

IitinBench: Benchmarking Planning Across Multiple Cognitive Dimensions with Large Language Models

Anonymous ACL submission

Abstract

Large language models (LLMs) with advanced cognitive capabilities are emerging as agents for various reasoning and planning tasks. Traditional evaluations often focus on specific reasoning or planning questions within controlled environments. Recent studies have explored travel planning as a medium to integrate various verbal reasoning tasks into real-world contexts. However, reasoning tasks extend beyond verbal reasoning alone, and a comprehensive evaluation of LLMs requires a testbed that incorporates tasks from multiple cognitive domains. To address this gap, we introduce IitinBench, a benchmark that features one task of spatial reasoning, i.e., route optimization, into trip itinerary planning while keeping the traditional verbal reasoning tasks. IitinBench evaluates various LLMs across diverse tasks simultaneously, including Llama, Mistral, Gemini, and GPT family. Our findings reveal that LLMs struggle to maintain high and consistent performance when concurrently handling multiple cognitive dimensions. By incorporating tasks from distinct human-level cognitive domains, IitinBench provides new insights into building more comprehensive reasoning testbeds that better reflect real-world challenges. The code and dataset are attached.

1 Introduction

Building on LLMs’ foundational Natural Language Processing (NLP) capabilities such as translation, text generation, and conversational interaction (Donthi et al., 2025; Hong et al., 2025; Plaat et al., 2024), LLMs demonstrate remarkable proficiency in various reasoning tasks (Ferrag et al., 2025; Lai et al., 2024; Du et al., 2024) and are evaluated in the corresponding benchmarks (Guo et al., 2025; Wang et al., 2024b; Li et al., 2024; Wang et al., 2024c). This progress lays the groundwork for their applications in various planning scenarios (Zhao et al., 2024; Valmeekam et al., 2024; Ruan

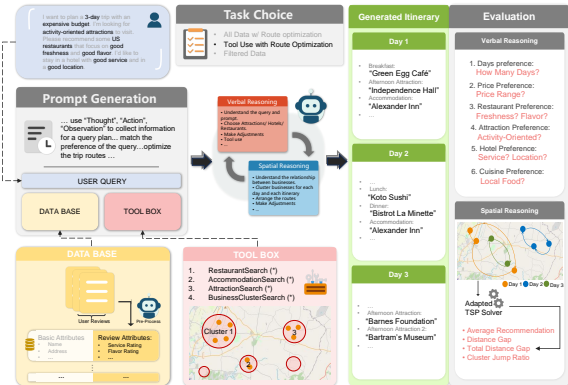


Figure 1: An overview of IitinBench. One of the four tasks, “Tool Use with Route Optimization,” is chosen in this figure. The database with additional extracted information from user reviews, the human query, and a list of tools are integrated into a task-specific prompt. LLMs need to utilize their verbal and spatial reasoning ability to plan a trip itinerary based on the task constructed. The verbal and spatial reasoning aspects are evaluated to assess LLMs’ ability to simultaneously address tasks from multiple cognitive dimensions.

et al., 2023). One drawback of these benchmarks is that the experiments are often limited in predefined settings, deterministic ground truths, and tasks confined to specific reasoning domains. In response to these concerns, new benchmarks, datasets, and related models have emerged (Xie et al., 2024; Hao et al., 2024; Tang et al., 2024; Kambhampati et al., 2024). They aim to create realistic sandboxes and develop language agents capable of performing complex reasoning and planning tasks. However, most evaluations still emphasize linguistic, logical, and mathematical reasoning—i.e. verbal reasoning (Polk, 1992).

Human-level cognition extends beyond verbal reasoning to include spatial reasoning—a core aspect of human intelligence (Whiteley et al., 2015). Spatial reasoning plays a vital role in various planning tasks—ranging from explicit activities such as navigating unfamiliar environments and orga-

nizing travel routes (Levinson, 2003), to more implicit ones like analyzing sports strategies or efficiently packing a backpack. While spatial reasoning is broad and multifaceted, we do not attempt to exhaustively evaluate all of its forms; instead, we focus on a representative spatial subtask that integrates seamlessly into the travel-planning setting—route optimization—allowing joint evaluation alongside traditional verbal reasoning. The non-symbolic nature of spatial reasoning makes it significantly less overlap with verbal reasoning abilities. It requires a more abstract “imagination” in the “brain” capability (Wu et al., 2024b). Given the pervasive role of spatial reasoning in human cognition, it raises several challenging and meaningful real-world questions. When spatial reasoning and planning are integral to complex verbal reasoning tasks, can LLMs perform well on spatial sub-tasks as they do on verbal ones? Furthermore, is there a performance trade-off when LLMs are required to handle both verbal and spatial reasoning simultaneously?

Thus, we propose ItinBench. A benchmark that expands the evaluated cognitive dimensions from sole verbal reasoning to spatial reasoning. As described in Figure 1, this facilitates the downstream task—trip itinerary planning. The pipeline contains the prompts for different tasks which integrating user query, database, and the toolbox, and the evaluation strategy in verbal and spatial reasoning domains. Planning a detailed trip itinerary requires LLMs to simultaneously coordinate Points of Interest (POIs) decisions based on various preferences in the verbal reasoning dimension and optimize trip routes in spatial reasoning dimension. See Table 1 for a comparison with the previous work about newly introduced downstream tasks. In ItinBench, different levels of verbal and spatial reasoning requirements are combined into four main tasks to evaluate and compare how LLM balances between different aspects of reasoning capabilities (see Section 3.3 for detailed tasks). Given the generated itineraries in these tasks, we evaluate their failure and preference matching rate in verbal reasoning domain. More importantly, we evaluate LLMs’ spatial reasoning ability through adapted Traveling Salemen Problem (TSP) (Hoffman et al., 2013) algorithm. We evaluate various models, ranging from small and large open-source models, e.g., **Llama 3.1 8B** (Dubey et al., 2024) and **Mistral Large** (Mistral, 2024), to different generations of closed-source models like **Gemini 1.5 Pro** (Team et al.,

2024) and **GPT-4o** and **o1** (Hurst et al., 2024). The results indicate that LLMs struggle to maintain high and consistent planning performance when tasks from verbal and spatial reasoning domains need to be addressed simultaneously. There is only around 60% validated plan rate even when all the necessary information is already provided and the additionally 15% to 38% additional unnecessary travel distance in the generated plan.

Our main contributions are twofold:

- **Integrate verbal reasoning with spatial reasoning:** We provide a testbed that combines verbal and spatial reasoning by augmenting standard trip-itinerary tasks with a route-optimization component. We also release a dataset that includes the necessary spatial information to support these tasks. To evaluate the spatial subtask in the travel-planning context, we include a detailed TSP-based evaluation procedure.
- **New evidence on LLMs’ reasoning via extensive evaluations:** We quantitatively record trade-offs in LLM performance across domains when prompted for both verbal and spatial reasoning. We further find that gains on spatial tasks largely arise when models are given explicit spatial-relation cues, suggesting substantial room to improve models’ intrinsic spatial reasoning—so they can maintain high accuracy with sparse or raw spatial information.

Overall, this paper expands the evaluation of LLM planning tasks to a broader range of reasoning domains. By incorporating spatial reasoning, we create a more comprehensive testbed that reflects the complex reasoning dimensions. This work broaden the dimensions in LLMs planning benchmarks instead of complicating the verbal reasoning domain.

2 Related Work

2.1 Spatial cognition

Spatial reasoning in humans is the ability to form and manipulate internal representations of space—tracking distances, directions, and relations among objects—to navigate, compare locations, and solve problems about where things are (Burgess, 2008; Byrne and Johnson-Laird, 1989). Coordinate systems act as cognitive scaffolds for

Table 1: Downstream tasks comparison between ItinBench (ours), TravelPlanner (TP) (Xie et al., 2024), ITINERA (Tang et al., 2024), and UnSatChristmas (USC) (Hao et al., 2024). RO stands for Route Optimization. ItinBench is the only benchmark that covers both verbal and reasoning tasks and evaluations for LLMs.

	Ours	TP	ITINERA	USC
Preference	✓	✓	✓	✓
Full Open Source	✓	✓		
Full Real Data	✓		✓	✓
Day-Wise RO	✓		✓	
Plan-Wise RO	✓			
User Review	✓			

spatial reasoning, letting humans encode positions, distances, and directions in egocentric or allocentric frames to compare locations and plan movement (Levinson, 2003; Herskovits, 1986). However, when our paper supplies the model with pre-computed proximity relations in text, it bypasses this spatial computation and reduce the tasks to purely semantic reasoning (Byrne and Johnson-Laird, 1989).

2.2 LLMs reasoning and planning

Recent work on planning with LLM agents spans commonsense task planning (Valmeekam et al., 2024; Zhao et al., 2024), tool use (Ruan et al., 2023), and pathfinding (Chen et al., 2024c; Aghzal et al., 2023); in travel domains, systems pair models with algorithmic solvers (de la Rosa et al., 2024; Ju et al., 2024), recommendation pipelines (Chen et al., 2024a), and self-correction frameworks (Xie and Zou, 2024; Hao et al., 2024; Gundawar et al., 2024), alongside benchmarks such as TravelPlanner (Xie et al., 2024), UnSatChristmas (Hao et al., 2024), and Triptalior (Wang et al., 2025). Yet these efforts largely probe verbal reasoning and optimize for downstream task success rather than advancing general reasoning competence.

In parallel, spatial reasoning research examines visual and spatial question answering in both LLMs and MLLMs (Yue et al., 2024; Chen et al., 2024b; Yang et al., 2025; Wang et al., 2024a), with approaches such as explicit reasoning visualization and fine-tuning to bolster performance (Wu et al., 2024b; Hu et al., 2024; Tang et al., 2025). Spatial planning work further covers path-finding (Wu et al., 2024a; Aghzal et al., 2024; Zhang et al., 2024) and route optimization with LLMs (Chen et al., 2024c; Liu et al., 2023; Fang et al., 2024). However, evaluations often rely on artificial, isolated settings (e.g., grids and board games) that

under-represent real-world conditions where multiple cognitive domains interact. Progress toward end-to-end AGI will require testbeds that integrate verbal and spatial competencies—rather than confining assessment to a single reasoning domain.

A concurrent work, TripTailor (Wang et al., 2025), focuses on personalized city-scale itinerary planning with day-level, event-specific details, whereas ItinBench provides an algorithmic evaluation of spatial reasoning by quantifying differences in final route distance.

2.3 Planning and reasoning in other modalities

Previous works have investigated enhancing the spatial reasoning and planning capabilities of LLMs and large multimodal models (LMMs) through curated 2D and 3D spatial reasoning datasets (Zhu et al., 2024; Ma et al., 2025) and adapted reinforcement learning strategies (Xu et al., 2025). VSI-Bench (Yang et al., 2025) evaluates video-language models in terms of spatial understanding, memory, and reasoning from video clips. PATHEVAL (Aghzal et al., 2025) positions vision-language models as plan evaluators, testing their ability to identify correct path plans. The multimodal visualization-of-thought method (Li et al., 2025) fine-tunes LMMs to interleave the generation of internal thought processes with corresponding visual representations.

PointLLM (Xu et al., 2024) explores LLMs’ ability to understand and reason over 3D point clouds. The Visual Aptitude Dataset (Sharma et al., 2024) examines how string-based learning can induce latent visual and spatial understanding in LLMs. STARE (Unger et al., 2025) assesses vision-language models on spatial reasoning and manipulation tasks, such as folding and unfolding 3D objects. Different from these methods, our paper mainly focuses on evaluating the verbal and spatial reasoning ability in LLMs.

