



UNIVERSITY OF TRENTO

Department of Information Engineering and Computer Science

Master's Degree in
Artificial Intelligence Systems

FINAL DISSERTATION

EXPLORING THE USE OF LLMs FOR AGENT PLANNING: STRENGTHS AND WEAKNESSES

Supervisor
Paolo Giorgini

Student
Davide Modolo
229297

Academic year 2023/2024

Contents

Abstract	1
1 Introduction	2
2 Background	4
2.1 Artificial Intelligence	4
2.2 Large Language Models	5
2.2.1 Self-Attention Mechanism	6
2.2.2 LLMs' Uncertainty	6
2.2.2.1 Expressing Uncertainty	7
2.2.2.2 Stable Explanations as Confidence Measures	7
2.2.2.3 Tokens' log-probability	7
2.3 Agents	8
2.3.1 BDI Architecture	10
2.3.1.1 Core Components of BDI	10
2.4 State of the Art	10
2.4.1 PDDL-Based Solutions	11
2.4.2 Reinforcement Learning Solutions	13
2.4.3 Planning with LLM	15
2.4.3.1 Chain-of-Thought Reasoning	15
2.4.3.2 Zero-Shot and Few-Shot Planning	16
3 Experiment Setting	18
3.1 Problem Definition	18
3.1.1 Our Task	18
3.2 Environment - Deliveroo.js	19
3.2.1 Server Configuration and Event Handling	21
3.3 Large Language Models Selection	22
3.3.1 Open Source Models	24
3.3.1.1 Challenges with Open Source Models	24
3.3.1.2 LLaMa 3.2	25
3.3.1.3 Gemma 2	25
3.3.1.4 DeepSeek-V3	25
3.3.2 Closed Source Models - OpenAI	26
3.3.2.1 GPT-4o-mini	27
3.3.2.2 GPT-4o	27
3.3.2.3 GPT-3.5-turbo	28
4 Agent Development	29
4.1 Development information	29
4.2 First Approach	29
4.2.1 Helping Parameters	30
4.2.2 Takeaways	31
4.3 Second Approach	32

4.3.1	Takeaways	33
4.4	Final Agent	33
4.4.1	Prompt Management System	33
4.4.2	Agent Refactoring	34
4.4.3	RAG Experiments	34
4.4.4	Stateless and Stateful Agents	34
4.4.5	Uncertainty Handling	35
4.4.6	Takeaways	35
4.5	Extra: Closest Cell to the Goal	36
5	Data Collection	38
5.1	Prompts	38
5.1.1	Pickup prompt	38
5.1.2	Deliver prompt	39
5.2	Prompt Creation Choices	40
5.2.1	Role Prompting	40
5.2.2	Map Encoding to Reduce Attention Sparsity	40
5.2.2.1	Emerging behavior: encoded question decoded answer	41
5.2.3	Emerging behavior: math capabilities	44
5.2.4	Question Structure	45
5.2.4.1	Structure of the Paper’s Prompt	45
5.2.4.2	Comparison with Our Approach	45
5.2.4.3	Multichoice Benchmarking	46
5.2.5	Goal positioning	46
5.2.5.1	Leveraging Prompt Caching	47
5.3	Uncertainty Visualization	47
5.3.1	Heatmaps	48
5.3.1.1	Example	50
5.3.2	Correctness Heatmaps	51
5.3.2.1	Example	51
6	Results Discussion	54
6.1	Map Orientation	54
6.1.1	Comparison of Orientations	55
6.2	Stateless	56
6.2.1	Pickup Goal at the Center	56
6.2.2	Deliver Goal at the Center	58
6.2.3	Pickup and Deliver Goals in Different Map Sections	58
6.2.4	Goal position comparison	61
6.3	Stateful	62
6.3.1	Path Visualization	62
6.4	“Pathfinding”	62
6.5	Stateless & Stateful - Performance Summary	64
6.6	Insights from the Closest Cell to the Goal Approach	65
6.7	Model Comparison	65
7	Conclusions	67
	Bibliography	68
	A Acronyms	73
	B Prompts	74

Abstract

With the recent advancements in Large Language Models, there has been a growing interest in developing agents capable of understanding and executing complex tasks. In this work, we explore the use of LLMs as agents that can navigate and complete logistics tasks, primarily focused on picking up and delivering parcels, within a web-based environment.

Our approach aims to evaluate the raw performance of an LLM without integrating additional frameworks or specialized optimization techniques, allowing us to assess its inherent generative capabilities in problem-solving. We analyze how effectively the agent can navigate different map layouts and complete assigned objectives, testing its adaptability across various goal configurations. A key aspect of our evaluation is the use of LLM uncertainty measures, derived from tokens' log probabilities, to gain deeper insights into the model's confidence in its decision-making process. These measures help us understand when the agent is uncertain and how that uncertainty correlates with performance in different parts of the scenario.

We demonstrate that the agent's performance improves when using newer LLM versions, reflecting the continuous advancements in these models. However, we also observe a decline in performance as the map size increases, suggesting that larger environments pose challenges that the model struggles to overcome. To structure our approach effectively, we design the prompt based on established literature, ensuring alignment with best practices in prompt engineering.

Furthermore, we experiment with two distinct agent configurations: a stateless agent, which makes decisions solely based on the current state of the environment, and a stateful agent, which retains memory of past interactions. By comparing these approaches, we highlight the strengths and limitations of each. The stateless agent benefits from simplicity and avoids memory-related constraints, but it may struggle in scenarios where the environment description requires too much attention. Conversely, the stateful agent provides improved continuity in decision-making but faces challenges related to context length limitations and potential inconsistencies in stored information.

The code and data used in this work are available on GitHub[5]. Appendix A provides a list of acronyms used throughout this thesis.

1 Introduction

The field of Artificial Intelligence has seen significant progress with the advent of Large Language Models, which leverage deep learning techniques to process and generate human-like text. These models, built on the Transformer architecture, have revolutionized natural language understanding and processing in many ways. In the context of this thesis, the ability of LLMs to act autonomously on planning and executing complex tasks has become a key research direction. Traditional planning approaches, such as those based on Reinforcement Learning and heuristic-based search algorithms (e.g. by using Planning Domain Definition Language), have long been used to tackle structured decision-making problems. However, with the rise of LLMs, new possibilities have emerged for solving planning problems in a more flexible and adaptable manner.

The general problem addressed in this work revolves around planning with LLMs, particularly in the domain of logistics tasks. Logistics problems often involve dynamic environments where an agent must plan and execute a sequence of actions to achieve a goal, such as picking up and delivering parcels in our case. Classical AI planning techniques require predefined rules, domain knowledge, and structured representations of the world, while reinforcement learning-based approaches demand extensive training on simulated environments. LLMs, on the other hand, introduce a novel paradigm where reasoning and decision-making emerge from large-scale pretraining on diverse textual data. The challenge is to determine whether LLMs can effectively function as planners without explicit search algorithms or fine-tuned optimization techniques.

