-	692
639	cial data. In <i>Computer Graphics Forum</i> , volume 35,	derstand: Benchmarking audiovisual human speech	693
640	pages 599–617. Wiley Online Library.	understanding in multimodal large language models.	694
		<i>arXiv preprint arXiv:2512.02231</i> .	695
641	Zhifeng Kong, Arushi Goel, Rohan Badlani, Wei Ping,		
642	Rafael Valle, and Bryan Catanzaro. 2024. Audio	OpenAI. 2025. Gpt-5.1 instant and gpt-5.1 thinking sys-	696
643	flamingo: A novel audio language model with few-	tem card addendum. https://cdn.openai.com/	697
644	shot learning and dialogue abilities. <i>arXiv preprint</i>	pdf/4173ec8d-1229-47db-96de-06d87147e07e/	698
645	<i>arXiv:2402.01831</i> .	5_1_system_card.pdf . Accessed: 2025-12-30.	699
646	Anurag Kumar, Ke Tan, Zhaoheng Ni, Pranay Manocha,		
647	Xiaohui Zhang, Ethan Henderson, and Buye Xu.	Mirjana Pejić Bach, Živko Krstić, Sanja Seljan, and	700
648	2023. Torchaudio-squim: Reference-less speech	Lejla Turulja. 2019. Text mining for big data analysis	701
649	quality and intelligibility measures in torchaudio.	in financial sector: A literature review. <i>Sustainability</i> ,	702
650	In <i>ICASSP 2023-2023 IEEE International Confer-</i>	11(5):1277.	703
651	<i>ence on Acoustics, Speech and Signal Processing</i>		
652	(<i>ICASSP</i>), pages 1–5. IEEE.	Ofir Press, Muru Zhang, Sewon Min, Ludwig Schmidt,	704
		Noah A Smith, and Mike Lewis. 2023. Measuring	705
653	B Shravan Kumar and Vadlamani Ravi. 2016. A survey	and narrowing the compositionality gap in language	706
654	of the applications of text mining in financial domain.	models. In <i>Findings of the Association for Computa-</i>	707
655	<i>Knowledge-Based Systems</i> , 114:128–147.	<i>tional Linguistics: EMNLP 2023</i> , pages 5687–5711.	708
656	Caorui Li, Yu Chen, Yiyan Ji, Jin Xu, Zhenyu Cui,		
657	Shihao Li, Yuanxing Zhang, Jiafu Tang, Zheng-	Revanth Gangi Reddy, Xilin Rui, Manling Li, Xudong	709
658	hao Song, Dingling Zhang, and 1 others. 2025a.	Lin, Haoyang Wen, Jaemin Cho, Lifu Huang, Mo-	710
659	Omnivideobench: Towards audio-visual understand-	hit Bansal, Avirup Sil, Shih-Fu Chang, Alexander	711
660	ing evaluation for omni mllms. <i>arXiv preprint</i>	Schwing, and Heng Ji. 2022. Mumuqa: Multime-	712
661	<i>arXiv:2510.10689</i> .	dia multi-hop news question answering via cross-	713
		media knowledge extraction and grounding . <i>Preprint</i> ,	714
		arXiv:2112.10728 .	715
662	Yizhi Li, Ge Zhang, Yinghao Ma, Ruibin Yuan, Kang		
663	Zhu, Hangyu Guo, Yiming Liang, Jiaheng Liu, Zekun	Alon Talmor, Ori Yoran, Amnon Catav, Dan Lahav,	716
664	Wang, Jian Yang, Siwei Wu, Xingwei Qu, Jinjie Shi,	Yizhong Wang, Akari Asai, Gabriel Ilharco, Han-	717
665	Xinyue Zhang, Zhenzhu Yang, Xiangzhou Wang,	naneh Hajishirzi, and Jonathan Berant. 2021. Mul-	718
666	Zhaoxiang Zhang, Zachary Liu, Emmanouil Bene-	timodalqa: Complex question answering over text,	719
667	tos, and 2 others. 2025b. Omnibench: Towards the	tables and images . <i>Preprint</i> , arXiv:2104.06039 .	720
668	future of universal omni-language models . <i>Preprint</i> ,		
669	arXiv:2409.15272 .	Zhijie Tan, Xu Chu, Weiping Li, and Tong Mo. 2024.	721
		Order matters: Exploring order sensitivity in mul-	722
		timodal large language models. <i>arXiv preprint</i>	723
		<i>arXiv:2410.16983</i> .	724
670	Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae		
671	Lee. 2023a. Visual instruction tuning . <i>Preprint</i> ,	Mohammed Majbah Uddin, Rahmat Ullah, and Mo-	725
672	arXiv:2304.08485 .	hammad Moniruzzaman. 2024. Data visualization in	726
673	Nelson F Liu, Kevin Lin, John Hewitt, Ashwin Paran-	annual reports—impacting investment decisions. <i>Inter-</i>	727
674	jape, Michele Bevilacqua, Fabio Petroni, and Percy	<i>national Journal for Multidisciplinary Research</i> ,	728
675	Liang. 2023b. Lost in the middle: How lan-	6(5).	729
676	guage models use long contexts. <i>arXiv preprint</i>		
677	<i>arXiv:2307.03172</i> .	Eric Wallace, Kai Xiao, Reimar Leike, Lilian Weng,	730
		Johannes Heidecke, and Alex Beutel. 2024. The	731
678	Timothy Liu and 1 others. 2022. Towards better char-	instruction hierarchy: Training llms to prioritize privi-	732
679	acterization of paraphrases. In <i>Proceedings of the</i>	leged instructions. <i>arXiv preprint arXiv:2404.13208</i> .	733
680	<i>60th Annual Meeting of the Association for Computa-</i>		
681	<i>tional Linguistics (Volume 1: Long Papers)</i> , pages	Lei Wang, Wanyu Xu, Yihuai Lan, Zhiqiang Hu, Yunshi	734
682	8592–8601.	Lan, Roy Ka-Wei Lee, and Ee-Peng Lim. 2023. Plan-	735
683	Microsoft, :, Abdelrahman Abouelenin, Atabak Ash-	and-solve prompting: Improving zero-shot chain-of-	736
684	faq, Adam Atkinson, Hany Awadalla, Nguyen Bach,	thought reasoning by large language models. <i>arXiv</i>	737
685	Jianmin Bao, and 1 others. 2025. Phi-4-mini tech-	<i>preprint arXiv:2305.04091</i> .	738
686	nical report: Compact yet powerful multimodal		
687	language models via mixture-of-loras . <i>Preprint</i> ,	xAI. 2025. Grok 4 model card. https://data.	739
688	arXiv:2503.01743 .	x.ai/2025-08-20-grok-4-model-card.pdf . Ac-	740
		cessed: 2025-12-30.	741

742 Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting
743 He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan,
744 Kai Dang, Bin Zhang, Xiong Wang, Yunfei Chu, and
745 Junyang Lin. 2025a. [Qwen2.5-omni technical report](#).
746 *Preprint*, arXiv:2503.20215.