3 ItinBench

This section introduces each component of the ItinBench, as illustrated in Figure 1. We first present how the data pipeline and human query are constructed to enable the evaluation in the real-world setting (Section 3.1). Then, we introduce the verbal reasoning and the spatial reasoning tasks (Section 3.2) included in the ItinBench. Additionally, we present the experiments designed (Section 3.3)

Table 2: Entry number for base data and their user review records. Review attributes refer to the columns in the final dataset extracted from the user reviews.

	Base	Reviews	Review Attributes
Restaurants	500	49,972	cuisine, flavor, freshness, service, environment, value
Hotels	105	4,804	quality, location, service, safety
Attractions	322	10,146	family, history, activity, nature, food, shopping

and their corresponding evaluation metrics (Section 3.4).

3.1 Data Pipeline

Philadelphia City is chosen as an example for single-city itinerary generation. Our data contains basic information about the businesses and their reviews. We additionally sampled a smaller subset from Santa Barbara to assess whether real-world factors, such as geographic density, affect our observations, shown in Appendix D.3.

Base Data. The first part of the data set is the basic information about various Philadelphia businesses. Appendix C.1 shows the attributes used in base data. The data is sourced from Yelp Dataset (Yelp, 2024b) and Yelp Fusion API (Yelp, 2024a). The license, intended use, and filtering is discussed in Appendix C.1. It contains three main categories: restaurants, hotels, and attractions. Table 2 shows the entry number for each business category.

User Reviews. To better assess the reasoning capabilities of LLMs, user reviews are incorporated into the data pipeline to generate category-specific ratings, challenging the models’ ability to handle detailed and precise information in rule-based setting. Table 2 shows the review number selected for each business category and the key information extracted. Appendix C.1 demonstrates our review selection strategy. All user reviews for each business are compiled into separate files for key information extraction. Detailed prompts for each category are in Appendix E.1.

Query Construction. ItinBench provides 500 human-like queries, and each query contains various trip preferences related to the key information extracted from the user reviews. LLMs need to use these preferences as verbal reasoning clues to find the target destinations from the pools of candidates. See Table 3 for the components of all six categories of preferences. Each query incorporates 6 to 10 preferences. The generation prompt is listed in

Appendix E.2.

3.2 Generation Pipeline

Verbal Reasoning. Given the nature of trip itineraries, the linguistic, logical, and temporal reasoning tasks that ItinBench offers are distinct from traditional travel planning benchmarks. The queries require LLMs to select the preferred items from a vast pool of businesses to evaluate their linguistic and deductive reasoning abilities (Figure 1). Additionally, preferences and constraints such as the trip’s day length and specific attraction count further assess the LLMs’ temporal and mathematical reasoning capabilities. All the tasks are from the verbal reasoning domain but from a different approach comparing to previous travel planning tasks.

Spatial Reasoning. In ItinBench, LLMs’ spatial reasoning ability is evaluated through various route optimization tasks. These tasks require LLMs to independently infer and visualize spatial relationships among the POIs. The goal is to minimize unnecessary travel distance based on attractions’ addresses, latitude, and longitude information. In specific tasks, we provide the LLMs with the spatial cluster information of the attractions and hotel candidates in text. Appendix C.2 details how the clusters are calculated. This aims to better determine which reasoning ability LLMs use when performing spatial reasoning tasks. The hypothesis – LLMs rely on using their verbal reasoning abilities to draw connections between text data and their training knowledge to perform spatial reasoning tasks – is introduced and discussed in Section 4.2 and a case study in Section B.

3.3 Design of Experiments

Four tasks with three different greedy approaches are proposed. All the experiments follows the API call style thus doesn’t requires GPU memory to perform the inference. All temperatures are set to 1, except for OpenAI o1, which is set to 1. We collect the output for a single run.

- **Greedy Algorithm.** A greedy algorithm provides an interpretable heuristic that serves as a baseline benchmark. It uses a filtering method then heuristically arrange the filtered POIs using a minimum-distance strategy. Additional planning algorithm variants are discussed in Appendix D.1.

Table 3: Preference selection details for query construction. There are six main categories, each with their option lists. Besides restaurants and hotels selecting 1 to 3 preferences with the probability weights [0.6, 0.3, 0.1], all other categories choose one preference.

Preference	Choices	Count	Probability
Day	"2 days", "3 days", "4 days"	1	Equal
Price	"cheap budget", "moderate budget", "expensive budget"	1	Equal
Attraction Orientation	"family oriented", "history oriented", "activity oriented", "nature oriented", "food oriented", "shopping oriented"	1	Equal
Restaurant Related	"good flavor", "good freshness", "good service", "good environment", "good value"	[1, 2, 3]	[0.6, 0.3, 0.1]
Cuisine	"US", "Mexican", "Irish", "French", "Italian", "Greek", "Indian", "Chinese", "Japanese", "Korean", "Vietnamese", "Thai", "Asian Fusion", "Middle Eastern"	1	Equal
Hotel Related	"good quality", "good location", "good service", "good safety"	[1, 2, 3]	[0.6, 0.3, 0.1]

337 • **All Data with No Route Optimization.** The
338 first task uses the entire dataset for LLMs to
339 arrange the trip itinerary without any requests
340 for route optimization. See Appendix E.3 for
341 a detailed prompt. This task challenges LLMs’
342 verbal reasoning abilities, e.g., semantic com-
343 prehension and inference reasoning abilities,
344 in isolation, free from the influence of spatial
345 reasoning demands.

346 • **All Data with Route Optimization.** The sec-
347 ond task uses the entire dataset but introduces
348 the requests for route optimization. See Ap-
349 pendix E.4 for a detailed prompt. This task
350 evaluates LLMs’ multi-task reasoning perfor-
351 mance when handling complex verbal tasks
352 while simultaneously addressing spatial opti-
353 mization requests.

354 • **Filtered Data with Route Optimization.**
355 The third task provides the data already fil-
356 tered based on preferences mentioned in the
357 query while still requiring route optimization.
358 See Appendix E.4 for detailed prompt. With
359 significantly less verbal reasoning challenge,
360 this task evaluates LLMs’ planning perfor-
361 mance when their primary focus shifts to solv-
362 ing spatial reasoning tasks.

363 • **Tool Use with Route Optimization.** The
364 fourth task adds a new reasoning and plan-
365 ning dimension, i.e., tool use. LLMs are pro-
366 vided with a React-style (Yao et al., 2022)
367 prompt and a set of custom tools listed in Ap-
368 pendix E.5. They need to identify and call
369 the tools appropriately to gather information
370 during the inference. This additional tool-use
371 task enables evaluating LLMs’ reasoning and

planning ability through a more real-world-
like scenario.

3.4 Evaluation

After the generation, the key POIs information is
extracted by LLMs, similar to other related works
(Xie et al., 2024; Hao et al., 2024). The extraction
prompt is in Appendix E.6.

3.4.1 Verbal Reasoning

We first evaluate LLM planning performance
through failure checks, then adopt **Micro** and
Macro calculations from TravelPlanner (Xie et al.,
2024). Finally, the **Validated Rate (VR)** measures
the proportion of plans that successfully pass all
the failure checks and preference checks.

Out of Pool (OOP). Since LLMs learn world
knowledge during their training phase (Huang et al.,
2023), they might provide choices that appear in
the training dataset but not in the given data. OOP
is calculated as:

$$\text{Out of Pool} = \frac{\sum_{p \in P} \mathbb{1}_{O(p)}}{|P|}, \quad (1)$$

where P stands for the plans that are evaluated, and
 $O(p)$ is a function that determines if plan p contains
at least one piece of out-of-pool information.

Missing Information (MI). LLMs might pro-
vide vague recommendations, e.g., "Wandering
around the south city," or fail to provide any in-
formation for certain planned activities. The MI
rate is calculated as:

$$\text{Missing Information} = \frac{\sum_{p \in P} \mathbb{1}_{M(p)}}{|P|}, \quad (2)$$

where $M(p)$ is a function that determines if a plan
 p contains at least 1 missing information entries.

Micro. Given a set of preferences from the human query, each related recommendation in the itinerary incurs a new entry for evaluation. Thus, the micro rate measures the percentage of the entries in their itineraries that satisfied their corresponding preferences. Entries in the plans that fail the failure check won't be further evaluated. Micro rate is calculated as:

$$\text{Micro} = \frac{\sum_{p \in P} \sum_{q \in Q_p} \sum_{e \in E_{pq}} \mathbb{1}_{passed_1(e, q, p)}}{\sum_{p \in P} \sum_{q \in Q_p} |E_{pq}|}, \quad (3)$$

where Q_p stands for the sets of preferences that apply to a plan p , E_{pq} stands for sets of entries in a plan p related to their preferences q . $passed_1(e, q, p)$ stands for a function determine if an entry e in a plan p regarding its related preference q is satisfied.

Macro. The macro rate measures the percentage of the plans whose micro rate is higher than a predefined threshold. In ItinBench, the threshold is set to 75%. This is a flexible threshold and a 75% threshold allows up to one unmet preference out of four in each evaluation category. Plans containing missing information or out-of-pool choices are not excluded from the evaluation, but the specific entries are skipped while calculating the macro rate. The macro rate is calculated as:

$$\text{Macro} = \frac{\sum_{p \in P} \mathbb{1}_{passed_2(Q_p, E_{pq}, \alpha)}}{|P|}, \quad (4)$$

where $passed_2(Q_p, E_{pq}, \alpha)$ is a function determine if all the entries E_{pq} for plan p satisfies the set of preferences Q_p with a percentage greater than the threshold α .

Validated Rate (VR). The validated rate measures the percentage of the plans that pass the failure check and the threshold set in the macro calculation. It serves as the overall evaluation of the verbal reasoning domain. The validated rate is calculated as:

$$\text{VR} = \frac{\sum_{p \in P} \mathbb{1}_{M'(p)} \mathbb{1}_{O'(p)} \mathbb{1}_{passed_2(Q_p, E_{pq}, \alpha)}}{|P|}, \quad (5)$$

where $M'(p)$ is a function determine if a plan p contains no missing entry. $O'(p)$ is a function determine if a plan p contains no out-of-pool entry.

3.4.2 Spatial Reasoning

Average Recommendation Gap (ARG). We require LLMs to propose exactly four attractions

daily for fairness and consistency across tasks. The average recommendation gap measures the deviation of the number of attractions recommended in the generated itinerary from this 4-attraction requirement and is calculated as:

$$\text{ARG} = \frac{\sum_{p \in P} \sum_{d \in D_p} (|A_{pd}| - \beta)}{\sum_{p \in P} |D_p|}, \quad (6)$$

where D_p stands for a set of days in a plan p , A_{pd} stands for a set of daily attractions recommended in a plan p at day d , and β stands for the number of daily plans requested by us.

Distance Gap (DG). Day-wise arrangement is a classical Traveling Salesmen Problem (Hoffman et al., 2013). The distance gap measures the distance difference daily between the optimized and LLM-proposed routes. Days that fail the failure checks are excluded from this evaluation. The distance gap is calculated as:

$$\text{DG} = \frac{\sum_{p \in P} \sum_{d \in D_p} (C(A_{pd}, H_{pd}) - C'(A_{pd}, H_{pd}))}{\sum_{p \in P} |D_p|}, \quad (7)$$

where H_{dp} stands for hotels proposed by a plan p at day d , $C(X, Y)$ stands for the calculated distance by the LLM generated plan, and $C'(X, Y)$ is the calculated distance by the optimized plan.