Current research on AI planning has extensively explored methods leveraging PDDL, which formalizes decision-making problems in a structured manner, allowing traditional planners to compute optimal action sequences. Reinforcement Learning, another key approach, enables agents to learn optimal strategies through trial and error, with the final goal to maximize a reward. More recently, researchers have investigated the potential of LLMs in planning, leveraging techniques such as Chain-of-Thought reasoning and few-shot prompting to guide models toward generating coherent plans. However, these methods still face challenges, particularly in maintaining consistency, handling uncertainty, and ensuring goal-directed behavior in complex environments.

This work focuses on a specific logistics task in which an agent, powered solely by an LLM, must pick up and deliver parcels in a simulated environment while considering uncertainty in its decision-making process. Unlike traditional planning approaches, no additional external frameworks are integrated; instead, the LLM itself serves as the core reasoning engine via its generative capabilities. A key aspect of the research is the evaluation of uncertainty, using measures derived from token log-probabilities to assess the model’s confidence in its choices. By analyzing how uncertainty impacts decision-making, we aim to identify the strengths and weaknesses of this approach.

The objective of this study is to understand the capabilities and limitations of LLMs when applied to planning problems, particularly in dynamic and uncertain environments. Through systematic experiments, we aim to evaluate how well an LLM can navigate and complete logistics tasks, whether it can adapt to varying scenarios, and how uncertainty influences its performance. The expected results include insights into the effectiveness of LLM-based planning using the models’ generative capabilities and potential systematic weaknesses.

This thesis is structured as follows:

- **Chapter 2 - Background** establishes the theoretical foundations necessary to understand the problem space and the methodologies explored in the research. It covers core AI concepts, the evolution of LLMs, and their architecture, with a particular focus on the Attention mechanism and token-based uncertainty estimation. Additionally, it analyzes traditional planning approaches such as rule-based systems, search-based techniques, and Reinforcement Learning with newer LLM-driven methods, providing a comprehensive view of the current state of the

art;

- **Chapter 3 - Experiment Setting** formalizes the objective of this thesis, describing the environment used to evaluate the LLM-driven agent, detailing the rules of the task, the constraints imposed by the system as well as the functioning of the web-based environment used for testing, and how agent interactions are structured. It explains the decisions behind the model selection, highlighting the difference between uncertainty computation in both open source and closed source models;
- **Chapter 4 - Agent Development** illustrates the iterative development process of the LLM-based agent, outlining key design decisions and the challenges encountered. It covers the evolution of the different implementation strategies, including the final agent divided in its stateful and stateless versions;
- **Chapter 5 - Data Collection** focuses on the impact of prompt design on the agent’s performance and decision-making capabilities. It explores various prompt engineering strategies based on the literature, examining how different wording structures and contextual information influence the LLM’s ability to understand and execute actions within the environment. Furthermore, it details how uncertainty can be visualized and analyzed to provide insights into the model’s behavior;
- **Chapter 6 - Results Discussion** shows and analyze the findings from the experimental phase, highlighting key trends, limitations and recurrent behaviors observed in the LLM’s performance. The discussion addresses the impact of different variables, such as map size, task complexity, and uncertainty estimation on the overall success rate of the agent. Comparative insights between stateless and stateful agents are provided, with their respective advantages and drawbacks;
- **Chapter 7 - Conclusions** summarizes the entire work, providing a synthesis of the key insights gained from evaluating LLMs in this scenario. It outlines current limitations and potential directions for future research, including improvements in uncertainty modeling and leveraging, hybrid approaches integrating different AI models, and the exploration of more structured approaches to limit the impact of uncertainty on decision-making.

2 Background

In this thesis, we analyze in detail the behavior of a Large Language Model as an agent performing a logistics task within a controlled environment.

Before presenting all the work carried out in detail, this chapter aims to provide a comprehensive explanation of all the theoretical foundations necessary to understand the steps presented in the following chapters. Starting from a brief introduction of Artificial Intelligence just to define the boundaries in which we are working, we will move to the core concepts. In particular, we want to highlight what an LLM is and how it works, with a special focus on the Attention mechanism and how the uncertainty of an LLM can be calculated. This will serve as a basis for correctly interpreting the analysis in Chapter 6.

There will also be a broader discussion on agents in a strict sense and *LLM agents* to better show the difference between our implementation and what is currently being discussed in the literature.

To better define the context of this thesis, we will also examine the main alternative approaches to solving a logistics problem currently recognized as the state of the art.

2.1 Artificial Intelligence

Artificial Intelligence is a very broad field, that can be resumed as the category of systems designed to perform tasks that traditionally require intelligence, such as *Natural Language Understanding* and visual perception.

In recent years, AI has rapidly evolved, driven by advances in Deep Learning¹, increased computational power and the easy availability of massive datasets (language models are even trained on the entirety of the internet). Early AI systems, including expert systems and early Machine Learning models, relied on manually crafted rules or statistical techniques. However, with the rise of Neural Networks, particularly Deep Learning models, AI has shifted toward self-learning systems capable of extracting complex patterns from raw data.

One of the key breakthroughs in this evolution was the development of deep neural networks (DNNs), particularly Convolutional Neural Networks (CNNs) for image processing, introduced by Krizhevsky et al. [16] and Recurrent Neural Networks (RNNs) for sequential data, including language modeling, introduced by Hochreiter et al. [11].

Despite their success, RNNs struggled with long-term dependencies due to vanishing gradients², leading to the development of the *Transformer* architecture (from Vaswani et al., ‘Attention Is All You Need’[31]), which eliminated recurrence in favor of self-attention mechanisms, significantly improving efficiency and scalability in natural language processing.

In general, two main categories of Machine Learning models can be defined: discriminative models and generative models.

Discriminative models Discriminative models are a class of machine learning models that aim to directly model the decision boundary between different classes in a dataset. Unlike generative models, which learn the underlying distribution of the data, discriminative models focus on learning the conditional probability of a target class given the input features. Classical models like Support Vector Machines³ and Conditional Random Fields⁴ have been widely used for text classification and sequence labeling tasks such as Named Entity Recognition (Lafferty et al. [17]). More recently, deep learning-based models like BERT (Devlin et al. [7]) have been invented, that leverage contextualized

¹https://en.wikipedia.org/wiki/Deep_learning

²https://en.wikipedia.org/wiki/Vanishing_gradient_problem

³https://en.wikipedia.org/wiki/Support_vector_machine

⁴https://en.wikipedia.org/wiki/Conditional_random_field

word representations to improve performance on tasks like sentiment analysis, intent detection, and slot filling.