747 Jin Xu, Zhifang Guo, Jinzheng He, Hangrui Hu, Ting
748 He, Shuai Bai, Keqin Chen, Jialin Wang, Yang Fan,
749 Kai Dang, and 1 others. 2025b. [Qwen2.5-omni](#)
750 [technical report](#). *arXiv preprint arXiv:2503.20215*.

751 Jin Xu, Zhifang Guo, Hangrui Hu, Yunfei Chu, Xiong
752 Wang, Jinzheng He, Yuxuan Wang, Xian Shi, Ting
753 He, Xinfa Zhu, Yuanjun Lv, Yongqi Wang, Dake
754 Guo, He Wang, Linhan Ma, Pei Zhang, Xinyu
755 Zhang, Hongkun Hao, Zishan Guo, and 19 others.
756 2025c. [Qwen3-omni technical report](#). *Preprint*,
757 arXiv:2509.17765.

758 Dongchao Yang, Songxiang Liu, Disong Wang,
759 Yuanyuan Wang, Guanglu Wan, and Helen Meng.
760 2025. [Omni-autothink: Adaptive multimodal rea-](#)
761 [soning via reinforcement learning](#). *arXiv preprint*
762 *arXiv:2512.03783*.

763 Yuan Yao, Tianyu Yu, Ao Zhang, Chongyi Wang, Junbo
764 Cui, Hongji Zhu, Tianchi Cai, Haoyu Li, Weilin
765 Zhao, Zhihui He, and 1 others. 2024. [Minicpm-v:](#)
766 [A gpt-4v level mllm on your phone](#). *arXiv preprint*
767 *arXiv:2408.01800*.

768 Hanrong Ye, Chao-Han Huck Yang, Arushi Goel, Wei
769 Huang, Ligeng Zhu, Yuanhang Su, Sean Lin, An-
770 Chieh Cheng, Zhen Wan, Jinchuan Tian, and 1 others.
771 2025. [Omnivinci: Enhancing architecture and data](#)
772 [for omni-modal understanding llm](#). *arXiv preprint*
773 *arXiv:2510.15870*.

774 Yuan Zhang, Jason Baldridge, and Luheng He. 2019.
775 [Paws: Paraphrase adversaries from word scrambling](#).
776 *arXiv preprint arXiv:1904.01130*.

777 Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei,
778 Nathan Scales, Xuezhi Wang, Dale Schuurmans,
779 Claire Cui, Olivier Bousquet, Quoc Le, and 1 oth-
780 ers. 2022. [Least-to-most prompting enables complex](#)
781 [reasoning in large language models](#). *arXiv preprint*
782 *arXiv:2205.10625*.

783 Ziwei Zhou, Rui Wang, and Zuxuan Wu. 2025. [Daily-](#)
784 [omni: Towards audio-visual reasoning with tem-](#)
785 [poral alignment across modalities](#). *arXiv preprint*
786 *arXiv:2505.17862*.

A Preliminary CMR Dataset Construction

Since MMQA and MuMuQA are formulated as short-answer question answering tasks, while FCMR follows a multiple-choice format, we handle these datasets separately during the preliminary CMR dataset construction stage.

We first preprocess MMQA to align its format with that of MuMuQA, thereby simplifying the overall pipeline design. Specifically, we treat the table input in MMQA instances as a text modality and extend the context column of each instance to include a textual representation of the given table. Next, we select instances from both MuMuQA and our preprocessed MMQA, that can be reformulated using our attribute-entity formulation. This step is necessary because certain instances in MMQA do not require cross-modal reasoning and thus cannot be reversed. We then prompt Gemini 2.0 Flash to decompose each question into its constituent entity and attributes, followed by question-generation step with the prompts provided in Figure 15 and 16. Upon manual inspection of the generated questions, we find that some instances explicitly include the entity or the answer within the question text; we discard such cases. Finally, we use an LLM-based validation step to further filter out questions with an incorrect reasoning order, using the prompt in Figure 17. We pair up instances of two direction—I-T, and T-I—using the generated questions and their original counterparts. Based on these pairs, we construct two sub-datasets: Original, which consists of instances from the original dataset, and Controlled, which includes both directional pairs.

For FCMR, the dataset construction process differs from MMQA and MuMuQA. Since FCMR is designed with a template-based structure and explicit multi-hop reasoning paths, we directly apply Gemini 3 Flash to generate controlled instances over the answer options. Using this procedure, we construct both Original and Controlled versions of the FCMR dataset. The detailed generation prompt is provided in Figure 18.

B Details of OMHBench Construction

The following describes the detailed design choices involved in constructing the benchmark. Example samples from the dataset are shown in Figure 22, Figure 23 and Figure 24.

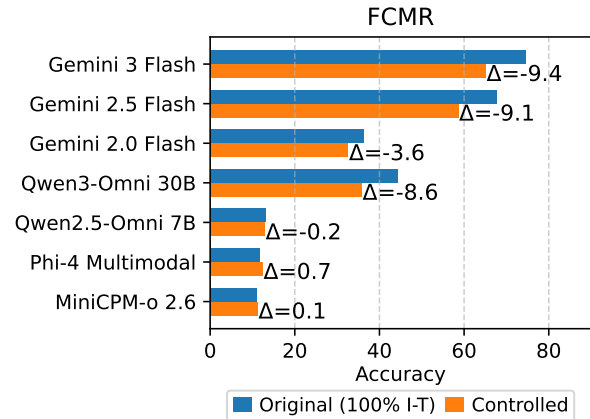


Figure 11: Performance comparison on the *original* FCMR datasets vs. their *controlled* variants. The random baseline is 12.5%, and due to the difficulty of the dataset, most open-source models perform close to this baseline under both settings.

B.1 Data Source

OMHBench comprises four domains: finance, economics, climate, and nutrition. The data for each domain were obtained from Yahoo Finance⁹, World Bank¹⁰, Open-Meteo¹¹, and U.S. Department of Agriculture (USDA)¹², respectively. All values are standardized to use consistent units across entities to ensure comparability. To avoid ambiguous answers, attributes with duplicate values across samples are excluded from the dataset.

Finance Domain The finance domain is constructed using annual financial statements from 23 publicly listed companies. The selected companies and their ticker symbols are: MSFT, NVDA, AVGO, QCOM, TXN, IBM, ADI, MU, KLAC, LLY, MRK, ABBV, TMO, PFE, GILD, BMY, ZTS, PG, KO, PEP, DIS, MDLZ, and HON. For each company, we extract 15 standardized financial attributes from their 2024 annual reports obtained from Yahoo Finance. These indicators span various aspects of corporate performance, including Total Revenue, Cost Of Revenue, Selling General And Administration, Cash And Equivalents, Receivables, Inventory, Other Current Assets, Other Non Current Assets, Net PPE, Goodwill, Payables, Long Term Debt, Other Non Current Liabilities, Depreciation, and Stock Based Compensation.