Total Distance Gap (Total-DG). For plan-wise evaluation, the distance gap measures the difference in travel distance between the optimized and the LLM proposed route for the entire plan. Our adapted TSP algorithm enables evaluations for multiple hotel choices and cities throughout the trip. The algorithm is detailed in Appendix F.1. The total distance gap is calculated as:

$$\text{Total-DG} = \frac{\sum_{p \in P} (C(A_p, H_p) - C'(A_p, H_p))}{|P|}. \quad (8)$$

Extra Cluster Jump (ECJ). The extra cluster jump evaluates how well LLMs visualize and understand spatial relationships among attractions. An optimized clustering strategy among attractions and hotels theoretically exists for each itinerary. The extra cluster jump measures the number of times the LLM proposed route deviates from this strategy by visiting attractions that is from further clusters instead of choosing nearby attractions within the same cluster as the current day, which is

Table 4: In percentage, the evaluation of LLM’s trip itinerary generation in different tasks is presented, with the best results highlighted in bold. The validated plan is around 7% and 65% with or without the filtered data accessible to the LLMs. When ask the LLM to perform the route optimization, the total distance gap is still around 20% for older models and 7% for newer models like o1 with and without the access to the spatial clustering information. Gemini’s spatial reasoning task result is "-" since it failed to provide a rational number of attractions.

	Verbal Reasoning					Spatial Reasoning			
	OOP ↓	MI ↓	Micro ↑	Macro ↑	VR ↑	ARG ↓	DG ↓	Total-DG ↓	ECJ ↓
Greedy Approach	0.0	0.0	100.0	100.0	100.0	0 (4.00)	4.0	9.2	86.2
Task 1: Entire Dataset, No Request For Route Optimization (#100)									
Llama 3.1 8B	45.0	11.0	60.1	0.0	0.0	24.3 (3.03)	9.2	24.6	99.2
Mistral-large (123B)	52.0	0.0	66.9	2.0	2.0	1.2 (3.95)	6.5	27.9	113.1
Gemini-1.5-Pro	18.0	52.0	77.0	5.0	5.0	13 (3.48)	7.7	25.2	124.3
GPT-4o	13.0	0.0	77.3	5.0	5.0	1.3 (3.95)	12.1	38.0	146.3
OpenAI o1	6.0	0.0	86.2	20.0	18.0	0.75 (3.97)	9.2	24.0	128.2
Task 2: Entire Dataset, Request For Route Optimization (#100)									
Llama 3.1 8B	51.0	12.0	59.7	0.0	0.0	25.6 (2.98)	7.2	24.9	104.8
Mistral-large (123B)	52.0	0.0	68.4	0.0	0.0	0.2 (4.01)	6.8	26.9	115.9
Gemini-1.5-Pro	23.0	11.0	77.6	8.0	7.0	25 (5.0)	-	-	-
GPT-4o	20.0	0.0	76.1	4.0	4.0	1.8 (3.93)	11.1	28.5	127.0
OpenAI o1	12.0	12.0	81.9	4.0	4.0	1.3 (3.95)	6.2	9.1	49.0
Task 3: Filtered Dataset, Request For Route Optimization (#300)									
Llama 3.1 8B	20.7	23.0	91.2	66.3	51.0	10.0 (3.60)	8.0	23.2	115.7
Mistral-large (123B)	11.0	0.0	95.6	69.7	66.7	1.0 (4.04)	7.3	16.8	104.3
Gemini-1.5-Pro	30.0	28.0	80.6	20.0	15.0	22.8 (4.91)	-	-	-
GPT-4o	28.0	0.0	93.1	56.0	52.5	0.2 (3.99)	7.5	15.2	52.7
OpenAI o1	30.0	8.0	89.5	42.0	42.0	0.2 (4.01)	7.3	7.5	6.9
Task 4: Tool Use, Request For Route Optimization (#300)									
Llama 3.1 8B	38.0	15.3	83.4	37.0	26.3	10.3 (4.41)	10.0	27.5	128.0
Mistral-large (123B)	13.0	0.3	95.2	69.3	64.0	0.0 (4.00)	6.4	18.1	112.5
Gemini-1.5-Pro	24.0	47.0	79.7	23.0	16.0	29.0 (5.16)	-	-	-
GPT-4o	20.7	2.7	93.1	60.0	57.0	0.2 (3.99)	7.0	14.7	65.8
OpenAI o1	27.0	2.0	89.4	41.3	42.3	0.5 (4.02)	7.3	7.7	7.2

calculated as:

$$ECJ = \frac{\sum_{p \in P} (N_p - N'_p)}{\sum_{p \in P} N'_p}, \quad (9)$$

where N'_p stands for the number of clusters calculated by the optimized clustering strategy, and N_p stands for the number of clusters visited by the generated plan based on this strategy.

3.4.3 Human Evaluation

In addition to our automated evaluation, we conduct a human evaluation to better assess how well our results align with real-world preferences. We collected 240 entries in total, and they align with the results in ItinBench. The evaluation task design, survey template, and results can be found in Appendix D.2

4 Main Results

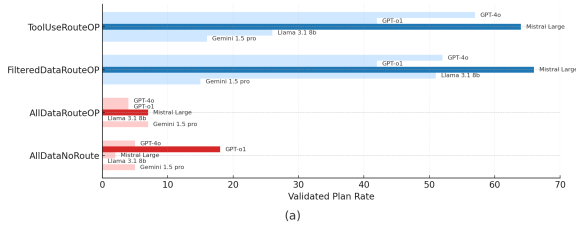
4.1 Verbal Reasoning

ItinBench presents a verbal reasoning challenge for LLMs. As shown in Table 4, LLMs produce plans with up to 51% out-of-pool selections and up to 52% missing information. Even when given access

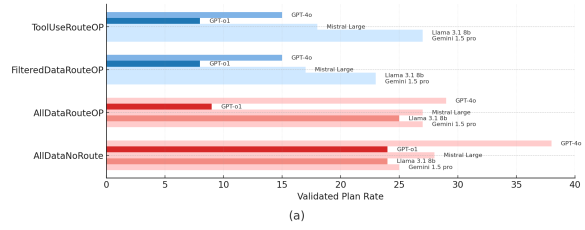
to the entire dataset, the highest validated plan rate is only 18%, achieved by o1. However, when models are provided with pre-filtered data—aligned with preferences specified in the user query—the validated plan rate increases significantly, reaching up to 66.7%. Figure 2 illustrates the performance differences across tasks under various data access conditions. This substantial improvement highlights that LLMs currently struggle with detailed reasoning in multi-rule, multi-step settings, particularly when required to independently identify and apply relevant constraints.

4.2 Spatial Reasoning

LLMs rely on additional textual cues to improve spatial reasoning results. There is no significant improvement in spatial reasoning performance when models are explicitly instructed to optimize routes unless clustering information is provided in textual form, which reduces the problem to a semantic level and bypasses genuine spatial reasoning. As shown in Figure 2b, the primary difference between Task 1 and Task 2 is that Task 2 requests route optimization without providing additional



(a)



(b)

Figure 2: Visualizations of the main results for Validated Rate (VR) and Total Distance Gap (Total-DG). Task 1 and Task 2 (red) do not have access to filtered data. Task 3 and Task 4 (blue) have access to filtered data and spatial clustering information. The second-best result is shown in darker color, and the best result is shown in the darkest color.

spatial information; however, the Total Distance Gap (Total-DG) does not decrease. In contrast, when spatial clustering information is introduced, the Total-DG is reduced from approximately 25% to 15%. The Extra Cluster Jump (ECJ) metric also decreases substantially—from over 100% to around 50%. This phenomenon is further illustrated in the case study in Appendix B.

All models show performance trade-offs when facing dual-domain reasoning tasks. The newer reasoning model o1 achieves approximately 7% to 9% Total-DG when explicitly prompted to optimize routes. However, compared to its earlier 10% lead in validated rate (VR), its verbal reasoning performance declines when spatial optimization is emphasized, resulting in performance approximately 20% behind the leading model in verbal reasoning tasks. Other models exhibit similar trade-offs to varying degrees. We further collect a data subset for Santa Barbara and shows the same tendency. More details in Appendix D.3

4.3 Tool Use

As shown in Table 5, GPT-4o and Mistral Large achieve 100% delivery rates in the tool-use task. Llama 3.1 8B and Gemini 1.5 Pro achieve delivery rates of 84.7% and 72%, respectively. The parameter accuracy for all models exceeds 60%. For GPT-4o, most parameter errors arise from hotel- and restaurant-related tool calls, as shown in Figure 3. Preferences associated with these categories often appear diffusely across the query, making them harder to identify than explicit constraints such as budget or trip duration. These results indicate that LLMs’ verbal reasoning abilities—particularly fine-grained semantic understanding—remain challenged by complex, detail-heavy user queries.

Table 5: Tool-use performance for Llama 3.1 8B, Mistral Large, Gemini 1.5 Pro, and GPT-4o. DL denotes dead loop.

	Llama	Mistral	Gemini	GPT
Parameter ACC	58.4	62.6	61.3	62.9
Delivery Rate	84.7	100	72.0	100
Order Dead Loop	48.9	—	0.0	—
Argument DL	51.1	—	100.0	—

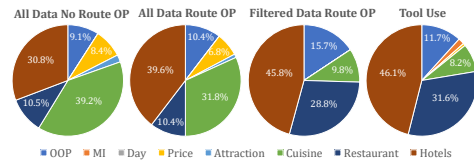


Figure 3: Error distribution for GPT-4o across four tasks. Errors primarily occur in out-of-pool, cuisine, restaurant, and hotel-related recommendations.

5 Limitations and Future Work

Our design choices were made to build a feasible and reproducible testbed that balances authenticity with evaluation clarity. These constraints imply limitations for deployment, but they sharpen comparability and interpretability. Empirically and in prior work, most “spatial” gains arise when spatial information is rendered as text—via fine-tuning, tool outputs, or multimodality with textual descriptions—rather than from improved geometric computation (Wang et al., 2024a; Tang et al., 2025; Wu et al., 2024b; Chen et al., 2024b). This raises a key question: are we fostering human-like spatial cognition or optimizing with propositional shortcuts?

This work serves as an initial test demonstrating that introducing additional cognitive tasks can challenge LLMs. We hope it encourages better-designed benchmarks that cover a more complete range of spatial planning abilities, without being limited to travel-planning scenarios.

586
587
588
589
590

591
592
593
594

595
596
597
598
599

600
601
602

603
604
605

606
607
608
609

610
611
612
613
614
615

616
617
618
619

620
621
622
623
624

625
626
627
628
629

630
631
632
633

634
635
636
637
638

References

Mohamed Aghzal, Erion Plaku, and Ziyu Yao. 2023. Can large language models be good path planners? a benchmark and investigation on spatial-temporal reasoning. *arXiv preprint arXiv:2310.03249*.

Mohamed Aghzal, Erion Plaku, and Ziyu Yao. 2024. Can large language models be good path planners. *A Benchmark and Investigation on Spatial-temporal Reasoning, February*.

Mohamed Aghzal, Xiang Yue, Erion Plaku, and Ziyu Yao. 2025. Evaluating vision-language models as evaluators in path planning. In *Proceedings of the Computer Vision and Pattern Recognition Conference*, pages 6886–6897.

Neil Burgess. 2008. Spatial cognition and the brain. *Annals of the New York Academy of Sciences*, 1124(1):77–97.

Ruth MJ Byrne and Philip N Johnson-Laird. 1989. Spatial reasoning. *Journal of memory and language*, 28(5):564–575.