Generative models Generative models learn the underlying data distribution to create new samples that resemble the original data. This category includes several architectures that have pushed the boundaries of AI-generated content. Variational Autoencoders (Kingma and Welling [14]) introduced a probabilistic approach to generating structured data, while Generative Adversarial Networks (Goodfellow et al. [9]) refined the concept by using two competing neural networks, a generator and a discriminator, to iteratively improve synthetic data generation. More recently, diffusion models (Ho et al. [10]) have surpassed GANs in generating high-quality images by modeling data transformations through iterative denoising processes. In the domain of text generation, models like GPT (Radford et al. [24]) demonstrated the power of large-scale, unsupervised pretraining. These Large Language Models predict the next token (that can be seen as a building block of a word) in a sequence based on vast amounts of textual data, learning contextual nuances and producing human-like responses.

2.2 Large Language Models

Large Language Models are a class of deep learning models that leverage the transformer architecture to generate coherent and contextually relevant text. These models have revolutionized natural language processing by achieving state-of-the-art performance on a wide range of tasks, including language modeling, translation, summarization, and question-answering.

The Transformer architecture, introduced by Vaswani et al. in the paper ‘Attention Is All You Need’ [31], is the foundation of LLMs. It consists of an encoder-decoder structure, where the encoder processes the input sequence and generates a sequence of hidden states, while the decoder generates the output sequence based on the encoder’s hidden states. The key innovation in transformers is the self-attention mechanism, which allows the model to weight the importance of different input tokens when generating the output. This mechanism enables transformers to capture long-range dependencies and contextual information more effectively than RNNs.

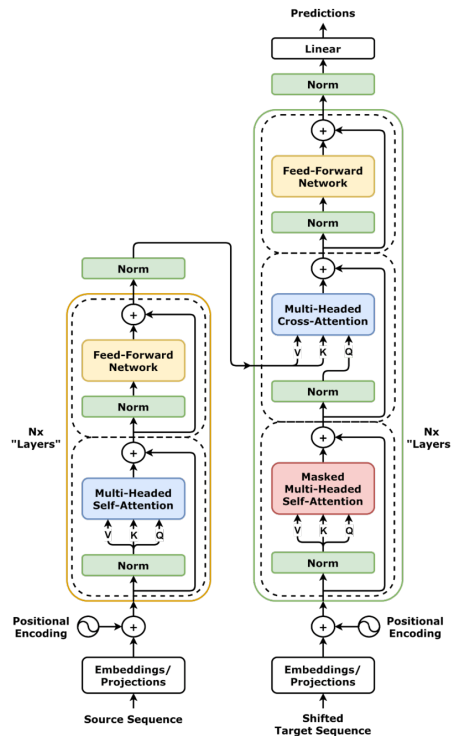


Figure 2.1: Transformer Architecture
Source: Vaswani et al., Attention Is All You Need [31]

2.2.1 Self-Attention Mechanism

The self-attention mechanism is a fundamental component of the transformer architecture (Figure 2.1), enabling the model to focus on specific parts of the input sequence when generating the output. It computes a weighted sum of the input tokens, where the weights are learned during training based on the relevance of each token to the current context.

The self-attention mechanism works in this way:

1. create 3 vectors from embeddings (*Query*, *Key*, *Value*) multiplying by 3 weight matrices learned during the training process;
2. calculate a score that determines how much focus goes to different parts of the input sentence as it encodes a word;
3. divide the score for more stable gradients and apply softmax;
4. multiply each value vector by the softmax score to keep the value of the word it focuses on, and sink other irrelevant words;
5. sum the weighted value vectors: this produces the output of the self-attention layer at this position.

The self-attention operation computes the relevance of each token in the input with respect to the query token using the scaled dot-product attention:

$$\text{Attention}(Q, K, V) = \text{softmax} \left(\frac{QK^T}{\sqrt{d_k}} \right) V$$

where d_k is the dimensionality of the key vectors, ensuring that the dot products do not grow too large as input size increases. The softmax function normalizes the scores into attention weights, which determine how much influence each token should have on the final representation.

Multi-head attention extends this mechanism by computing multiple sets of Q, K, V matrices in parallel, allowing the model to capture different aspects of contextual relationships:

$$\text{MultiHead}(Q, K, V) = \text{Concat}(\text{head}_1, \dots, \text{head}_h) W^O$$

where each attention head independently applies scaled dot-product attention, and the outputs are concatenated and linearly projected using W^O (Weight matrix). This improves the model's ability to encode complex dependencies and contextual meaning.

The attention mechanism allows the model to focus on different parts of the input sequence based on the current context, enabling it to capture long-range dependencies.

2.2.2 LLMs' Uncertainty

Despite their impressive capabilities, LLMs are inherently probabilistic and can generate responses that are syntactically correct yet factually inaccurate. Understanding and quantifying this uncertainty is crucial for evaluating the reliability of generated text, especially in high-stakes applications such as medical diagnosis, legal advice, or automated decision-making.

For example, if an LLM generates an answer to a yes/no question with probabilities:

$$P(\text{Yes}) = 0.51, P(\text{No}) = 0.49$$

then the model is nearly uncertain, and this information should be communicated rather than presenting "Yes" as a definitive response.

A key consequence of uncertainty is the phenomenon of *hallucination*, where the model generates confident but factually incorrect or fabricated information [13]. Hallucinations arise when:

- the model lacks knowledge about a specific query but still generates an answer;
- the training data contains conflicting or misleading patterns;

- the model overgeneralizes from limited training examples.

Mitigating hallucinations involves uncertainty-aware generation techniques, the most common one is *Retrieval-Augmented Generation* [18], which enhances the prompt with additional context from a knowledge base to improve the model’s factual accuracy.

The literature is studying different approaches to quantify uncertainty in LLMs, noting that this is a still active field of research.

2.2.2.1 Expressing Uncertainty

A study titled ‘Can LLMs Express Their Uncertainty? An Empirical Evaluation of Confidence Elicitation in LLMs’ [35] investigates methods for eliciting confidence from LLMs without accessing their internal parameters or fine-tuning. The researchers propose a framework comprising three components:

- Prompting Strategies: Techniques to elicit verbalized confidence from the model;
- Sampling Methods: Generating multiple responses to assess variability;
- Aggregation Techniques: Computing consistency across responses to determine confidence levels.

The study evaluates these methods on tasks such as confidence calibration and failure prediction across various datasets and LLMs.

Key findings indicate that LLMs often exhibit overconfidence when verbalizing their certainty, possibly mirroring human confidence expression patterns. Additionally, as model capabilities increase, both calibration and failure prediction performance improve, though they remain suboptimal. They show that implementing strategies like human-inspired prompts and assessing consistency among multiple responses can mitigate overconfidence.