⁹<https://finance.yahoo.com/>

¹⁰<https://data.worldbank.org/>

¹¹<https://open-meteo.com/>

¹²<https://fdc.nal.usda.gov/>

Statistics	Number
Dataset	
OMHBench-Connect	3,072
OMHBench-Reasoning	3,072
Total samples	6,144
Reasoning path	
I-S-T	1,024
I-T-S	1,024
S-I-T	1,024
S-T-I	1,024
T-I-S	1,024
T-S-I	1,024
Domain	
Finance	1,536
Economics	1,536
Climate	1,536
Nutrition	1,536
Operation types	
Operation combinations	33
Image type	
Image color	20
Image font	20
Plot library	2
Speech type	
Speech voice	27
Text type	
Generation LLM variants	3

Table 4: Key statistics of OMHBench.

Economics Domain The economics domain is constructed using country-level economic attributes obtained from the World Bank World. All values correspond to annual data for the year 2024. The dataset includes 18 countries: ARG, AUS, BRA, CHE, DEU, EGY, ESP, FRA, GBR, IDN, IND, ITA, MEX, NLD, NOR, SAU, SWE, and ZAF. For each country, we collect 18 standardized economic indicators representing major components of national economic activity. The selected indicators include Personal remittances, paid, Personal remittances, received, Total reserves, excluding gold, Final consumption expenditure, Gross fixed capital formation, Gross capital formation, Agriculture, forestry, and fishing, value added, Manufacturing, value added, Industry, including construction, value added, Services, value added, Gross Value Added (GVA) at basic prices, GNI, Gross savings, Taxes

less subsidies on products, Merchandise imports, Commercial service imports, Merchandise exports, and Commercial service exports.

Climate Domain The climate domain is constructed using monthly meteorological data collected from major cities around the world. All climate data correspond to the year 2024 and are obtained from Open-Meteo. The dataset includes 20 representative cities: Karachi, Addis Ababa, Cairo, Nairobi, Los Angeles, Tokyo, Ho Chi Minh City, Ulaanbaatar, Chicago, Singapore, Toronto, Shanghai, Manila, London, Lagos, Chengdu, Beijing, Dubai, Rome, and Mumbai. These cities were selected to cover diverse geographic regions and climate conditions. For each city, we collect 12 climate attributes corresponding to the maximum wind speed for each month from January to December.

Nutrition Domain The nutrition domain is constructed using food composition data obtained from the U.S. Department of Agriculture (USDA). The dataset includes 24 food items: Potatoes, mashed, dehydrated, granules without milk, dry form; Sorghum flour, whole-grain; Wheat, KAMUT khorasan, uncooked; PAPA JOHN’S 14" The Works Pizza, Original Crust; Lasagna with meat sauce, frozen, prepared; Potatoes, mashed, home-prepared, whole milk added; Frankfurter, turkey; Seeds, sesame seed kernels, dried (decorated); Pork, cured, ham – water added, slice, bone-in, separable lean and fat, unheated; Nuts, cashew nuts, raw; T.G.I. FRIDAY’S, chicken fingers, from kids’ menu; Pork, cured, ham with natural juices, shank, bone-in, separable lean only, unheated; Broccoli, cooked, boiled, drained, with salt; Pork, cured, ham and water product, shank, bone-in, unheated, separable lean only; Bologna, beef; Teff, uncooked; Nuts, pecans; HOT POCKETS Ham ’N Cheese Stuffed Sandwich, frozen; Pork, cured, ham – water added, slice, boneless, separable lean only, heated, pan-broil; Pork sausage, link/patty, fully cooked, microwaved; DENNY’S, chicken strips; Pork, cured, ham and water product, rump, bone-in, separable lean only, heated, roasted; Pork, cured, ham with natural juices, spiral slice, boneless, separable lean only, unheated; Kielbasa, fully cooked, unheated; For each food item, 19 nutritional attributes are collected. These include Ash, Protein, Lysine, Methionine, Isoleucine, Leucine, Valine, Phenylalanine, Threonine, Histidine, Arginine, Tyrosine, Alanine, Glycine, Serine, Proline,

Operation sequence	# Instances
<i>Connect</i>	
Lookup–Comparison–Retrieval	3,072
<i>Reasoning</i>	
Ranking–Ranking–Mean	96
Ranking–Ranking–Summation	96
Ranking–Range–Mean	96
Ranking–Range–Summation	96
Ranking–Comparison–Mean	96
Ranking–Comparison–Summation	96
Ranking–Proximity–Mean	96
Ranking–Proximity–Summation	96
Range–Ranking–Mean	96
Range–Ranking–Summation	96
Range–Range–Mean	96
Range–Range–Summation	96
Range–Comparison–Mean	96
Range–Comparison–Summation	96
Range–Proximity–Mean	96
Range–Proximity–Summation	96
Comparison–Ranking–Mean	96
Comparison–Ranking–Summation	96
Comparison–Range–Mean	96
Comparison–Range–Summation	96
Comparison–Comparison–Mean	96
Comparison–Comparison–Summation	96
Comparison–Proximity–Mean	96
Comparison–Proximity–Summation	96
Proximity–Ranking–Mean	96
Proximity–Ranking–Summation	96
Proximity–Range–Mean	96
Proximity–Range–Summation	96
Proximity–Comparison–Mean	96
Proximity–Comparison–Summation	96
Proximity–Proximity–Mean	96
Proximity–Proximity–Summation	96

Table 5: Operation sequences used to construct OMHBench-Connect and OMHBench-Reasoning.

931 Tryptophan, Cystine, and Glucose.

932 B.2 Table Triplet Formation

933 In the first stage of the dataset generation frame-
934 work, namely the Table Triplet Formation stage,
935 we construct table triplets consisting of three tables
936 that share the same set of entities but contain differ-
937 ent attributes. Each table contains 10 entities and 3
938 attributes, resulting in a table size of 10×3 .

939 When forming each table triplet, we ensure that
940 the selected attributes are distinct across the three
941 tables while referring to the same set of entities. In
942 addition, to facilitate subsequent conversion into
943 chart-based representations, we constrain the value
944 ranges of the attributes such that the ratio between
945 the maximum and minimum values does not ex-
946 ceed 30. This restriction prevents extreme scale
947 differences across attributes and ensures stable vi-
948 sualization and comparison across tables.

B.3 Multi-Hop QA Construction

949 OMHBench constructs each question as a three-
950 hop reasoning chain over entities and their at-
951 tributes. The first two hops select or filter sets
952 of entities, while the final hop produces a scalar
953 numerical answer. Each hop applies a specific op-
954 eration to a designated attribute, and all operations
955 are deterministic and rule-based, ensuring full con-
956 trol over the reasoning path. Below, we describe
957 the concrete mechanics of each operation in detail.
958 The valid sequences of operations used to construct
959 multi-hop questions are summarized in Table 5.
960

961 **Lookup** Lookup anchors the reasoning chain to
962 a single entity. From a base table containing ten
963 entities, one entity is randomly selected, and the
964 value of a specified attribute is retrieved. A strict
965 uniqueness constraint is enforced: the selected at-
966 tribute value must occur exactly once in the table.
967 If the same value appears for multiple entities, the
968 instance is discarded. The output of this operation
969 is a uniquely identified entity and its corresponding
970 attribute value.