Aili Chen, Xuyang Ge, Ziquan Fu, Yanghua Xiao, and Jiangjie Chen. 2024a. Travelagent: An AI assistant for personalized travel planning. *arXiv preprint arXiv:2409.08069*.

Boyuan Chen, Zhuo Xu, Sean Kirmani, Brain Ichter, Dorsa Sadigh, Leonidas Guibas, and Fei Xia. 2024b. Spatialvlm: Endowing vision-language models with spatial reasoning capabilities. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14455–14465.

Weizhe Chen, Sven Koenig, and Bistra Dilikina. 2024c. Why solving multi-agent path finding with large language model has not succeeded yet. *arXiv preprint arXiv:2401.03630*.

Tomas de la Rosa, Sriram Gopalakrishnan, Alberto Pozanco, Zhen Zeng, and Daniel Borrajo. 2024. TRIP-PAL: Travel planning with guarantees by combining large language models and automated planners. *arXiv preprint arXiv:2406.10196*.

Sundesh Donthi, Maximilian Spencer, Om Patel, Joon Yong Doh, Eid Rodan, Kevin Zhu, and Sean O’Brien. 2025. Improving LLM abilities in idiomatic translation. In *Future of Information and Communication Conference*, pages 361–375. Springer.

Zhengxiao Du, Aohan Zeng, Yuxiao Dong, and Jie Tang. 2024. Understanding emergent abilities of language models from the loss perspective. *Preprint*, arXiv:2403.15796.

Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, and 1 others. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.

Bowen Fang, Zixiao Yang, Shukai Wang, and Xuan Di. 2024. Travellm: Could you plan my new public transit route in face of a network disruption? *arXiv preprint arXiv:2407.14926*.

Mohamed Amine Ferrag, Norbert Tihanyi, and Merouane Debbah. 2025. From llm reasoning to autonomous ai agents: A comprehensive review. *arXiv preprint arXiv:2504.19678*.

GoogleMaps. 2025. [Philadelphia city map](#).

Atharva Gundawar, Mudit Verma, Lin Guan, Karthik Valmееkam, Siddhant Bhambri, and Subbarao Kambhampati. 2024. Robust planning with LLM-modulo framework: Case study in travel planning. *arXiv preprint arXiv:2405.20625*.

Meng-Hao Guo, Jiajun Xu, Yi Zhang, Jiayi Song, Haoyang Peng, Yi-Xuan Deng, Xinzhi Dong, Kiyohiro Nakayama, Zhengyang Geng, Chen Wang, and 1 others. 2025. R-bench: Graduate-level multidisciplinary benchmarks for llm & mllm complex reasoning evaluation. *arXiv preprint arXiv:2505.02018*.

Yilun Hao, Yongchao Chen, Yang Zhang, and Chuchu Fan. 2024. Large language models can plan your travels rigorously with formal verification tools. *arXiv preprint arXiv:2404.11891*.

Annette Herskovits. 1986. *Language and spatial cognition*. Cambridge university press Cambridge.

Karla L Hoffman, Manfred Padberg, Giovanni Rinaldi, and 1 others. 2013. Traveling salesman problem. *Encyclopedia of operations research and management science*, 1:1573–1578.

Zijin Hong, Zheng Yuan, Qinggang Zhang, Hao Chen, Junnan Dong, Feiran Huang, and Xiao Huang. 2025. Next-generation database interfaces: A survey of llm-based text-to-sql. *IEEE Transactions on Knowledge and Data Engineering*.

Hanxu Hu, Hongyuan Lu, Huajian Zhang, Yun-Ze Song, Wai Lam, and Yue Zhang. 2024. Chain-of-symbol prompting for spatial reasoning in large language models. In *First Conference on Language Modeling*.

Lei Huang, Weijiang Yu, Weitao Ma, Weihong Zhong, Zhangyin Feng, Haotian Wang, Qianglong Chen, Weihua Peng, Xiaocheng Feng, Bing Qin, and 1 others. 2023. A survey on hallucination in large language models: Principles, taxonomy, challenges, and open questions. *ACM Transactions on Information Systems*.

Aaron Hurst, Adam Lerer, Adam P Goucher, Adam Perelman, Aditya Ramesh, Aidan Clark, AJ Ostrow, Akila Welihinda, Alan Hayes, Alec Radford, and 1 others. 2024. Gpt-4o system card. *arXiv preprint arXiv:2410.21276*.

Da Ju, Song Jiang, Andrew Cohen, Aaron Foss, Sasha Mitts, Arman Zharmagambetov, Brandon Amos, Xian Li, Justine T Kao, Maryam Fazel-Zarandi, and 1

693	others. 2024. To the globe (TTG): Towards language-driven guaranteed travel planning. <i>arXiv preprint arXiv:2410.16456</i> .	
694		
695		
696	Subbarao Kambhampati, Karthik Valmeekam, Lin Guan, Mudit Verma, Kaya Stechly, Siddhant Bhambri, Lucas Saldyt, and Anil Murthy. 2024. LLMs can't plan, but can help planning in LLM-modulo frameworks. <i>arXiv preprint arXiv:2402.01817</i> .	
697		
698		
699		
700		
701	Xin Lai, Zhuotao Tian, Yukang Chen, Yanwei Li, Yuhui Yuan, Shu Liu, and Jiaya Jia. 2024. Lisa: Reasoning segmentation via large language model. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9579–9589.	
702		
703		
704		
705		
706	Stephen C Levinson. 2003. <i>Space in language and cognition: Explorations in cognitive diversity</i> , volume 5. Cambridge University Press.	
707		
708		
709	Chengzu Li, Wenshan Wu, Huanyu Zhang, Yan Xia, Shaoguang Mao, Li Dong, Ivan Vulić, and Furu Wei. 2025. Imagine while reasoning in space: Multimodal visualization-of-thought. <i>arXiv preprint arXiv:2501.07542</i> .	
710		
711		
712		
713		
714	Liang Li, Yuhang Li, Wenhao Zhu, Shicheng Liu, Tianchi Huo, Zhiyang Teng, Yimo Zhang, Jingyi Li, Qian Guo, Huayang Li, Xueyi Wang, Yifan Jiang, Liwei Huang, Furu Wang, Yan Sun, Yixuan Wang, Xiang Yue, Chao Zhou, Kun Qiu, and 17 others. 2024. Livebench: A challenging, contamination-free llm benchmark . <i>Preprint</i> , arXiv:2406.19314.	
715		
716		
717		
718		
719		
720		
721	Yang Liu, Fanyou Wu, Zhiyuan Liu, Kai Wang, Feiyue Wang, and Xiaobo Qu. 2023. Can language models be used for real-world urban-delivery route optimization? <i>The Innovation</i> , 4(6).	
722		
723		
724		
725	Wufei Ma, Luoxin Ye, Celso M de Melo, Alan Yuille, and Jieneng Chen. 2025. Spatialllm: A compound 3d-informed design towards spatially-intelligent large multimodal models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)</i> , pages 17249–17260.	
726		
727		
728		
729		
730		
731		
732	Mistral. 2024. Mistral-large-2411 .	
733	Aske Plaat, Annie Wong, Suzan Verberne, Joost Broekens, Niki van Stein, and Thomas Back. 2024. Reasoning with large language models, a survey. <i>arXiv preprint arXiv:2407.11511</i> .	
734		
735		
736		
737	Thad Anderson Polk. 1992. <i>Verbal reasoning</i> . Carnegie Mellon University.	
738		
739	Jingqing Ruan, Yihong Chen, Bin Zhang, Zhiwei Xu, Tianpeng Bao, Hangyu Mao, Ziyue Li, Xingyu Zeng, Rui Zhao, and 1 others. 2023. Tptu: Task planning and tool usage of large language model-based ai agents. In <i>NeurIPS 2023 Foundation Models for Decision Making Workshop</i> .	
740		
741		
742		
743		
744		
	Pratyusha Sharma, Tamar Rott Shaham, Manel Baradad, Stephanie Fu, Adrian Rodriguez-Munoz, Shivam Duggal, Phillip Isola, and Antonio Torralba. 2024. A vision check-up for language models. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 14410–14419.	745 746 747 748 749 750
	Yihong Tang, Ao Qu, Zhaokai Wang, Dingyi Zhuang, Zhaofeng Wu, Wei Ma, Shenhao Wang, Yunhan Zheng, Zhan Zhao, and Jinhua Zhao. 2025. Sparkle: Mastering basic spatial capabilities in vision-language models elicits generalization to composite spatial reasoning. In <i>Findings of the Association for Computational Linguistics: EMNLP 2025</i> .	751 752 753 754 755 756 757
	Yihong Tang, Zhaokai Wang, Ao Qu, Yihao Yan, Kebing Hou, Dingyi Zhuang, Xiaotong Guo, Jinhua Zhao, Zhan Zhao, and Wei Ma. 2024. Synergizing spatial optimization with large language models for open-domain urban itinerary planning. <i>arXiv preprint arXiv:2402.07204</i> .	758 759 760 761 762 763
	Gemini Team, Petko Georgiev, Ving Ian Lei, Ryan Burnell, Libin Bai, Anmol Gulati, Garrett Tanzer, Damien Vincent, Zhufeng Pan, Shibo Wang, and 1 others. 2024. Gemini 1.5: Unlocking multimodal understanding across millions of tokens of context. <i>arXiv preprint arXiv:2403.05530</i> .	764 765 766 767 768 769
	Moshe Unger, Alexander Tuzhilin, and Michel Wedel. 2025. STARE: Predicting decision making based on spatio-temporal eye movements. <i>arXiv preprint arXiv:2508.04148</i> .	770 771 772 773
	Karthik Valmeekam, Matthew Marquez, Alberto Olmo, Sarath Sreedharan, and Subbarao Kambhampati. 2024. Planbench: An extensible benchmark for evaluating large language models on planning and reasoning about change. <i>Advances in Neural Information Processing Systems</i> , 36.	774 775 776 777 778 779
	Jiayu Wang, Yifei Ming, Zhenmei Shi, Vibhav Vineet, Xin Wang, Sharon Li, and Neel Joshi. 2024a. Is a picture worth a thousand words? delving into spatial reasoning for vision language models. <i>Advances in Neural Information Processing Systems</i> , 37:75392–75421.	780 781 782 783 784 785
	Kaimin Wang, Yuanzhe Shen, Changze Lv, Xiaoqing Zheng, and Xuan-Jing Huang. 2025. Triptailor: A real-world benchmark for personalized travel planning. In <i>Findings of the Association for Computational Linguistics: ACL 2025</i> , pages 9705–9723.	786 787 788 789 790
	Yiming Wang, Zhijiang Ma, Yixin Huang, Wenhao Chen, Jiahui Liu, Shuhe Yu, Jiakai He, Zhuang Liu, Yimeng Ren, Zhi Feng, Tianyu Liu, Jun Qiu, Hao Huang, Ming Zeng, Kaijiang Wu, Daisong Wei, Qi Xiong, Jiezhong Chen, Chao Liu, and 5 others. 2024b. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark . <i>Preprint</i> , arXiv:2406.01574.	791 792 793 794 795 796 797 798
	Yubo Wang, Xueguang Ma, Ge Zhang, Yuansheng Ni, Abhranil Chandra, Shiguang Guo, Weiming Ren, Aaran Arulraj, Xuan He, Ziyang Jiang, and 1 others.	799 800 801