2.2.2.2 Stable Explanations as Confidence Measures

In the pursuit of reliable uncertainty quantification in Large Language Models, the paper ‘Cycles of Thought: Measuring LLM Confidence through Stable Explanations’ [2] introduced a novel framework that assesses model confidence through the stability of generated explanations.

Their approach posits that the consistency of explanations accompanying an answer can serve as a proxy for the model’s certainty. Instead of assigning a single probability to an answer, the method generates multiple explanations for the same question and treats each explanation-answer pair as a distinct classifier. A posterior distribution is then computed over these classifiers, allowing for a principled estimation of confidence based on explanation stability. If the model’s explanations remain stable across different reasoning paths, it suggests high confidence in the answer. Conversely, significant variation in explanations signals uncertainty. Empirical evaluations across multiple datasets demonstrated that this framework enhances confidence calibration and failure prediction, outperforming traditional baselines.

However, there are some potential drawbacks. The method requires generating multiple explanations, which increases computational cost and latency. Additionally, it can be sensitive to prompt variations, and may misinterpret repetitive patterns as high confidence.

2.2.2.3 Tokens’ log-probability

The paper ‘Robots That Ask For Help: Uncertainty Alignment for Large Language Model Planners’ [28] introduces the KnowNo framework, which is the one we took inspiration from to quantify the uncertainty of the agent in this thesis.

The KnowNo framework leverages Conformal Prediction⁵, a statistical method that provides formal guarantees on the reliability of predictions to assess uncertainty.

In the paper, they ask the LLM to generate a set of four actions for a given prompt (since the `logit.bias` parameter in the OpenAI API was limited to five tokens at that time, more on this in Section 3.3.2), and then they append a “no-op” action to the set. This will not be the case of this thesis, since the actions will always be the same, but we will use the same math behind the uncertainty calculation.

⁵https://en.wikipedia.org/wiki/Conformal_prediction

Then, they ask the model for the action to select, adding the bias to the tokens representing each action. Then, they use the log-probabilities of the tokens (referring to the actions) to compute the uncertainty of the model.

KnowNo computes uncertainty evaluating the “validity” of each option: CP calculates a confidence interval based on previous data, and from this, a set of valid actions is generated (based on their scaled log-probability). This set can include one or more actions, and the size of this set is indicative of the level of uncertainty:

- **Singleton:** If CP narrows down the options to just one action, this indicates low uncertainty, and the robot can proceed confidently with the task. The model is highly certain that this action is the most appropriate next step;
- **Multiple Options:** When CP leaves multiple possible actions in the valid set, this may indicate high uncertainty. In such cases, KnowNo triggers the robot to request human assistance. This allows the robot to seek clarification when it is unsure, thereby avoiding errors that might arise from acting on uncertain predictions.

A simplified version of the KnowNo flow can be seen in Figure 2.2. Technically speaking, the computation of the uncertainty can be summarize in 5 steps:

1. give each action a single-token label (eg. A), B), C), D), E));
2. use the `logit_bias` parameter in the API to force the model to only answer using these labels;
3. get the log-probabilities of the tokens and scale them: this results in a “confidence” value for each token;
4. filter the resulting set of option with a threshold computed with CP;
5. either the result will be a singleton (denoting no uncertainty) or a set of options.

On a final note, in the paper they state that this framework has the advantage of being model-agnostic, as it can be applied to LLMs out-of-the-box without requiring any fine-tuning, thanks to the “caution” that is given if the resulting filtered set of options is not a singleton.

2.3 Agents

As widely explained in the book ‘An Introduction to Multiagent Systems’ [34], we can summarize the definition of an agent as an autonomous entity that perceives its environment through sensors and acts upon it through effectors, making decisions based on its perceptions and objectives in order to achieve specific goals.

This definition highlights several key aspects of agents:

- **Autonomy:** Agents operate without direct human intervention, controlling their own actions;
- **Perception and Action:** They interact with the environment via sensors (perception) and actuators (action execution);
- **Decision-making:** Agents select actions based on their internal model, goals, and the state of the environment;
- **Non-determinism and Adaptability:** Since environments are generally non-deterministic, agents must be prepared for uncertainty and potential failures in action execution;
- **Preconditions and Constraints:** Actions are subject to certain conditions that must be met for successful execution.