971 **Comparison** Comparison filters an entity set us-
972 ing a strict inequality condition on a specified at-
973 tribute. Entities are sorted by the attribute, and a
974 threshold value is constructed such that exactly
975 a predefined number of entities satisfy either a
976 “larger than” or “smaller than” condition. Instances
977 in which ties occur at the decision boundary are
978 discarded. This operation outputs a reduced entity
979 set of fixed size.

980 **Ranking** Ranking selects entities based on their
981 ordinal position under a given attribute. Entities are
982 sorted in ascending or descending order, and a fixed
983 number of top or bottom entities are selected. No
984 explicit threshold values are involved. The output
985 is a subset of entities.

986 **Range** Range selects entities whose attribute val-
987 ues fall within a contiguous interval. Entities are
988 first sorted by the target attribute, and a consecu-
989 tive segment is selected. The interval boundaries
990 are defined to ensure that the selected entities are
991 uniquely determined. The output is a subset of
992 entities.

993 **Proximity** Proximity selects entities whose at-
994 tribute values are closest to a reference value. En-
995 tities are ranked by their absolute distance to the
996 reference, and the closest entities are selected. The
997 output is a subset of entities.

998 **Retrieval** Retrieval is applied to a uniquely de- 1047
999 termined entity. Given a target attribute, the corre- 1048
1000 sponding attribute value is retrieved from the table 1049
1001 and returned as the final answer. 1050

1002 **Summation** Summation aggregates a numerical 1051
1003 attribute over a filtered entity set by summing all 1052
1004 corresponding values. The output is a scalar nu- 1053
1005 merical value. 1054

1006 **Mean** Mean computes the arithmetic average of 1055
1007 a numerical attribute over a filtered entity set. In- 1056
1008 stances in which the resulting mean is non-integer 1057
1009 are discarded during dataset construction. The out-
1010 put is a scalar numerical value.

1011 **B.4 Omni-Modal Context Generation**

1012 **Image Modality** We convert tabular data into 1058
1013 image modality using widely adopted visualization 1059
1014 libraries, Matplotlib and Seaborn. A total of ten 1060
1015 chart types are generated: vertical bar, horizontal 1061
1016 bar, vertical stacked bar, horizontal stacked bar, 1062
1017 lollipop, line, scatter, heatmap, bubble, and tile. 1063

1018 To enhance visual diversity, we randomly 1064
1019 select one of 20 fonts for each image, uniformly 1065
1020 sampled across all images. The fonts used 1066
1021 are: Arimo[wght], FiraSansCondensed-Regular, 1067
1022 OpenSans-Regular, RobotoSlab[wght], WorkSans- 1068
1023 VariableFont[wght], CALIBRI, Kosugi-Regular, 1069
1024 OpenSansHebrew-Regular, SourceSansPro- 1070
1025 Regular, arial, EBGaramond-VariableFont[wght], 1071
1026 Lato-Regular, OpenSansHebrewCondensed- 1072
1027 Regular, Tinos-Regular, tahoma, FiraSans-Regular, 1073
1028 NotoSans-Regular, Roboto-Regular, Ubuntu- 1074
1029 Regular, and times. 1075

1030 In addition, chart elements are colored using 1076
1031 a fixed palette of 20 distinct colors: #4E79A7, 1077
1032 #A0CBE8, #F28E2B, #FFBE7D, #59A14F, 1078
1033 #8CD17D, #B6992D, #F1CE63, #499894, 1079
1034 #86BCB6, #E15759, #FF9D9A, #79706E, 1080
1035 #BAB0AC, #D37295, #FABFD2, #B07AA1, 1081
1036 #D4A6C8, #9D7660, and #D7B5A6. These design 1082
1037 choices allow us to construct a rich and diverse set 1083
1038 of images. 1084

1039 **Text Modality** We define a total of 24 represen- 1085
1040 tative text scenarios across four domains. For the 1086
1041 finance domain, the scenarios include Analyst Re- 1087
1042 port, News Article, Blog Post, Email Newsletter, 1088
1043 Executive Summary, and Meeting Minutes. For the 1089
1044 economics domain, the scenarios include Analyst 1090
1045 Report, News Article, Blog Post, Email Newsletter, 1091
1046 Executive Summary, and Meeting Minutes. For 1092

the climate domain, the scenarios include Research 1047
Report, Business Report, City Marketing, News 1048
Article, Blog Post, and Magazine. For the nutri- 1049
tion domain, the scenarios include Research Report, 1050
Quality Assurance Log, Dietary Guidelines, Ingre- 1051
dient Encyclopedia, Blog Post, and Magazine. 1052

For each scenario, we design a tailored situa- 1053
tional prompt and use large language models to con- 1054
vert structured tabular data into natural language 1055
text. Figure 19 illustrates an example of the prompt 1056
used for text scenario generation. 1057

Speech Modality We define 22 representative 1058
speech scenarios across four domains. For the fi- 1059
nance domain, the scenarios include Meeting, Pod- 1060
cast, Seminar, Audit, and News Debate. For the 1061
economics domain, the scenarios include Meet- 1062
ing, Podcast, Seminar, News Debate, and Global 1063
Summit. For the climate domain, the scenarios 1064
include Weather Forecast, Meeting, Airport Con- 1065
trol Tower, Sports Event Briefing, Business Risk 1066
Briefing, and Green Energy Assessment. For the 1067
nutrition domain, the scenarios include Meeting, 1068
Lab Briefing, Documentary, Ingredient Safety Au- 1069
dit, Conference, and Podcast. Each scenario con- 1070
sists of a four-speaker dialogue. One speaker 1071
serves as the moderator, while the remaining three 1072
speakers are each responsible for different at- 1073
tributes. Similar to the text modality, we first 1074
use Large Language Models to convert structured 1075
tabular data into textual scripts tailored to each 1076
scenario. These scripts are then converted into 1077
speech using the Kokoro-82M TTS model (Hex- 1078
grad, 2025). To enhance acoustic diversity, we uti- 1079
lize 27 distinct voices provided by the Kokoro-82M 1080
model, including 14 female and 13 male voices.¹³ 1081
The female voices include: af_heart, af_alloy, 1082
af_aoede, af_bella, af_jessica, af_kore, af_nova, 1083
af_river, af_sarah, af_sky, bf_alice, bf_emma, 1084
bf_isabella, bf_lily. The male voices include: 1085
am_adam, am_echo, am_eric, am_fenrir, am_liam, 1086
am_michael, am_onyx, am_puck, am_santa, 1087
bm_daniel, bm_fable, bm_george, bm_lewis. This 1088
voice configuration ensures a high degree of varia- 1089
tion in the generated speech data. Figure 20 illus- 1090
trates an example of the prompt used for speech 1091
scenario generation. 1092

¹³The full list of available voices is provided at <https://huggingface.co/hexgrad/Kokoro-82M/blob/main/VOICES.md>