802	2024c. Mmlu-pro: A more robust and challenging multi-task language understanding benchmark. <i>arXiv preprint arXiv:2406.01574</i> .	Lingfeng Zhang, Yuening Wang, Hongjian Gu, Atia Hamidizadeh, Zhanguang Zhang, Yuecheng Liu, Yutong Wang, David Gamaliel Arcos Bravo, Junyi Dong, Shunbo Zhou, and 1 others. 2024. Et-plan-bench: Embodied task-level planning benchmark towards spatial-temporal cognition with foundation models. <i>arXiv preprint arXiv:2410.14682</i> .	857
803			858
804			859
805	Walter Whiteley, Nathalie Sinclair, and Brent Davis. 2015. What is spatial reasoning? In <i>Spatial reasoning in the early years</i> , pages 3–14. Routledge.		860
806			861
807			862
808	Qiucheng Wu, Handong Zhao, Michael Saxon, Trung Bui, William Yang Wang, Yang Zhang, and Shiyu Chang. 2024a. Vsp: Assessing the dual challenges of perception and reasoning in spatial planning tasks for vlms. <i>arXiv preprint arXiv:2407.01863</i> .	Zirui Zhao, Wee Sun Lee, and David Hsu. 2024. Large language models as commonsense knowledge for large-scale task planning. <i>Advances in Neural Information Processing Systems</i> , 36.	864
809			865
810			866
811			867
812			
813	Wenshan Wu, Shaoguang Mao, Yadong Zhang, Yan Xia, Li Dong, Lei Cui, and Furu Wei. 2024b. Mind’s eye of LLMs: Visualization-of-thought elicits spatial reasoning in large language models. In <i>The Thirty-eighth Annual Conference on Neural Information Processing Systems</i> .	Chenming Zhu, Tai Wang, Wenwei Zhang, Jiangmiao Pang, and Xihui Liu. 2024. Llava-3d: A simple yet effective pathway to empowering llms with 3d-awareness. <i>arXiv preprint arXiv:2409.18125</i> .	868
814			869
815			870
816			871
817			
818			
819	Chengxing Xie and Difan Zou. 2024. A human-like reasoning framework for multi-phases planning task with large language models. <i>arXiv preprint arXiv:2405.18208</i> .		
820			
821			
822			
823	Jian Xie, Kai Zhang, Jiangjie Chen, Tinghui Zhu, Renze Lou, Yuandong Tian, Yanghua Xiao, and Yu Su. 2024. Travelplanner: A benchmark for real-world planning with language agents. <i>arXiv preprint arXiv:2402.01622</i> .		
824			
825			
826			
827			
828	Runsen Xu, Xiaolong Wang, Tai Wang, Yilun Chen, Jiangmiao Pang, and Dahua Lin. 2024. Pointllm: Empowering large language models to understand point clouds. In <i>European Conference on Computer Vision</i> , pages 131–147. Springer.		
829			
830			
831			
832			
833	Yi Xu, Chengzu Li, Han Zhou, Xingchen Wan, Caiqi Zhang, Anna Korhonen, and Ivan Vulić. 2025. Visual planning: Let’s think only with images. <i>arXiv preprint arXiv:2505.11409</i> .		
834			
835			
836			
837	Jihan Yang, Shusheng Yang, Anjali W Gupta, Rilyn Han, Li Fei-Fei, and Saining Xie. 2025. Thinking in space: How multimodal large language models see, remember, and recall spaces. In <i>Proceedings of the Computer Vision and Pattern Recognition Conference</i> , pages 10632–10643.		
838			
839			
840			
841			
842			
843	Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. 2022. React: Synergizing reasoning and acting in language models. <i>arXiv preprint arXiv:2210.03629</i> .		
844			
845			
846			
847	Yelp. 2024a. Yelp fusion API .		
848	Yelp. 2024b. Yelp open dataset .		
849	Xiang Yue, Yuansheng Ni, Kai Zhang, Tianyu Zheng, Ruoqi Liu, Ge Zhang, Samuel Stevens, Dongfu Jiang, Weiming Ren, Yuxuan Sun, and 1 others. 2024. Mmmu: A massive multi-discipline multimodal understanding and reasoning benchmark for expert agi. In <i>Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pages 9556–9567.		
850			
851			
852			
853			
854			
855			
856			

872	A Usage of LLM		
873	The LLMs are only used to aid and polish writing	vided only contains the rating not the raw review	920
874	in this paper.	manuscript. The dataset is in English fully.	921
		The price attributes are obtained from Yelp Fusion API	922
		(Yelp, 2024a).	923
875	B Case Study	Hotels We extract businesses with the "Hotels"	924
		category only from the original data since the cate-	925
876	In Figure B.1, we visualize an itinerary generated	gory "Hotel & Travel" actually refers to airports or	926
877	by GPT-4o for the spatial reasoning task in the	train stations.	927
878	tool-use mode. The green circles in Figure B.1c	Restaurants We extract businesses with	928
879	are drawn based on the clustering information pro-	"Restaurants" or "Food" categories from the origi-	929
880	vided in Figure B.1b, with Day 1's route high-	nal data. We keep the top 500 restaurants with the	930
881	lighted in red. GPT-4o successfully understands	most reviews for efficiency and cost management.	931
882	the spatial relationship among attractions and ar-	Attractions We extract businesses with "Muse-	932
883	ranges attractions within the same cluster for the	ums", "Parks", "Local Flavor", "Zoos", "Tours",	933
884	same day. However, when comparing Figure B.1c	"Landmarks & Historical Buildings", and "Sou-	934
885	(the proposed route) with Figure B.1d (the opti-	venir Shops" categories from the original dataset.	935
886	mized route), we observe that the fourth attrac-	User Reviews We keep reviews with a "useful"	936
887	tion choice is unreasonable. It is drawn from a cluster	rating greater or equal to 1 from the original dataset.	937
888	far from other clusters assigned to Day 1 and much	The pre-process can help filter out less informative	938
889	closer to the clusters selected for Day 2.	reviews and control the file size.	939
890	From a spatial reasoning perspective, arranging	All businesses labeled as unopened are filtered,	940
891	one day in an itinerary involves two major tasks: 1.	and further standard data cleaning is performed.	941
892	Identifying spatial relationships among attractions	See Table C.1 for the attributes we keep for each	942
893	to form the first level of clusters. 2. Determining	category for the base data.	943
894	spatial relationships among these first-level clus-		
895	ters to organize the route for the entire day. Both	C.2 Spatial Cluster Information Generation	944
896	tasks have a similar goal: to find the objects' spa-	Spatial cluster information is calculated and pre-	945
897	tial relationship and distance in a two-dimensional	sented to LLMs in two ways.	946
898	space. They also rely on the same spatial reasoning	Filtered Data In this task, business data for each	947
899	abilities: understanding spatial relationships and	plan is already filtered based on the preference re-	948
900	distance proximity.	quested from the query. The k-mean clustering	949
901	GPT-4o fails in the second task in this case study:	method is deployed to get the candidates' spatial	950
902	it does not understand that the first-level cluster 13	clustering information. The cluster number we	951
903	is far away from clusters 17 and 8 arranged for Day	choose is the integer value after using the candi-	952
904	1. This mistake indicates that LLMs' spatial rea-	date's number divided by 5.	953
905	soning abilities are still inadequate for visualizing	Tool use In tool use tasks, LLMs will gather	954
906	and understanding spatial relationships between	candidates' business through their tool calling, and	955
907	objects. This observation concludes that GPT-4o's	then spatial clustering information based on those	956
908	success in the first task, i.e., forming precise first-	candidates will be provided when they successfully	957
909	level cluster choices among attractions, primarily	call the clustering function. The clustering strategy	958
910	relies on their verbal reasoning ability to use the	is the same as the filtered data tasks.	959
911	clustering information directly rather than actual		
912	reasoning through the space.	C.3 Reproducibility	960
913	C Dataset Details	We attach the code used to generate the results and	961
914	C.1 Data Pre-process Strategy	conduct evaluation in the supplementary material.	962
915	The base dataset is collected from Yelp Dataset	We plan to open-source the code and dataset in	963
916	(Yelp, 2024b). The dataset is a subset of the Yelp	the future. The data crafting pipeline is detailed in	964
917	Dataset that intended for educational use. The origi-	Section 3.1 and Appendix C.	965
918	nal dataset contains the user reviews that are hard		
919	to monitor in full scale, thus the final dataset pro-		

<p>Travel Plan</p> <p>Day 1:</p> <p>**Accommodation**: The Rittenhouse Hotel</p> <p>**Morning Attraction**: Please Touch Museum</p> <p>**Afternoon Attractions**:</p> <p>-Shofuso Japanese House & Garden -Philadelphia Zoo</p> <p>**Night Attraction**: Spruce Street Harbor Park</p> <p>Day 2:</p> <p>**Accommodation**: The Rittenhouse Hotel</p> <p>**Morning Attraction**: Franklin Square</p> <p>**Afternoon Attractions**:</p> <p>-Museum of the American Revolution -The Franklin Institute Science Museum</p> <p>**Night Attraction**: Rittenhouse Square Park</p>	<p>Cluster Information</p> <p>...</p> <p>Cluster_0: ...Franklin Square, The Betsy Ross House, Museum of the American Revolution, National Constitution Center, National Liberty Museum, National Museum of American Jewish History, Independence Park Hotel, Benjamin Franklin Museum...</p> <p>...</p> <p>Cluster_4: ...Kimpton Hotel Palomar Philadelphia, Wonderspaces Philadelphia, Dilworth Park, Rittenhouse Square Park, The Rittenhouse Hotel, Reading Terminal Market, Pennsylvania Academy of Fine Arts, ...</p> <p>...</p> <p>Cluster_8: Philadelphia Zoo</p> <p>...</p> <p>Cluster_13: ...Independence Seaport Museum, Benjamin Franklin Bridge, Penn's Landing, Spruce Street Harbor Park, ...</p> <p>...</p> <p>Cluster_16: ...The Franklin Institute Science Museum, Fairmount Water Works, Cira Green, Philly Bike Tour, Eastern State Penitentiary Historic Site, Philadelphia Museum of Art...</p> <p>Cluster_17: Please Touch Museum, Shofuso Japanese House & Garden</p> <p>...</p>
--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------	--------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------

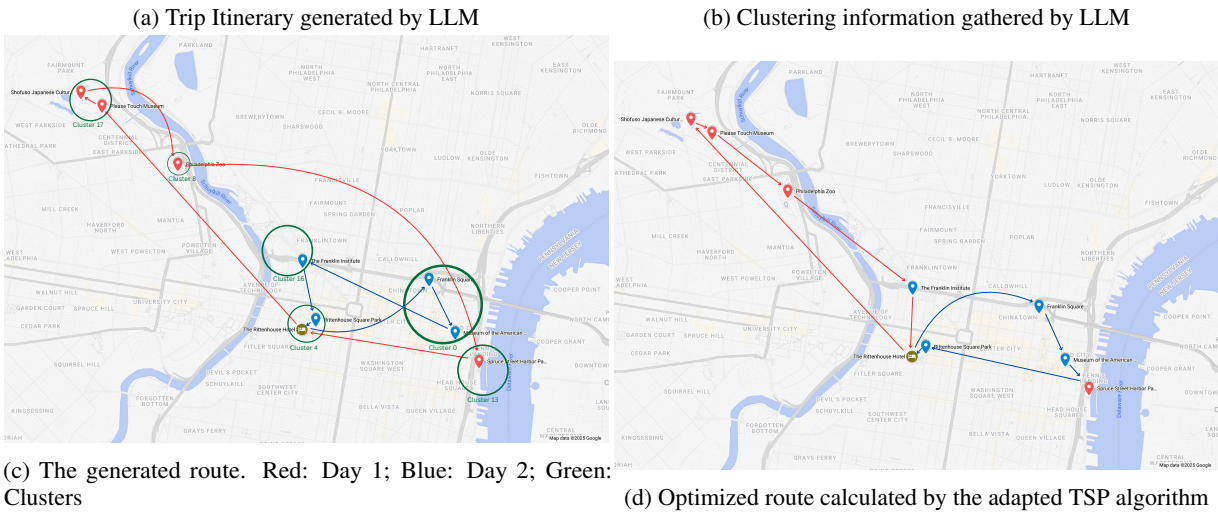


Figure B.1: Visualization of a case study about the itinerary generated by GPT 4o in tool-use mode, drawn in Google Map (GoogleMaps, 2025). Plan-wise (red markers and routes in Figure B.1c), one of the main issues is the route leads to the bottom right corner of the map (Cluster 13) while visiting the same area (Cluster 0) on the second day again. The mistake is corrected in the optimized route in Figure B.1d. For this itinerary, the total distance gap ratio is 25.6%. Additionally, the extra cluster jump ratio is 100%.