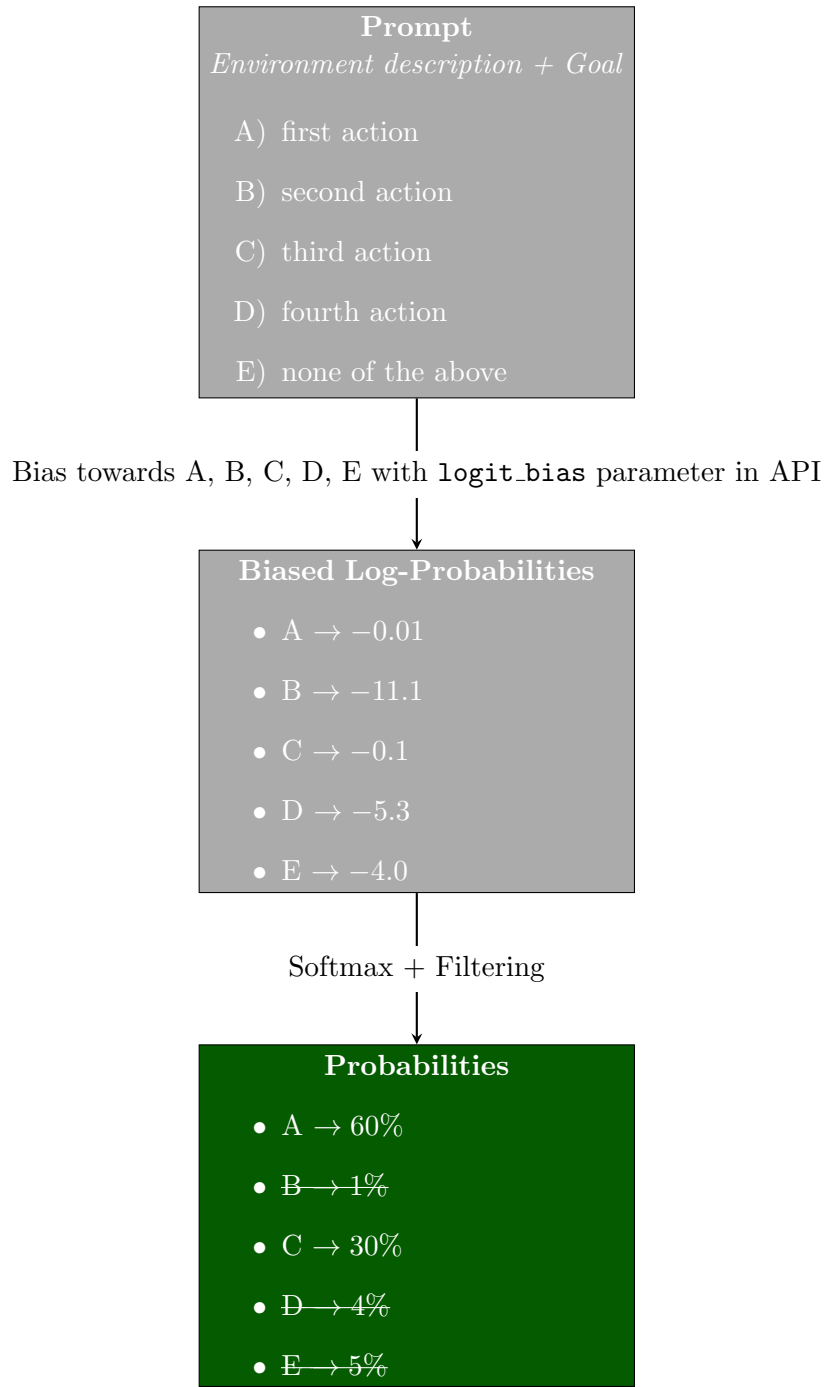


Figure 2.2: KnowNo Uncertainty Computation

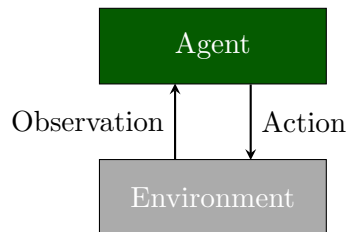


Figure 2.3: Agent Design Scheme
Source: redesign of a scheme in [34]

Thus, an agent’s fundamental challenge is deciding which actions to perform in order to best satisfy its objectives, given the constraints and uncertainties of its environment.

As shown in Figure 2.3, an agent is some entity that perceives the environment and reacts to it. The setting can be anything from a simple thermostat to a complex system like a self-driving car. The idea is that the agent is able to react to a change in the environment and take actions to achieve its goals.

To align our agent with the definition above, we can map some of its concept to what this thesis will analyze:

- **Autonomy:** the agent will choose its action based on the prompt built using the environment information only;
- **Perception and Action:** what the server sends about the current state of the environment can be seen as the perception of the agent, while the action it can take will be given in the prompt in a specific way;
- **Decision-making:** the decision-making process will be the generation of the text by the LLM, weighted by the uncertainty.
- **Non-determinism and Adaptability:** to emulate the non-determinism of the environment, the state received by the server will be used “raw” in the prompt, without any hard processing or parsing;
- **Preconditions and Constraints:** being in a “limited” map with a fixed number of walkable cells, is itself a constraint the agent must consider.

2.3.1 BDI Architecture

The Belief-Desire-Intention (BDI) architecture is a widely adopted framework in artificial intelligence for modeling rational agents. It was formally developed by Rao and Georgeff in 1995 [27] and has been implemented in several architectures. BDI provides a structured approach to practical reasoning, allowing agents to function effectively in dynamic and unpredictable environments.

2.3.1.1 Core Components of BDI

BDI agents operate based on three key components:

- **Belief:** Represents the agent’s knowledge about the world, including past events and observations;
- **Desire (Goals):** Defines the agent’s objectives or preferred end states;
- **Intention:** Represents the commitments of an agent toward achieving specific goals through selected plans.

BDI has been extensively used in fields like robotics, automated planning, and multi-agent systems.

2.4 State of the Art

A logistic problem is a fundamental challenge in the field of Artificial Intelligence, since depending on the complexity of the specific problem, it can contain tasks such as route optimization, supply chain management, and delivery scheduling. These problems arise in various domains, including transportation, e-commerce, and manufacturing, where efficient resource allocation and decision-making are critical. Given the complexity of modern logistics, AI has emerged as a powerful tool for finding optimal or near-optimal solutions.

Traditional research techniques, such as linear programming and heuristics, have been widely employed. However, with the increasing availability of data and computational power, Machine Learning and Deep Learning methods have become more prevalent. These methods can predict demand, optimize routes dynamically, and enhance decision-making under uncertainty based on the data. Additionally, Reinforcement Learning has gained attention for its ability to learn optimal strategies.

In the recent years with the explosion of Large Language Models, many researchers started to include and test them to different fields, including planning and logistics.

2.4.1 PDDL-Based Solutions

Planning Domain Definition Language (PDDL) is a human-readable format for problems in automated planning that gives a description of the possible states of the world, a description of the set of possible actions, a specific initial state of the world, and a specific set of desired goals.

Source: Wikipedia⁶

The fundamental distinction between a PDDL-based solution and any Machine Learning/Deep Learning approach lies in the very nature of how problems are defined and solved.

In a PDDL-based system, the problem must be explicitly encoded using a formal, structured language that describes the initial state, the goal state, and the set of available actions. This formal encoding serves as a blueprint for the planner, which then performs the computationally intensive task of exploring a vast search space. The planner systematically generates and evaluates possible action sequences, using algorithms to determine an optimal path from the initial state to the goal state. This process is highly deterministic, with each action being considered in the context of its direct impact on reaching the goal.

While effective in structured, static environments with well-defined parameters, this approach is inherently time-consuming and computationally demanding. The planner must traverse a potentially enormous state space, guided by heuristics to prune less relevant possibilities, but still constrained by the rigid formalism of PDDL. Because of this, it can struggle with real-time decision-making, particularly in situations where the environment is dynamic, uncertain, or rapidly changing.

PDDL Code

```
1 (define (domain bit-toggle)
2   (:requirements :strips :negative-preconditions)
3   (:predicates
4     (bit ?b)                                ; predicate meaning
5                                           ; bit ?b is set (true)
6   )
7
8   (:action setbit
9     :parameters (?b)
10    :precondition (not (bit ?b))           ; can only set a bit if
11                                           ; it is not already set
12    :effect (bit ?b)                       ; setting the bit to true
13  )
14
15  (:action unsetbit
16    :parameters (?b)
17    :precondition (bit ?b)                 ; can only unset a bit if
18                                           ; it is currently set
19    :effect (not (bit ?b))                 ; setting the bit to false
20  )
21 )
```

Listing 2.1: Domain file example for a bit toggle problem

⁶https://en.wikipedia.org/wiki/Planning_Domain_Definition_Language

PDDL Code

```

1 (define (problem bit-toggle-full-problem)
2   (:domain bit-toggle-full)
3   (:objects
4     b1 b2 b3
5   )
6   (:init)                                ; Initially all bits are unset (
7     false)
8   (:goal                                  ; It can be any combination of T/F
9     (and (bit b1) (bit b2) (not(bit b3)))
10  )
11 )

```

Listing 2.2: Problem file example for a bit toggle problem

With the increasing number of variables (actions or predicates), the number of arcs and nodes grows exponentially. A little example that makes this problem easy to visualize is the Domain where we can have N possible bits, that can be turned to **true** or **false** (Domain file in Listing 2.1) and the Problem where everything starts at **false** and we want a specific final combination (Problem file in Listing 2.2).