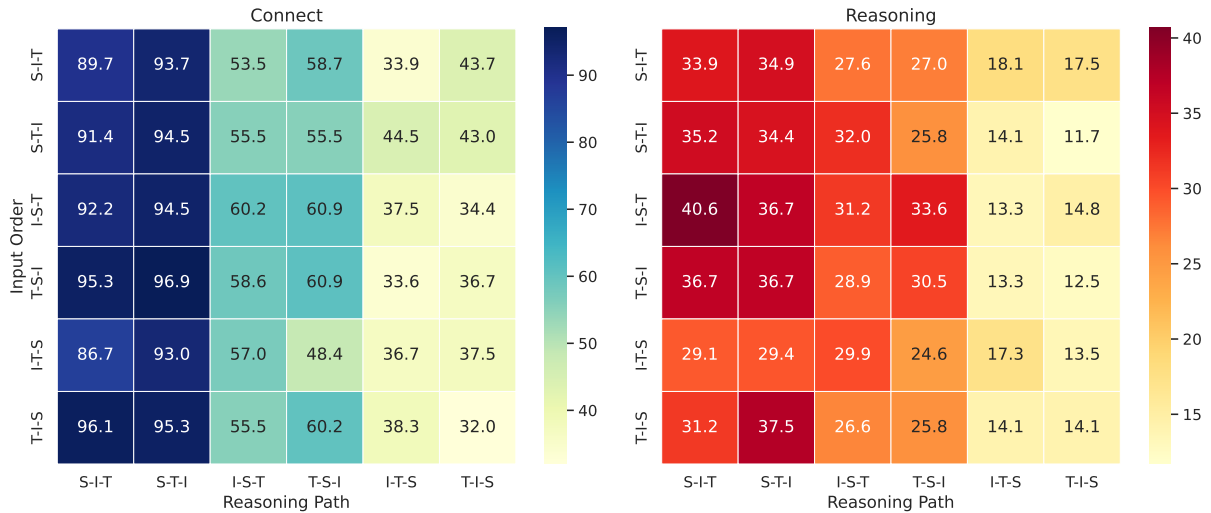


Figure 12: Heatmaps of performance across six input modality orders and six reasoning paths for OMHBench-Connect and OMHBench-Reasoning, evaluated using the Gemini 3 Flash.

C Significance of Input Modality Order

Prior work has shown that altering input order can affect a model’s behavior (Chen et al., 2024; Tan et al., 2024). Moreover, models tend to prioritize information presented in earlier tokens, potentially influencing the reasoning process (Liu et al., 2023b; Wallace et al., 2024). To examine the impact of input order on performance metrics, we conducted a systematic experiment on the OMHBench with Gemini 3 Flash. We enumerated all six input permutations and evaluated accuracy over all 36 combinations of reasoning paths and input orders, as shown in Figure 12. We observe non-negligible performance variation across input orders. For OMHBench-Connect the accuracy varies by up to 12.5 percentage points across input orders while holding the reasoning path fixed (T-S-I). For OMHBench-Reasoning, the corresponding variation is at most 11.5 percentage points (S-I-T). While the overall trends across reasoning paths are broadly consistent, these accuracy fluctuations induced by input-order can still add noise to comparative analysis. Therefore, we randomize input order in our main experiments to reduce this source of variance.

D Experiments Environment

All experiments were conducted on a machine equipped with Intel Xeon Gold 6338 CPU (2.00 GHz), and an NVIDIA A100-SXM4 GPU with 80 GB of memory. The system ran Ubuntu 22.04.4 LTS with CUDA compilation tools release 12.4.

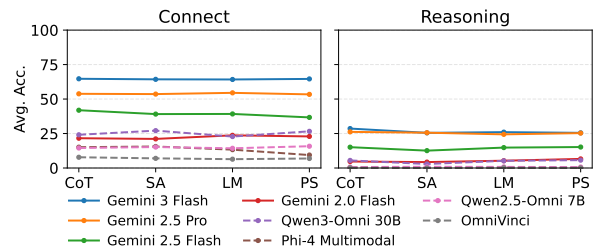


Figure 13: Performance of four prompting methods on OMHBench: Chain-of-Thought (CoT), Self-Ask (SA), Least-to-Most (LM), and Plan-and-Solve (PS). They yield limited gains, calling for dedicated future research.

We used Python 3.10.18 and PyTorch 2.6.0+cu124 as the core software environment. During both dataset generation and evaluation, the random seed was fixed to 42 to ensure reproducibility.

Model	Accuracy by Reasoning Path (%)						Avg. (Acc.)	PBS
	S-I-T	S-T-I	I-S-T	T-S-I	I-T-S	T-I-S		
<i>Proprietary Models</i>								
Gemini 3 Flash (Think)	97.5	98.4	75.4	75.0	60.2	63.5	78.3	32.2
Gemini 3 Flash (Non-Think)	93.0	94.5	60.9	57.8	43.0	39.1	64.7	8.6
Gemini 2.5 Pro (Think)	94.5	96.9	66.4	71.1	55.5	50.8	72.5	25.0
Gemini 2.5 Pro (Non-Think)	65.6	68.8	50.0	58.6	41.4	38.3	53.8	6.2
Gemini 2.5 Flash (Think)	82.0	85.9	50.8	54.7	26.6	21.9	53.6	4.7
Gemini 2.5 Flash (Non-Think)	65.6	69.5	37.5	39.1	24.2	15.6	41.9	2.3
Gemini 2.5 Flash-lite (Think)	49.2	60.9	38.3	35.2	5.5	4.7	32.3	0.0
Gemini 2.5 Flash-lite (Non-Think)	32.8	37.5	28.9	26.6	2.3	0.8	21.5	0.0
Gemini 2.0 Flash	28.9	33.6	26.6	29.7	4.7	6.2	21.6	0.0
Gemini 2.0 Flash-lite	35.9	32.8	21.1	11.7	2.3	2.3	17.7	0.0
<i>Open-Source Models</i>								
Qwen3-Omni 30B (Think)	75.8	77.0	46.7	49.6	16.0	16.0	46.8	2.3
Qwen3-Omni 30B (Non-Think)	35.0	44.7	17.8	33.8	6.2	7.2	24.1	0.0
Phi-4 Multimodal	26.6	23.6	21.5	18.4	0.6	0.0	15.1	0.0
Qwen2.5-Omni 7B	22.7	20.9	19.3	20.5	2.0	1.8	14.5	0.0
Qwen2.5-Omni 3B	12.7	17.6	15.6	14.6	1.2	2.0	10.6	0.0
OmniVinci	14.8	8.6	14.8	7.0	0.8	0.6	7.8	0.0
MiniCPM-o 2.6	8.0	10.9	7.4	8.4	1.2	0.2	6.0	0.0
Omni-AutoThink	7.6	6.6	8.0	6.1	0.6	0.0	4.8	0.0

Table 6: Accuracies and Path Balance Scores (PBSs) across six reasoning paths in **OMHBench-Connect**, including both thinking and non-thinking variants. Avg denotes macro-averaged accuracy. PBSs measure robustness to reasoning path variations.