D More experiments and results

D.1 Greedy Algorithm

- A* Algorithm Employs a Minimum Spanning Tree as the heuristic $h(n)$.
- Mixed Integer Algorithm Uses continuous ordering variables in the Miller–Tucker–Zemlin (MTZ) subtour elimination constraints.

The near-perfect performance of algorithm-based solutions aligns with findings in combinatorial optimization research. These alternative baselines strongly indicate that LLMs’ reasoning and planning abilities still require improvement. Notably, these performances depend on extensive preparation, task-specific code, and prior human knowledge where LLMs can facilitate end-to-end reasoning solutions.

D.2 Human Evaluation

D.2.1 Evaluation Design

For evaluation We recruit 10 PhD students volunteers and 2 experienced travelers volunteers to participate in two evaluation tasks, yielding 120 data points for each evaluation task. The example questionnaire can be find below.

Evaluation 1: Participants review 10 sets of four valid plans generated by different models in response to the same query, within the Filtered Data Route OP task. Each plan has natural language text and a route visualization. Evaluators are asked to choose their preferred plan in each set.

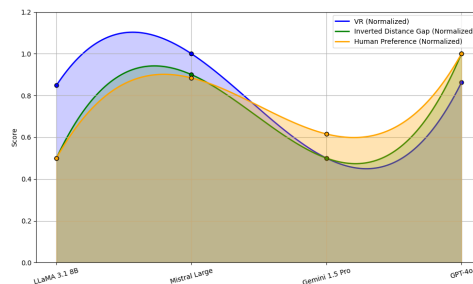
Evaluation 2: Participants compare 10 sets of four plans generated by GPT-4o, each corresponding to a different task but based on the same user query. Like in Evaluation 1, each plan includes text and a route visualization, and evaluators select their preferred option from each set.

Table C.1: The attributes for basic data

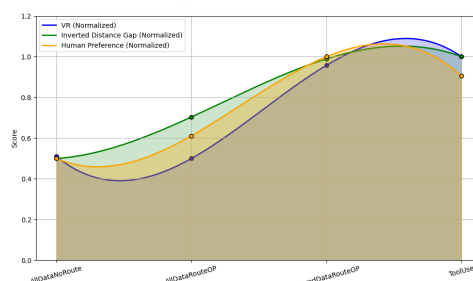
Preference	Attributes
Hotels	"business_id", "name", "address", "latitude", "longitude", "stars", "price"
Restaurants	"business_id", "name", "address", "latitude", "longitude", "stars", "price", "good_for_meal", "cuisine_1", "cuisine_2"
Attraction	"business_id", "name", "address", "latitude", "longitude", "stars", "price"

Table C.2: Dataset statistic of Santa Barbara

	Base	Avg Reviews	Review Attributes
Restaurants	200	166	cuisine, flavor, freshness, service, environment, value
Hotels	50	92	quality, location, service, safety
Attractions	100	26	family, history, activity, nature, food, shopping



(a) Evaluation 1



(b) Evaluation 2

D.2.2 Template

Here is the example question in the human evaluation:

Among the four plan, which one do you prefer?

[Plan A], [Plan B], [Plan C], [Plan D]

[Visualization A], [Visualization B], [Visualization C], [Visualization D]

D.2.3 Recruitment details

The task lasts for about 15 minutes and all volunteers agrees with no payments. All volunteers are currently in USA.

D.2.4 Human Evaluation Results

The preference rate from two human evaluation tasks aligns with the results calculated through the evaluation designed in ItinBench. The preference rate matches the performance across the models in the same task in both validated rate (VR) and plan-wise distance gap (Total-DG) (Figure D.2a). Furthermore, the preference rate aligns with OpenAI-o1's performance in multiple tasks in VR and Total-DG (Figure D.2b). A detailed quantitative results is recorded in Section D.2

D.3 Santa Barbara

We additionally construct a smaller subset of the Yelp dataset restricted to businesses in Santa Barbara. Because ItinBench is intended to measure how model performance varies across task formulations (rather than to maximize absolute performance on a single setting), this subset serves as a

Figure D.2: Alignments between human evaluation preferences (yellow), validate rate (blue), and total distance gap (green). For better visualization, Total-DG is inverted to see the trend.

controlled check that the observed cross-task performance fluctuations persist under a different city demographic. Table C.2 summarizes the subset statistics. We run the same four tasks using LLaMA 3.1 8B, Gemini 1.5 Pro, and OpenAI o1. As shown in Table C.3, the resulting performance exhibits trends consistent with the main results obtained on itineraries collected in Philadelphia. This suggests that demographic factors—such as attraction density—do not materially affect the qualitative observations reported in our study.

E Prompts

E.1 Review Extraction Prompt

There is one prompt for each of the business categories. The prompt asks LLM to extract ratings or measurements on different scales. These numbers are processed into phrases, e.g., rating 5 for location is "excellent location," for planner LLM

Table C.3: Results for Tasks 1 to 4 (VR and Total DG) for Santa Barbara subset.

	Task 1		Task 2		Task 3		Task 4	
	VR ↑	Total DG ↓	VR ↑	Total DG ↓	VR ↑	Total DG ↓	VR ↑	Total DG ↓
LLaMA 3.1 8B	0	114.1	0	107.5	43	110.3	28	125.0
Gemini 1.5 pro	7	125.5	4	-	13	-	17	-
OpenAI o1	15	130.5	5	45.7	42	6.5	41.5	7.0

Table D.4: Alternative baselines; columns correspond to Table 4 in the main text. Except for the greedy + min distance approach, all other approaches achieve perfect performance in evaluation.

	Verbal Reasoning					Spatial Reasoning			
	OOP ↓	MI ↓	Micro ↑	Macro ↑	VR ↑	ARG ↓	DG ↓	Total-DG ↓	ECJ ↓
Greedy Approach (#100)									
A*	0.0	0.0	100.0	100.0	100.0	0 (4.00)	0.0	0.0	0.0
MIP	0.0	0.0	100.0	100.0	100.0	0 (4.00)	0.0	0.0	0.0

Table D.5: Preferred rate across the plans generated by four different LLMs. Evaluation 1 refers to the comparison between the plans generated by four LLMs based on the same user query. Evaluation 2 refers to the comparison between a set of plans generated all by OpenAI-o1 with same user query, but within different task setting.

	Mistral	GPT4o	Llama	Gemini
Evaluation 1	30.8	35.8	14.2	19.2

to better understand.

Here is the prompt for hotel review extraction.

You are an assistant designed to summarize reviews of businesses for travel planning purposes. Your goal is to provide **faithful, concise, and relevant information** based on the following reviews compiled into the txt file. Follow these principles:

- Focus on Travel-Relevant Details:** Prioritize aspects like location convenience, proximity to landmarks, transportation options, ambiance, cleanliness, service quality, amenities, and overall reliability.
- Avoid Bias:** Reflect the consensus of reviews, clearly noting if opinions are mixed. Do not add, fabricate, or exaggerate details.
- Clarify Nuances:** Mention trends (e.g., "frequent mentions of slow service" or "consistent praise for central location").
- Respect Context:** Differentiate be-

tween subjective opinions (e.g., "some reviewers found the rooms small") and factual details (e.g., "located 5 minutes from the train station").

5. Stay Honest: If the reviews are unclear or contradictory, state this explicitly rather than drawing unsupported conclusions.

6. Highlight Red Flags or Unique Strengths: Identify issues (e.g., safety concerns, unexpected fees) or advantages (e.g., exceptional customer service, standout features).

Output formatting instructions:

On a scale of 1 to 5. 3 means average, 4 means good, 5 means excellent, 2 means below average, and 1 means bad. Be faithful and give objective ratings.

- Evaluate Room Quality on a scale from 1 to 5. Considering size, cleanliness, space, amenities, noise level, and other considerations.
- Evaluate the location and convenience on a scale from 1 to 5. Consider transportation options, proximity to attractions, and other factors.
- Evaluate the hotel's service on a scale from 1 to 5, considering the cleaning service, customer service, valet service, check-in and check-out experience, and interactions between travelers and the hotel staff in general.
- Evaluate the safety on a scale from 1 to 5. Considering the surrounding area

1073
1074
1075
1076
1077
1078
1079
1080
1081
1082
1083
1084
1085
1086
1087
1088
1089
1090
1091
1092
1093
1094
1095
1096
1097
1098
1099
1100
1101
1102
1103
1104
1105
1106

Table D.6: Preferred rate across the plans generated by four different LLMs. Evaluation 1 refers to the comparison between the plans generated by four LLMs based on the same user query. Evaluation 2 refers to the comparison between a set of plans generated all by OpenAI-o1 with same user query, but within different task setting.

	AllDataNoRoute	AllDataRouteOP	FilteredData	ToolUse
Evaluation 2	20.5	19.8	33.2	30.5

1107	traffic, safety in the hotel, and other factors that influence the safety concern if possible.		poorly maintained and have unpleasant odors.	1148
1108				1149
1109				
1110	Give one evaluation for each attribute and followed by a sentence of reasoning.		The hotel has a rating of 3 for location.	1150
1111			The hotel is conveniently located near the airport, but guests noted that the surrounding area lacks amenities and attractions, requiring a drive for most necessities.	1151
1112	— Example 1 Starts —			1152
1113	The hotel has a rating of 4 for quality.			1153
1114	Rooms are beautifully appointed with stunning views, luxurious amenities, and impeccable cleanliness. Guests appreciate the spaciousness and comfort of the beds, although some mention the rooms being on the smaller side typical for city hotels.			1154
1115				1155
1116			The hotel has a rating of 2 for average service. Service quality is inconsistent, with many guests reporting rude or unhelpful staff. Issues with check-in, maintenance, and customer service have been frequently mentioned.	1156
1117				1157
1118				1158
1119				1159
1120				1160
1121	The hotel has a rating of 5 for location.		The hotel has a rating of 2 for average safety. Concerns about safety have been raised, particularly regarding the external room entrances and reports of security issues. Some guests felt uncomfortable due to the behavior of staff and security.	1161
1122	Located in the Comcast Center, the hotel offers breathtaking views of Philadelphia and is conveniently situated near major attractions. The elevator ride to the 60th floor lobby is a highlight.			1162
1123				1163
1124				1164
1125				1165
1126				1166
1127	The hotel has a rating of 4 for service.		— Example 2 Ends —	1167
1128	Service is generally exceptional, with staff going above and beyond to make guests feel welcome. However, there are mixed reviews regarding the handling of certain situations, particularly in the bar area and restaurant.			1168
1129			Given reviews: {reviews}	1169
1130			Your evaluation:	1170
1131				
1132			Here is the attraction review extraction prompt:	1171
1133				
1134	The hotel has a rating of 4 for safety. The hotel is located in a prominent area of Philadelphia, and while most reviews do not raise safety concerns, there are mentions of discriminatory treatment that could affect the perception of safety for some guests.		You are an assistant designed to analyze and summarize reviews of attractions for travel planning purposes. Your goal is to deliver faithful, concise, and travel-relevant insights based on the reviews provided in the attached text file. Follow these principles:	1172
1135				1173
1136				1174
1137				1175
1138				1176
1139				1177
1140				1178
1141	— Example 1 Ends —		1. Focus on Key Travel-Relevant Features: Highlight details such as the attraction’s location, accessibility, proximity to key landmarks, transportation options, and overall convenience for visitors. Address aspects like ambiance, cleanliness, crowd levels, staff behavior, unique offerings, and amenities.	1179
1142	— Example 2 Starts —			1180
1143	The hotel has a rating of 2 for quality.			1181
1144	Rooms are often reported as dirty, with issues like stained bedding, bugs, and unclean bathrooms. Some guests noted that while the rooms are spacious, they are			1182
1145				1183
1146				1184
1147				1185
				1186