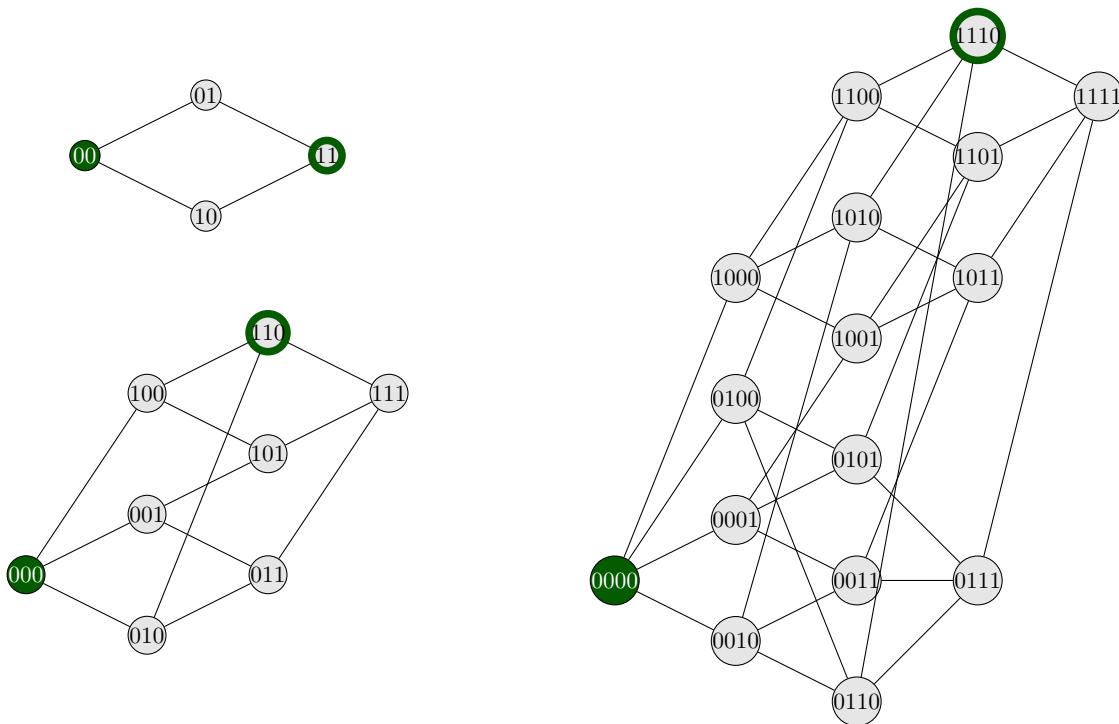
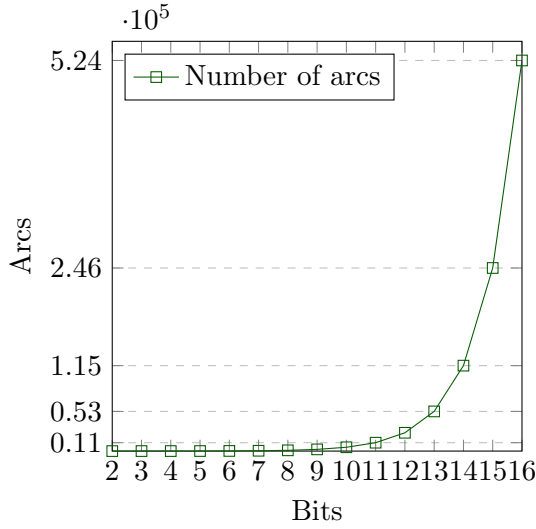
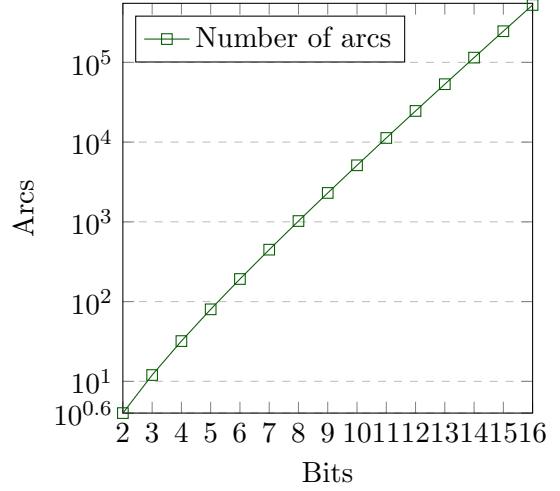


Figure 2.4: Graphs for bit-toggle problem with 2, 3, and 4 bits

As we can see in the plot Figure 2.5, the number of arcs (example of graphs for 2, 3 and 4 bits in Figure 2.4) grows exponentially with the number of bits, as well as the number of states obviously. This shows how even a simple problem with a simple solution can become time-intensive and not suitable for real-time applications.



(a) Arcs over number of bits



(b) Arcs over number of bits(Log Scale)

Figure 2.5: Arcs per Bit

PDDL Code

```

1           ; Found Plan (output)
2 (setbit b2)
3 (setbit b1)

```

Listing 2.3: Plan for the bit toggle problem (110), solved by LAMA-first planner

However, a PDDL approach is more explainable, since all the information is provided by the user and the output result is a sequence of actions (example at Listing 2.3). This makes it easier to understand and debug the solution, as each step is explicitly defined. Of course, there might be different paths to reach the goal, and the planner might choose one based on heuristics or optimization criteria. This transparency in the decision-making process is one of the key advantages of using PDDL for planning problems.

Literature An example of a problem related to the one presented in this thesis, solved using PDDL, can be found in the paper ‘An AI Planning Approach to Emergency Material Scheduling Using Numerical PDDL’ by Yang et al. [36].

In their work, they utilize PDDL 2.1 that allows to model the scheduling problem, incorporating factors such energy consumption constraints. Their approach employs the Metric-FF planner to generate optimized scheduling plans that minimize total scheduling time and transportation energy usage. However, while this demonstrates the applicability of AI planning to emergency logistics, their model simplifies the real-world scenario by assuming predefined transport routes, limited vehicle types, and abstract representations of congestion effects. This highlights a broader limitation of PDDL in capturing the full complexity of dynamic and uncertain environments often encountered in emergency response situations.

2.4.2 Reinforcement Learning Solutions

Reinforcement Learning is a branch of machine learning focused on making decisions to maximize cumulative rewards in a given situation. Unlike supervised learning, which relies on a training dataset with predefined answers, RL involves learning through experience. In RL, an agent learns to achieve a goal in an uncertain, potentially complex environment by performing actions and

receiving feedback through rewards or penalties.