Model	Accuracy by Reasoning Path (%)						Avg. (Acc.)	PBS
	S-I-T	S-T-I	I-S-T	T-S-I	I-T-S	T-I-S		
<i>Proprietary Models</i>								
Gemini 3 Flash (Think)	55.9	58.8	49.8	49.6	40.0	42.6	49.4	8.6
Gemini 3 Flash (Non-Think)	35.9	39.8	31.2	29.7	16.4	18.8	28.6	1.6
Gemini 2.5 Pro (Think)	53.9	51.6	52.3	47.7	41.4	46.1	48.8	10.9
Gemini 2.5 Pro (Non-Think)	28.9	32.8	30.5	28.9	14.8	21.1	26.2	0.0
Gemini 2.5 Flash (Think)	32.0	30.5	17.2	24.2	10.9	10.9	21.0	0.0
Gemini 2.5 Flash (Non-Think)	22.7	24.2	17.2	14.8	6.2	5.5	15.1	0.0
Gemini 2.5 Flash-lite (Think)	18.8	21.1	15.6	8.6	0.0	0.0	10.7	0.0
Gemini 2.5 Flash-lite (Non-Think)	7.8	9.4	3.9	5.5	0.0	0.8	4.6	0.0
Gemini 2.0 Flash	4.7	11.7	4.7	6.2	0.8	0.0	4.7	0.0
Gemini 2.0 Flash-lite	3.9	5.5	3.9	2.3	0.8	0.0	2.7	0.0
<i>Open-Source Models</i>								
Qwen3-Omni 30B (Think)	27.3	28.5	14.1	14.6	2.7	2.7	15.0	0.0
Qwen3-Omni 30B (Non-Think)	9.8	10.5	4.7	6.4	0.8	0.6	5.5	0.0
Phi-4 Multimodal	0.6	0.4	0.2	0.0	0.2	0.2	0.3	0.0
Qwen2.5-Omni 7B	0.4	1.0	1.0	0.6	0.2	1.2	0.7	0.0
Qwen2.5-Omni 3B	0.8	0.6	0.2	0.0	0.4	0.2	0.4	0.0
OmniVinci	0.6	0.2	0.2	0.4	0.0	0.0	0.2	0.0
MiniCPM-o 2.6	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0
Omni-AutoThink	0.4	0.2	0.4	0.2	0.0	0.0	0.2	0.0

Table 7: Accuracies and Path Balance Scores (PBSs) across six reasoning paths in **OMHBench-Reasoning**, including both thinking and non-thinking variants. Avg denotes macro-averaged accuracy. PBSs measure robustness to reasoning path variations.

Think step by step and provide the final answer at the end.

Based on the given question above, build a chain of sub-questions and intermediate answers to solve it.

First, determine if follow-up questions are needed. If yes, output "Are follow up questions needed here: Yes."

Then, explicitly state a "Follow up:" question and provide an "Intermediate answer:" for it.

Repeat this process until you have enough information to determine the final answer.

Finally, conclude your reasoning.

To solve the complex problem above, first decompose it into a series of simpler sub-questions.

Then, solve each sub-question in order, using the information from previous answers to reach the final conclusion.

Your response should follow this structure:

1. Decompose the problem: "To answer the question, we need to know: [List sub-questions]"
2. Solve sequentially: Solve each sub-question one by one.
3. Provide the final answer.

Let's first understand the problem, extract relevant variables and their corresponding numerals, and devise a complete plan.

Then, let's carry out the plan, calculate intermediate variables (pay attention to correct numerical calculation and commonsense), solve the problem step by step, and show the answer.

Figure 14: Prompt formats used for each prompting method: Zero-Shot CoT, Self-Ask, Least-to-Most, and Plan-and-Solve.

Question: {question}
Context: {context}
Entities: {entities}
Answers: {answers}
Question Type: {question_type}

****Mathematical Formulation****
Attribute a_I: Value of some attribute of an entity exclusive to image
Attribute a_T: Value of some attribute of an entity exclusive to text
Relation: Relation connecting Attribute and Entity or two Entities

The question given is formulated as Image-to-Text (I-T) type or Text-to-Image (T-I) type.
I-T question gives value of attribute a_I in the question and asks for a value of attribute a_T for the answer.
T-I question gives value of attribute a_T in the question and asks for a value of attribute a_I for the answer.

Your task is to identify the a_I and a_T from the given entity and question from given context.
Your final answer should be in the form of json containing two items, a_I and a_T.

Figure 15: Prompt for extracting attributes for an entity from each question in MMQA and MuMuQA.

Question: {question}
Question Type: {question_type}
reversed_question: {reversed_question}

Attribute a_I: Value of some attribute of an entity exclusive to image
Attribute a_T: Value of some attribute of an entity exclusive to text

Question Types:
I-T type: a_I -> entity -> a_T.
T-I type: a_T -> entity -> a_I.

Given two multimodal questions, you need to determine whether given Question is type of {question_type}, AND reversed_question is type of {reverse_type}.
After thinking step-by-step, give your final answer as Yes or No.

Figure 17: Prompt for validating the reasoning order of generated questions in MMQA and MuMuQA.

Question: {question}
Context: {context}
Entities: {entities}
Answers: {answers}
Question Type: {question_type}
a_I: {a_I}
a_T: {a_T}

You are given a QA pair involving multimodal reasoning with attributes from two different modalities: image (a_I) and text (a_T).
A relation connects these attributes or the attribute and an entity.

The question given is formulated as I-T type: a_I -> entity -> a_T.
For instance, if question is "What did the government of the person in the image with the grey tie do?", With a_I = color_of_tie, e = person, and a_T (answer) = conduct_investigation

Your task is to reverse the original question, generating a new question (T-I) that:
- Asks for the value of the attribute that was originally given in the question.
- Uses the original answer as a part of the generated question. (e.g. position in the image, color)

Think step by step and give your final answer in the form of json containing 2 items: reversed_question, answer_reversed_question

Figure 16: Prompt for generating question from our formulation in MMQA and MuMUQA.

The provided options are multi-hop options that, given text, table, and chart as inputs, always reason from Image (Chart) to Text.

Each option follows an explicit reasoning path, defined as an ordered sequence of modality relations.
In the original formulation used in this task, the options reason from Chart to Text.

I want to change ONLY the reasoning path so that the option instead reasons from Text to Chart, while keeping the content/meaning as unchanged as possible.

Do NOT add new facts. Do NOT remove important constraints.
Preserve entities, numbers, and conditions.
Return in EXACTLY the following format (one line per item, keep the tags):

[option1]: <rewritten option1>
[option2]: <rewritten option2>
[option3]: <rewritten option3>

Here are the inputs:
option1: {option1}
option2: {option2}
option3: {option3}

Figure 18: Prompt for controlling the reasoning path of each option in FCMR.