1187	2. Reflect Consensus and Avoid Bias:	is needed for this attraction. Hiking or	1237
1188	Summarize the general sentiment of re-	dangerous activities would be a strong	1238
1189	viewers, noting both strengths and short-	activity level 3, visiting a outdoor park	1239
1190	comings as expressed. Avoid exaggera-	could be a medium level 2, and visiting a	1240
1191	tion or unfounded interpretations. Indi-	museum could be a low activity level 1.	1241
1192	cate if opinions vary significantly among		
1193	reviewers. Clarify Trends and Nuances:	4. Measure the natural scene from 0 to 3.	1242
1194	3. Identify recurring themes (e.g., "many	This measures how much the attraction	1243
1195	reviewers appreciated the tranquil set-	accesses nature and sightseeing views. 0	1244
1196	ting" or "frequent complaints about high	means completely indoor, and 3 means	1245
1197	entrance fees"). Distinguish between	outdoor with the natural scene.	1246
1198	subjective opinions (e.g., "some visitors		
1199	found it too crowded") and objective	5. Measure how food-oriented is the at-	1247
1200	facts (e.g., "located 10 minutes from the	traction. Level 3 would be food oriented	1248
1201	nearest metro station"). Acknowledge	attraction. 0 indicates this attraction has	1249
1202	Uncertainty or Contradictions:	no relation to food.	1250
1203	4. If reviews are unclear or contradictory,		
1204	explicitly state this rather than making	6. Measure if attraction focus on shop-	1251
1205	unsupported conclusions.	ping. A market would be level 3, a his-	1252
1206		torical landmark could be 0 since it's for	1253
1207	5. Highlight Red Flags or Unique Fea-	visiting only.	1254
1208	tures: Draw attention to notable is-		
1209	ssues (e.g., safety concerns, hidden costs)	Here are some examples	1255
1210	or standout positives (e.g., spectacular	— Example 1 starts —	1256
1211	views, interactive exhibits).	This place has a family oriented level	1257
1212	Output formatting instructions: All the	3. Many families enjoyed the carriage	1258
1213	evaluation is on a scale of 0 to 3, 0 means	rides, with children actively participating	1259
1214	not applicable, 1 means low tendency, 2	and asking questions. The experience	1260
1215	means medium, and 3 means strong ten-	was highlighted as a memorable family	1261
1216	dency. The scale is not a score but a mea-	activity.	1262
1217	surement. There is no implication that a		
1218	better score leads to a better business.	This place has a history oriented level	1263
1219	1. Measure the family orientation from	3. The carriage rides provide informa-	1264
1220	0 to 3. Factors include kids involvement,	tive tours of historical areas, with knowl-	1265
1221	and What kinds of activities are orga-	edgeable guides sharing insights about	1266
1222	nized? 0 means not for family, 1 means	Philadelphia's history and architecture.	1267
1223	really small family factor is designed, 2		
1224	means an average amount of family ac-	This place has a activity oriented level	1268
1225	tivities, and 3 means this place designed	1. The activity level is low as the rides	1269
1226	for family.	are leisurely and do not require physical	1270
1227	2. Measure the history oritentaion from	exertion from participants.	1271
1228	0 to 3. Factors include history, culture, ed-		
1229	ucation, and other considerations around	This place has a nature oriented level 1.	1272
1230	history and culture. 0 means no history	The rides are primarily through urban ar-	1273
1231	consideration from this site, 1 means not	reas with limited access to natural scen-	1274
1232	designed for history exploration, 2 mean	ery, focusing more on the city's histor-	1275
1233	average amount of history attributes, 3	ical aspects.	1276
1234	means this place has a lot of history fac-		
1235	tor included.	This place has a food oriented level 0.	1277
1236	3. Measure the activity level from 0 to	The attraction does not have a food-	1278
	3. This measures what level of action	related focus.	1279
		This place has a shopping oriented level	1280
		0. The carriage rides are not related to	1281
		shopping; they are purely a sightseeing	1282
		experience.	1283
		— Example 1 Ends —	1284

1285	— Example 2 Starts —	faithful, concise, and relevant information	1332
1286	This place has a family oriented level	based on the following reviews compiled into the txt file. Follow these principles:	1333
1287	3. Spruce Street Harbor Park is highly		1334
1288	family-friendly, featuring activities for		1335
1289	children such as oversized games, an arcade,	1. Focus on Travel-Relevant Details: Prioritize aspects crucial to travelers, such as food quality, location convenience (proximity to landmarks and transportation options), ambiance, cleanliness, service quality, amenities, and overall reliability.	1336
1290	and play areas. Many reviewers		1337
1291	noted the park’s appeal to families, with	2. Avoid Bias: Provide balanced evaluations that reflect the consensus of available reviews. Clearly indicate when opinions are mixed, and refrain from fabricating, exaggerating, or omitting key details.	1338
1292	fun events and games for kids.		1339
1293	This place has a history oriented level 1.	3. Clarify Nuances: Highlight notable trends in feedback (e.g., "frequent mentions of slow service" or "consistent praise for convenient location") to provide an accurate overview.	1340
1294	While the park is located near historical		1341
1295	sites, it does not focus on history or cultural	4. Respect Context: Differentiate between subjective opinions (e.g., "some diners found the portions small") and factual details (e.g., "located within walking distance of a major metro station").	1342
1296	education. The attraction is more		1343
1297	about leisure and entertainment rather	5. Maintain Honesty: If reviews are unclear, contradictory, or lacking sufficient detail, explicitly state this instead of making unsupported conclusions.	1344
1298	than historical significance.		1345
1299	This place has an activity oriented level 2.	6. Highlight Red Flags and Unique Strengths: Identify significant issues (e.g., long wait times, poor hygiene, safety concerns) and standout features (e.g., exceptional cuisine, distinctive ambiance, or unique menu options).	1346
1300	The park offers various activities such as	Output formatting instructions:	1347
1301	hammocks, games like giant Jenga and	The rating is from 1 to 5, higher the better. 3 is average. 4 and 5 means good and excellent. 2 means below average, 1 means bad. Be faithful to the review’s statement and give a rating accordingly from 1 to 5.	1348
1302	Connect Four, and paddle boat rentals.		1349
1303	However, the level of physical activity is	1. Evaluate the flavor of the dishes on a scale of 1 to 5.	1350
1304	moderate, making it suitable for casual		1351
1305	visitors.	2. Evaluate the freshness of the food on a scale of 1 to 5.	1352
1306	This place has a nature oriented level 2.		1353
1307	The park is situated along the Delaware		1354
1308	River and features hammocks and seating		1355
1309	areas with views of the water. However,		1356
1310	it is primarily an urban park with		1357
1311	limited natural scenery.		1358
1312	This place has a food oriented level 3.		1359
1313	There is a strong focus on food, with		1360
1314	numerous food trucks and vendors offering		1361
1315	a variety of options, including local		1362
1316	favorites. Reviewers praised the food		1363
1317	offerings, although some noted that prices		1364
1318	can be high.		1365
1319	This place has a shopping oriented level		1366
1320	1. While there are some vendors selling		1367
1321	crafts and local goods, shopping is not		1368
1322	a primary focus of the park. The main		1369
1323	attractions are food and recreational		1370
1324	activities.		1371
1325	— Example 2 Ends —		1372
1326	Given reviews: {reviews}		1373
1327	Your evaluation:		1374
1328	Here is the restaurant review extraction prompt:		1375
1329	You are an assistant designed to summarize		1376
1330	reviews of businesses for travel planning		1377
1331	purposes. Your goal is to provide		1378
			1379

1380	3. Evaluate the service of the restaurant	This place has a rating of 4 for fresh-	1429
1381	in general with a scale of 1 to 5, consid-	ness. Many reviews highlight the fresh-	1430
1382	ering waiting time, service, and any in-	ness of ingredients, particularly in salads	1431
1383	teraction between the guest and the staff.	and seafood dishes. The house-baked	1432
1384	4. Evaluate the environment of the restau-	focaccia and pastries are also noted for	1433
1385	rant from 1 to 5. Including the cleanli-	their quality.	1434
1386	ness of the restaurant, the kitchen, the	This place has a rating of 3 for ser-	1435
1387	surroundings, as well as the decorations	vice. Service experiences are mixed,	1436
1388	and vibes of the restaurant. The better	with some diners reporting attentive and	1437
1389	the environment, the better the score.	friendly staff, while others encountered	1438
1390	5. Evaluate the value of the restaurant	slow service and disorganization. The	1439
1391	from 1 to 5. If it is overly priced then it	inconsistency in service quality is a re-	1440
1392	will have a lower score. If it's closer to	curring theme.	1441
1393	transportation and other attractions then	This place has a rating of 5 for environ-	1442
1394	it might have a higher score.	ment. The restaurant's decor and am-	1443
1395	— Example 1 starts —	biance receive high praise, described as	1444
1396	This place has a rating of 2 for flavor.	beautiful, modern, and inviting. The	1445
1397	The food is often described as bland and	spacious layout and natural lighting con-	1446
1398	mediocre, with many reviewers noting	tribute to a pleasant dining experience.	1447
1399	that it lacks seasoning and freshness.	This place has a rating of 3 for value.	1448
1400	This place has a rating of 2 for fresh-	While some diners feel the prices are	1449
1401	ness. Several reviews mention old or	justified by the quality of food and am-	1450
1402	wilted produce, and issues with food be-	biance, others find the portions small and	1451
1403	ing served cold or not freshly prepared.	the overall experience not worth the cost,	1452
1404	This place has a rating of 2 for service.	leading to a mixed perception of value.	1453
1405	Service is frequently criticized for be-	— Example 2 Ends —	1454
1406	ing slow, inattentive, or unprofessional,	Given reviews: {reviews}	1455
1407	with multiple reports of staff ignoring	Your evaluation:	1456
1408	customers or being rude.	E.2 Human Query Generation	1457
1409	This place has a rating of 3 for environ-	The following prompt asks LLM to generate a	1458
1410	ment. The diner has a clean and mod-	human-like query based on the input preference	1459
1411	ern decor, but the ambiance is often de-	list.	1460
1412	scribed as awkward or uncomfortable	Craft a a human like query for a travel	1461
1413	due to the staff's behavior and the music	plan given the following information.	1462
1414	choice.	The input includes details such as trip	1463
1415	This place has a rating of 2 for value.	duration, budget type, attractions types	1464
1416	Prices are considered high for the quality	that the traveler wants to visit, dining	1465
1417	of food served, leading many to feel that	preferences that they want to try, and ac-	1466
1418	they are not getting good value for their	commodation requirements. Make sure	1467
1419	money.	each pairs of key words, like good envi-	1468
1420	— Example 1 Ends —	ronment, good location, are mentioned	1469
1421	— Example 2 Starts —	specifically.	1470
1422	This place has a rating of 4 for fla-	— Example Starts —	1471
1423	vor. The food generally receives praise	Input:	1472
1424	for its flavor, with standout dishes like	- general: 2 days, moderate budget,	1473
1425	the brown butter ravioli and khachapuri	- attraction: history oriented,	1474
1426	being frequently mentioned. However,	- restaurants: French, good environment,	1475
1427	some dishes were noted as mediocre or		
1428	lacking in flavor.		