*Source: GeegksforGeeks*⁷

Reinforcement Learning is a learning setting, where the learner is an Agent that can perform a set of actions depending on its state in a set of states and the environment.

It works by defining:

- **Environment:** the world in which the agent operates;
- **Agent:** the decision-maker that interacts with the environment;
- **Actions:** the possible moves the agent can make;
- **Rewards:** the feedback the agent receives for its actions;
- **Policy:** the strategy the agent uses to select Actions.

In performing action a in state s , the learner receive an immediate reward $r(s, a)$. In some states, some actions could be not possible or valid.

The task is to learn a policy π (a full specification of what action to take at each state) allowing the agent to choose for each state the action maximizing the overall reward, including future moves.

To deal with this delayed reward problem, the agent has to trade-off exploitation and exploration:

- **Exploitation:** the agent chooses the action that it knows will give some reward;
- **Exploration:** the agent tries alternative actions that could end in bigger rewards.

When considering a logistics problem, Reinforcement Learning naturally comes to mind. This is because defining a reward function is relatively straightforward: it could be measured in terms of packages delivered per minute, per step, or a similar metric. Additionally, the entire process can be simulated in a virtual environment, allowing multiple parallel simulations to accelerate the agent's learning process. As illustrated in Figure 2.6, the structure of the Reinforcement Learning framework closely resembles the agent-based model depicted in Figure 2.3. In both cases, the agent interacts with its environment, receives feedback in the form of rewards, and, in RL, continuously refines its policy to optimize future performance.

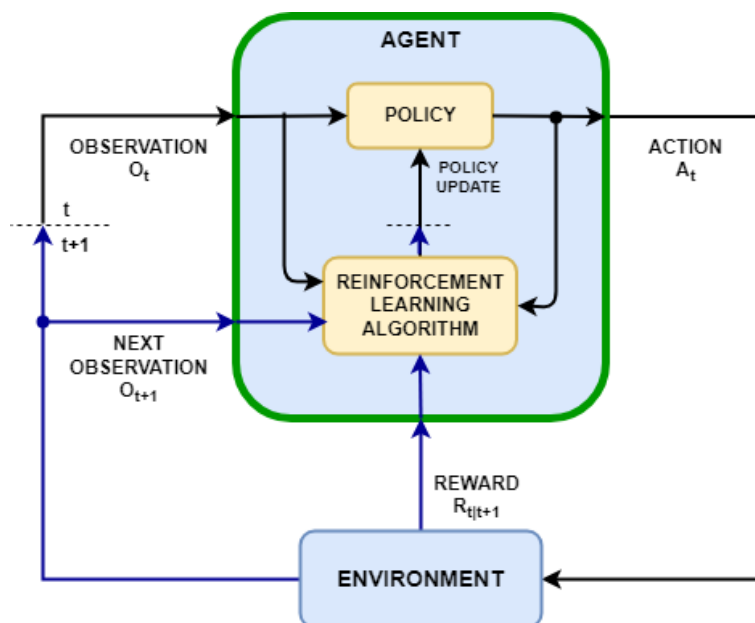


Figure 2.6: RL Agent Scheme

*Source: Mathworks*⁸

⁷<https://www.geeksforgeeks.org/what-is-reinforcement-learning/>

However, RL has its own set of challenges. The most common one is the convergence to a local optima in the reward function. This means that the agent might become stuck in suboptimal strategy that is not the best one. Moreover, RL is not fully explainable, meaning that we can't understand why the agent took a specific action in a specific situation.

Another issue with RL is the cost of training. To learn, the RL agent needs to perform a large number of actions to explore the environment and learn the best strategy. This can be computationally expensive and time-consuming, especially for complex problems with many variables and states. Moreover, once the agent is trained, its adaptability to new environments or situations is limited, as it is optimized for a specific reward function and environment configuration.

Literature An example of a problem similar to the one presented in this thesis, solved using Reinforcement Learning, can be found in the paper ‘DeliverAI: a distributed path-sharing network based solution for the last mile food delivery problem’ by Ashman et al. [20].

They aimed at solving the last-mile delivery problem by developing a distributed path-sharing network based on Reinforcement Learning. Their approach uses a multi-agent system to optimize delivery routes and schedules, considering factors such as traffic congestion, delivery time windows, and vehicle capacity.

However, their model simplifies the real-world scenario by assuming fixed delivery locations and known traffic patterns, which may not accurately reflect the dynamic and uncertain nature of real-world logistics environments. Moreover, their approach requires extensive training and tuning to achieve optimal performance.

2.4.3 Planning with LLM

LLMs are trained on vast amounts of textual data and have demonstrated remarkable performance across a wide range of language tasks, from translation and summarization to reasoning and problem-solving (emerging abilities). This success has naturally led researchers to explore whether these models can be repurposed for more complex, multi-step decision-making problems that require planning.

The key idea is that the same abilities that allow LLMs to understand and generate language can be harnessed to decompose a planning task into intermediate steps, reason about the consequences of actions, and even generate entire action sequences with minimal or no task-specific training.

2.4.3.1 Chain-of-Thought Reasoning

One of the most influential ideas for using LLMs in planning is Chain-Of-Thought prompting. Instead of asking the model to jump directly from a problem statement to a final answer, CoT prompting encourages the model to “think aloud” by generating intermediate reasoning steps. This decomposition can help in planning problems where the solution involves multiple, logically connected steps.

This was introduced by Wei et al. [33], who demonstrated that prompting the LLM to ‘answer step by step’ led to improved performance on mathematical problems compared to requesting only the final answer. They also showed that this step-by-step approach could be applied to other fields, ultimately giving rise to Chain-of-Thought reasoning models.

“Reasoning” Models Reasoning-focused LLMs are trained to generate multiple Chain-of-Thought steps, exploring different solution paths before selecting the most optimal one, often using Reinforcement Learning [6] techniques such as Reinforcement Learning from Human Feedback or self-consistency methods during the training.

This approach enhances both accuracy and explainability (to some extent), as the model articulates its reasoning process while still operating as a generative AI system. Expanding on this concept, reasoning models can integrate external tools, memory, and API calls, forming what is commonly referred to as an **LLM Agent**, capable of autonomous decision-making and real-world interaction.