<p>You are a financial analyst writing a professional analyst report based on the provided data. Write a clear, structured, and insight-driven report that includes the following:</p> <ul style="list-style-type: none"> - Executive summary of the key findings - Detailed analysis of all relevant figures - Interpretation of trends and possible implications <p>Use a formal and objective tone suitable for investors and stakeholders. Ensure that every single value from the data is explicitly mentioned. No Markdown: Output strictly as plain text. Do NOT use any Markdown formatting syntax. Do not bold company names or values. Use only standard paragraph breaks for structure. Make sure to cover all information from the given data, whether it relates to companies, financial metrics, or other quantitative or categorical indicators. Present the content in paragraph format with optional section headers. All units are standardized. Do not explicitly mention or append the units.</p>	<p>You are a journalist writing a news article based on the provided data. Write an informative, concise, and neutral report that would appear in a financial news section. The article should:</p> <ul style="list-style-type: none"> - Lead with a headline that reflects the core message - Summarize key highlights in the opening paragraph - Follow with detailed paragraphs that explain each data point or trend - Include all information from the data, regardless of type (companies, financials, etc.) <p>Ensure that every single value from the data is explicitly mentioned. No Markdown: Output strictly as plain text. Do NOT use any Markdown formatting syntax. Do not bold company names or values. Use only standard paragraph breaks for structure. Maintain an objective tone throughout the piece. Use paragraph-based prose, not a script or dialogue. All units are standardized. Do not explicitly mention or append the units.</p>
--	--

Figure 19: Example prompts used to generate the text.

<p>Role: Financial Meeting Scriptwriter</p> <p>Task: Convert the provided financial data table into a meeting transcript. The listener must identify the attribute ONLY by recognizing the speaker's voice from the introduction.</p> <p>Context: - [Speaker1]: The Moderator. - [Speaker2]: Handles the metric in the 2nd Column. - [Speaker3]: Handles the metric in the 3rd Column. - [Speaker4]: Handles the metric in the 4th Column.</p> <p>Strict Guidelines: 1. Intro (The Voice Mapping Phase): - [Speaker1] MUST ask "Who is in charge of [Header Name]?" - The corresponding Speaker MUST reply "That is me" or "I am." - Repeat this for all 3 metrics. This is the ONLY time headers are mentioned.</p> <p>2. Body (The Test): - [Speaker1] simply announces the Company Name (e.g., "Let's look at Company X."). - [Speaker2], [Speaker3], [Speaker4] must immediately report their numbers in a RANDOM sequence. - CRITICAL: The text MUST NOT contain clues like "My asset value is..." or "The debt is...". - Use generic phrases: "I have...", "On my sheet...", "It is...".</p> <p>3. Randomization: - For every company, the order of Sp2, Sp3, Sp4 MUST be shuffled. - Example: Co A (2->3->4), Co B (4->2->3), Co C (3->4->2).</p>	<p>4. Format: Plain text only. NO Markdown. All units are standardized. Do not explicitly mention or append the units.</p> <p>Example for Learning:</p> <p>[Input Data] Company,Receivables,Payables,Cash X,1000,500,200 Y,3000,100,500</p> <p>[Desired Output] [Speaker1]: Let's get started. Who is tracking Receivables? [Speaker2]: That would be me. [Speaker1]: And who has the Payables figures? [Speaker3]: I have those numbers. [Speaker1]: Lastly, who is monitoring Cash? [Speaker4]: I'm ready with that.</p> <p>[Speaker1]: Great. Let's look at Company Y first. [Speaker4]: On my end, it seems moderate. I'm seeing 500. [Speaker2]: Well, looking at my sheet, it's quite high. It's sitting at 3000. [Speaker3]: My figure is actually surprisingly low. It's just 100.</p> <p>[Speaker1]: Okay, moving on to Company X. [Speaker3]: The number on my list is 500. [Speaker4]: It's a bit lower on my side. The value is 200. [Speaker2]: I have a round number here. It shows exactly 1000.</p> <p>---</p> <p>Real Task: Generate the script based on this Real Data: {table_data_json}</p>
---	---

Figure 20: Example prompt used to generate the speech script.

<p># Role You are an expert evaluator for Multi-hop Question Answering models. Your task is to analyze a candidate model's response against a specific Question and Ground Truth Answer, diagnosing exactly where the reasoning chain failed.</p> <p># Input Data Question: {question} Ground Truth Answer: {ground_truth_answer}</p> <p>Candidate Model Response: {other_model_response START} {model_response} {other_model_response END}</p> <p># Reasoning Logic Chain This multi-hop QA problem relies on a specific chain of reasoning connecting attributes and entities. The correct logical path is defined as follows:</p> <ol style="list-style-type: none"> 1. Attribute 1 Value: {attr_1_val} 2. Entity 1: {ent_1} 3. Attribute 2 Value: {attr_2_val} 4. Entity 2: {ent_2} 5. Attribute 3 Value (Final Answer): {attr_3_val} 	<p># Evaluation Rules Compare the "Candidate Model Response" against the "Reasoning Logic Chain" above. Identify the *first* step where the model failed. Assign a failure case based on the following taxonomy:</p> <ul style="list-style-type: none"> * Case 1: The model correctly recognized 'Attribute 1 Value' but failed to identify 'Entity 1'. * Case 2: The model correctly recognized 'Attribute 1 Value' and 'Entity 1', but failed to recognize 'Attribute 2 Value'. * Case 3: The model correctly recognized 'Attribute 1 Value', 'Entity 1', and 'Attribute 2 Value', but failed to identify 'Entity 2'. * Case 4: The model correctly recognized everything up to 'Entity 2', but failed to derive the final 'Attribute 3 Value'. * Case 5: The model successfully followed the entire chain and the final answer is correct. <p># Equivalence Guidelines * Numerical Formatting: Ignore commas or formatting differences (e.g., "18,592" is equal to "18592"). * Units: Ignore unit discrepancies if the core value is correct (e.g., "26 million" is considered equal to "26" if the target value is 26).</p> <p># Response Format Provide a brief explanation of your analysis followed by the final case number.</p> <p>Explain: <Explanation of the reasoning and the specific point of failure> Final case: <1/2/3/4/5></p>
---	---

Figure 21: Prompt used to Step-by-Step Failure Analysis.

OMHBench-Connect	[Reasoning Path: I-T-S]	Lookup-Comparison-Retrieval																				
Question: What is the Total Revenue of the company whose Inventory is larger than that of the company with Goodwill equal to 40303?																						
<table border="1"> <caption>Company Financial Data</caption> <thead> <tr> <th>Company</th> <th>Selling General And Administration</th> <th>Goodwill</th> <th>Cost Of Revenue</th> </tr> </thead> <tbody> <tr> <td>L</td> <td>37190</td> <td>17534</td> <td>41744</td> </tr> <tr> <td>S</td> <td>8595</td> <td>45853</td> <td>25177</td> </tr> <tr> <td>G</td> <td>23305</td> <td>40303</td> <td>3849</td> </tr> <tr> <td>K</td> <td>14730</td> <td>69527</td> <td>17851</td> </tr> </tbody> </table>	Company	Selling General And Administration	Goodwill	Cost Of Revenue	L	37190	17534	41744	S	8595	45853	25177	G	23305	40303	3849	K	14730	69527	17851	<p>KEY FINDINGS Overall liquidity is concentrated in a few entities, with company T holding cash and equivalents of 13947 while maintaining very low inventory of 1289 ... For company G, inventory is 7016. cash and equivalents are 9482, and depreciation is 2896, combining high inventory and strong cash with comparatively modest depreciation ... RISK FACTORS The combination of very high inventory at K of 10851, very low cash and equivalents</p>	<p>[Speaker1]: Let's start. Is the Receivables ledger open? [Speaker2]: Open and ready. ... [Speaker1]: Do we have the Total Revenue records? [Speaker3]: Yes, ready. [Speaker1]: And the Net PPE confirmation? [Speaker4]: Prepared. ... [Speaker1]: Company K. Go. [Speaker3]: Confirmed at 63627.</p>
Company	Selling General And Administration	Goodwill	Cost Of Revenue																			
L	37190	17534	41744																			
S	8595	45853	25177																			
G	23305	40303	3849																			
K	14730	69527	17851																			
Image	Text	Speech																				
Answer: 63627																						

Figure 22: Example of OMHBench-Connect I-T-S Instance.