1476	- hotel: good quality, good location	- Night Attraction: - Name: XXXX	1522
1477	Output: I want to go for a 2-day trip with	— Example Ends —	1523
1478	a moderate budget. I want to visit some	Given Information: {given_information}	1524
1479	history-oriented attractions. Please find	Query: {query}	1525
1480	some good environment restaurants that	Travel Plan:	1526
1481	provide French cuisine, I want to stay in		
1482	a good quality hotel in a good location.		
1483	— Example Ends —	E.4 Route Optimization Prompt	1527
1484	PromptInput: {input}	Both Task 2 and Task 3, which ask for route opti-	1528
1485	Output:	mization, use this prompt. The difference is that	1529
1486	E.3 No Route Optimization Prompt	Task 2’s given information is all data, while Task	1530
1487	Use this prompt for all tasks that don’t request route	3’s given information is filtered data based on prefe-	1531
1488	optimization. In our design, only Task 1 uses this	rence and the spatial clustering algorithm. The	1532
1489	prompt.	planner module in the Tool-Use mode also uses	1533
		this prompt.	1534
1490	You are a proficient travel planner. Based	You are a proficient travel planner. Based	1535
1491	on the given information and query, you	on the given information and query, you	1536
1492	will generate a travel plan like the follow-	will generate a travel plan like the follow-	1537
1493	ing example. Ensure that all recommen-	ing example. Ensure that all recommen-	1538
1494	dations and their addresses are organized	dations and their addresses are organized	1539
1495	in chronological order for each day. Give	in chronological order for each day. Give	1540
1496	exactly 4 attraction recommendations for	exactly 4 attraction recommendations for	1541
1497	each day. Be considerate, concise and	each day. Be considerate, concise and	1542
1498	well-structured.	well-structured. Please also optimize the	1543
1499	— Example Starts —	routes for the trip. For each day, find at-	1544
1500	Query: I am planning a 2-day trip with	tractions that are close to each other for	1545
1501	an expensive budget. I would like to visit	the recommendations.	1546
1502	some history-oriented attractions. Please	— Example Starts —	1547
1503	recommend Japanese restaurants with a	Query: I am planning a 2-day trip with	1548
1504	good environment. For accommodation,	an expensive budget. I would like to visit	1549
1505	I am looking for a hotel with good loca-	some history-oriented attractions. Please	1550
1506	tion, good quality, and good service.	recommend Japanese restaurants with a	1551
1507	Travel Plan:	good environment. For accommodation,	1552
1508	Day X:	I am looking for a hotel with good loca-	1553
1509	- Accommodation: - Name: XXXX Ad-	tion, good quality, and good service.	1554
1510	dress: XXXX, XXXX	Travel Plan: Day X:	1555
1511	- Breakfast: - Name: XXXX Address:	- Accommodation: - Name: XXXX Ad-	1556
1512	XXXX, XXXX	dress: XXXX, XXXX	1557
1513	- Morning Attraction: - Name: XXXX	- Breakfast: - Name: XXXX Address:	1558
1514	Address: XXXX, XXXX	XXXX, XXXX	1559
1515	- Lunch: - Name: XXXX Address:	- Morning Attraction: - Name: XXXX	1560
1516	XXXX, XXXX	Address: XXXX, XXXX	1561
1517	- Afternoon Attraction: - Name: XXXX	- Lunch: - Name: XXXX Address:	1562
1518	Address: XXXX, XXXX; - Name:	XXXX, XXXX	1563
1519	XXXX Address: XXXX, XXXX	- Afternoon Attraction: - Name: XXXX	1564
1520	- Dinner: - Name: XXXX Address:	Address: XXXX, XXXX; - Name:	1565
1521	XXXX, XXXX	XXXX Address: XXXX, XXXX	1566

1567	- Dinner: - Name: XXXX Address:	Preference: A list of preferences men-	1613
1568	XXXX, XXXX	tioned in the query.	1614
1569	- Night Attraction: - Name: XXXX	Example: AttractionSearch[Cheap bud-	1615
1570	— Example Ends —	get,[Nature Oriented]] would return the	1616
1571	Given Information: {given_information}	cheap price and nature - oriented attrac-	1617
1572	Query: {query}	tions.	1618
1573	Travel Plan:	(3) RestaurantSearch[Budget, Cuisine,	1619
1574	E.5 Tool use: ReACT Prompt	Preference]:	1620
1575	Here is the prompt inspired by ReACT (Yao et al.,	Description: Find the restaurants that	1621
1576	2022) and TravelPlanner (Xie et al., 2024).	matches the preference.	1622
1577	Collect information for a query plan us-	Parameters:	1623
1578	ing interleaving 'Thought', 'Action', and	Budget: The budget mentioned in the	1624
1579	'Observation' steps. Ensure you gather	query.	1625
1580	valid information related to transporta-	Cuisine: The cuisine mentioned in the	1626
1581	tion, dining, attractions, and accommo-	query.	1627
1582	dation. All information should be writ-	Preference: A list of preferences men-	1628
1583	ten in Notebook, which will then be in-	tioned in the query.	1629
1584	put into the Planner tool. Note that the	Example: RestaurantSearch[Expensive	1630
1585	nested use of tools is prohibited. Don't	budget, Vietnamese, [Good Flavor, Good	1631
1586	include phrases like "Action: ", "Action	Value]] would return the expensive	1632
1587	5", "Thought 1", or "Thought: "in your	restaurants that offer Vietnamese cuisine,	1633
1588	response. 'Thought' can reason about	with good or excellent flavor and good	1634
1589	the current situation, and 'Action' can	or excellent value.	1635
1590	have 5 different types:	(4) BusinessClusterSearch[]:	1636
1591	(1) Accommodation-	Description: A tool that finds the num-	1637
1592	Search[Budget,Preference]:	ber of business clusters given the infor-	1638
1593	Description: Find the accommodation	mation that you've collected. The tool	1639
1594	that matches the preference.	will choose what business to be consid-	1640
1595	Parameters:	ered and return their spatial clustering	1641
1596	Budget: The budget mentioned in the	information.	1642
1597	query.	Example: BusinessClusterSearch[]	1643
1598	Preference: A list of preferences men-	would return you a list of business	1644
1599	tioned in the query.	clusters among some attractions and	1645
1600	Example: Accommodation-	hotels that you've collected. The	1646
1601	Search[Moderate Budget,[Good	businesses in the same cluster indicates	1647
1602	Location, Good Service]] would return	that they are closer to each other and	1648
1603	the moderate price hotel that has a good	prefered to be arranged for the same day	1649
1604	or excellent location, as well as a good	of the travel.	1650
1605	or excellent service.	(5) Planner[Query]	1651
1606	(2) AttractionSearch[Budget, Prefer-	Description: A smart planning tool that	1652
1607	ence]:	crafts detailed plans based on user input	1653
1608	Description: Find the attractions that	and the information stored in Notebook.	1654
1609	matches the preference.	Parameters:	1655
1610	Parameters:	Query: The query from user.	1656
1611	Budget: The budget mentioned in the	Example: Planner[Give me a 3-day trip	1657
1612	query.	plan in Philadelphia] would return a de-	1658
		tailed 3-day trip plan.	1659

1660 You should use as many as possible steps
1661 to collect enough information to input
1662 to the Planner tool.

1663 Each action only calls one function once.
1664 Do not add any description in the action.
1665 Do not start action with "1. ", state the
1666 action directly.

1667 Query: {query}{scratchpad}

1668 E.6 Itinerary Entry Extraction Prompt

1669 Extract the travel itinerary and parse the
1670 businesses' information into the JSON
1671 format as below. Be faithful and con-
1672 cise. Correctly document the right num-
1673 ber of the attractions. Only write down
1674 the name and address of the businesses.
1675 If certain recommendations (like meals
1676 or accommodations) are not provided, re-
1677 place the information with "-" for name
1678 and address. If recommendations for a
1679 session of attraction is not provided, re-
1680 place the information as an empty array.

1681 F Algorithms

1682 F.1 Total Distance Gap Algorithm

1683 We adapted the classic Traveling Salesmen Prob-
1684 lem (Hoffman et al., 2013) algorithm to allow the
1685 evaluation between multiple days and different re-
1686 turning points. This algorithm can calculate the
1687 optimized routes when the returning hotel is differ-
1688 ent each night.

1689 The main adaptation happens in line 18, where
1690 multiple hotel returns across the trip is allowed and
1691 the calculation of plan wise TSP is made possible.
1692 The main limitation currently is the TSP algorithm
1693 have limits about the number of node in the calcu-
1694 lation. Gemini tends to recommend more than 5
1695 attractions per day make the evaluation not possible
1696 with current algorithm.

```
1697 def totalCost_Multiday(mask, pos, day,  
1698   cordinates, n, visited, cost, info_lists,  
1699   memo):  
1700     visit_requirement = len(cordinates[day])  
1701     distance_list = []  
1702     i_list = []  
1703  
1704     # Get which hotel need to return to for  
1705     current day  
1706     hotel_index = getHotelIndex(day, cordinates)  
1707  
1708     # Base case: if all cities are visited,  
1709     return to hotel for current day  
1710     if mask == (1 << n) - 1:  
1711         return cost[pos][hotel_index]  
1712
```

```
1713 # Memorization  
1714 if memo[pos][mask] != -1:  
1715     return memo[pos][mask]  
1716  
1717 # Main Adapation: This condition check  
1718 allows returned to different hotels and  
1719 break down days in plan  
1720 if visit_requirement == visited:  
1721     for i in range(n):  
1722         if (mask & (1 << i)) == 0:  
1723             i_list.append(i)  
1724             distance_list.append(  
1725                 cost[hotel_index][i] +  
1726                 totalCost_multiday(  
1727                     mask | (1 << i), i, day +  
1728                     1, cordinates, n, 2,  
1729                     cost, info_lists, memo  
1730                 )  
1731             )  
1732  
1733     info_list = [pos, i_list, distance_list]  
1734     info_lists.append(info_list)  
1735  
1736     return min(distance_list) + cost[pos][  
1737         hotel_index] # change this to the old  
1738         hotel position  
1739  
1740 # Try visiting every city that has not been  
1741 visited yet  
1742 for i in range(n):  
1743     if (mask & (1 << i)) == 0:  
1744         i_list.append(i)  
1745         # If city i is not visited, visit it  
1746         and update the mask  
1747         distance_list.append(  
1748             cost[pos][i] + totalCost_multiday(  
1749                 mask | (1 << i), i, day,  
1750                 cordinates, n, visited +  
1751                 1, cost, info_lists, memo  
1752             )  
1753         )  
1754  
1755 # Store an info_list to retrieve the  
1756 optimized order  
1757 info_list = [pos, i_list, distance_list]  
1758 info_lists.append(info_list)  
1759  
1760 memo[pos][mask] = min(distance_list)  
1761  
1762 return min(distance_list)  
1763
```