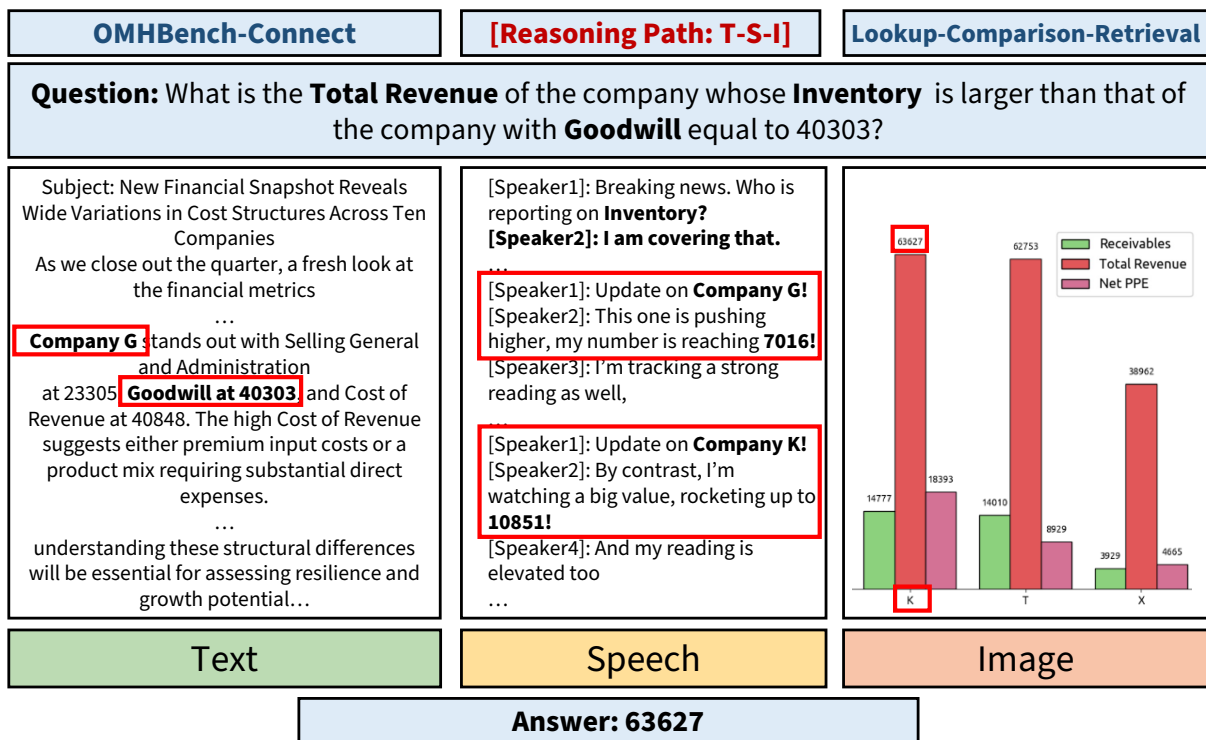


Figure 23: Example of OMHBench-Connect T-S-I Instance.

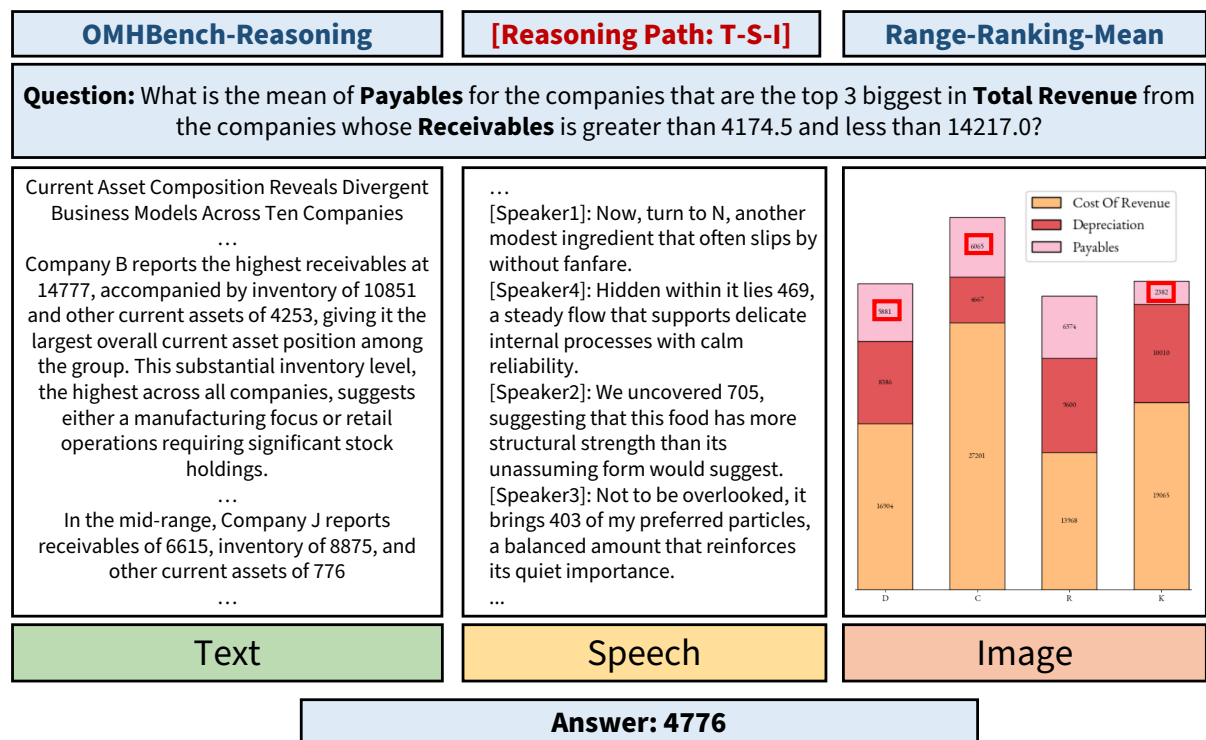


Figure 24: Example of OMHBench-Reasoning T-S-I Instance.