Most recent and famous reasoning models released to the public have been developed by many companies, both big and small, such as:

⁸<https://it.mathworks.com/help/reinforcement-learning/ug/create-agents-for-reinforcement-learning.html>

- **OpenAI:** o1⁹, o1-mini¹⁰ and o3-mini¹¹ are reasoning models designed to enhance logical problem-solving capabilities. o1 is specialized in complex problems across various domains, offering robust reasoning skills. Building upon this foundation, o3-mini provides a more cost-effective and faster alternative;
- **DeepSeek:** DeepSeek-R1¹², is a notable AI model from a startup¹³. Released in early 2025, DeepSeek-R1 is recognized for its powerful reasoning and coding skills, achieved at a fraction of the development cost compared to other leading models. Its open-source nature and efficiency have made it a significant player in the AI landscape.

2.4.3.2 Zero-Shot and Few-Shot Planning

In zero-shot planning, LLMs generate action sequences by utilizing their extensive pretraining on text and code, effectively inferring plausible step-by-step solutions to given tasks. Few-shot planning further enhances their capabilities by providing LLMs with a small set of demonstrations, enabling them to generalize patterns and improve their action sequencing capabilities.

However, while LLMs can produce reasonable plans, their direct applicability to embodied environments remains challenging. Huang et al. [12] in ‘Language Models as Zero-Shot Planners: Extracting Actionable Knowledge for Embodied Agents’ highlight the limitations of naive LLM planning, noting that LLMs struggle with real-world constraints, action feasibility, and long-horizon dependencies. Their work demonstrates that these shortcomings can be mitigated by leveraging the world knowledge embedded within LLMs and applying structured guidance, such as constraints on action generation and feedback-based refinements.

Similarly, Silver et al. [30] in ‘PDDL Planning with Pretrained Large Language Models’ extend this inquiry to classical AI planning domains by evaluating few-shot prompting of LLMs on problems expressed in the Planning Domain Definition Language (PDDL). Their findings reveal mixed results: while LLMs can generate syntactically correct PDDL plans in certain domains, they often fail due to a lack of explicit access to transition models and logical constraints inherent to planning problems. Nonetheless, their study also introduces a hybrid approach where LLMs are used to initialize heuristic-based search planners, demonstrating that even imperfect LLM-generated plans can improve the efficiency of traditional AI planning methods.

These findings collectively suggest that while LLMs alone are not yet fully capable of robust autonomous planning, their ability to extract and apply commonsense knowledge makes them valuable tools for augmenting structured planning frameworks. By integrating LLM-generated outputs with classical search-based methods, researchers have shown improvements in planning efficiency and problem-solving robustness, highlighting a promising direction for future research at the intersection of language models and automated planning.

Literature In the paper ‘Exploring and Benchmarking Planning Capabilities of Large Language Models’ by Bohnet et al. [3], the authors systematically analyze the planning capabilities of LLMs through a novel benchmarking suite that includes both classical planning tasks (expressed in PDDL) and natural language-based planning problems. Their work highlights the limitations of LLMs in planning, particularly their tendency to generate suboptimal or incorrect plans despite their strong language understanding capabilities. To address these shortcomings, they explore various methods to improve LLM-based planning (including many-shot in-context learning, fine-tuning with optimal plans, and the use of chain-of-thought reasoning techniques such as Monte Carlo Tree Search (MCTS) and Tree-of-Thought (ToT)). The results indicate that, while LLMs struggle with planning in zero-shot and few-shot settings, their performance significantly improves when provided with structured demonstrations and reasoning strategies. Moreover, fine-tuning on high-quality plan data leads to near-perfect accuracy in some cases, even with relatively small models. However, challenges remain in

⁹<https://openai.com/o1/>

¹⁰<https://openai.com/index/openai-o1-mini-advancing-cost-efficient-reasoning/>

¹¹<https://openai.com/index/openai-o3-mini/>

¹²<https://github.com/deepseek-ai/DeepSeek-R1>

¹³<https://www.deepseek.com/>

out-of-distribution generalization, where models fail to generalize effectively to novel scenarios without additional training. Their analysis also identifies key failure modes in LLM planning, such as constraint violations, failure to reach goal states, and incorrect action sequences, emphasizing the need for better training data curation and reasoning frameworks.

Literature In ‘Generalized Planning in PDDL Domains with Pretrained Large Language Models’ by Silver et al. [29], the authors investigate whether LLMs, specifically GPT-4, can serve as generalized planners, not just solving a single planning task, but synthesizing programs that generate plans for an entire domain. They introduce a pipeline where GPT-4 is prompted to summarize the domain, propose a general strategy, and then implement it in Python. Additionally, they incorporate automated debugging, where GPT-4 iteratively refines its generated programs based on validation feedback. Their evaluation on seven PDDL domains demonstrates that GPT-4 can often generate efficient domain-specific planning programs that generalize well from only a few training examples. The study also finds that automated debugging significantly improves performance, while the effectiveness of Chain-of-Thought summarization is domain-dependent. Notably, GPT-4 outperforms previous generalized planning approaches in some cases, particularly when leveraging semantic cues from domain descriptions. However, limitations remain, especially in handling domains requiring deeper structural reasoning or non-trivial search processes.

3 Experiment Setting

In this chapter, we provide a comprehensive and in-depth description of the experimental framework designed to evaluate the performance of our LLM-driven agent.

We begin by formally defining the problem, ensuring that our study is framed within a well-structured and precise context. We also outline the specific aspects of the problem that our research aims to investigate, clarifying our objectives and highlighting the choices we made in our work.

Following this, we offer a thorough explanation of the environment used to simulate the delivery platform; this section provides a detailed overview of the web-based system that serves as the operational space for our agent. We describe the structure of the platform, its key features, and how it functions as a testbed for evaluating autonomous agents.

Finally, we discuss the selection of various Large Language Models used in our experiments, including both the models that were actively tested and those that were considered but ultimately not included in our evaluations.

3.1 Problem Definition

As widely explained in Section 2.4.3, the recent advancements in Large Language Models have demonstrated their impressive capabilities across a wide range of tasks. Their ability to process and reason about complex problems opened new avenues for research, particularly in fields such as planning and logistics. Given the power and versatility of these models, we are motivated to further explore their potential in tackling planning and logistic challenges, evaluating their ability to comprehend and solve such problems autonomously.

In this work, our primary focus is on assessing the inherent strengths and weaknesses of LLMs when used in their raw form, without integrating any additional planning frameworks, heuristic search algorithms, or explicit reasoning mechanisms on top of them. Unlike conventional approaches that rely on dedicated pathfinding algorithms, rule-based systems, or carefully structured reinforcement learning paradigms, our objective is to investigate how well an LLM can independently interpret and navigate a logistic scenario using its generative abilities alone.

One of the key aspects we wish to emphasize is that our approach remains purely generative. In other words, rather than embedding domain-specific logic or fine-tuned strategies within the model, we allow the LLM to operate autonomously, generating its own understanding of the environment and devising its own strategies for completing the given tasks.

3.1.1 Our Task

Specifically, our problem formulation is centered around asking the LLM to provide only the *next step* that moves the agent closer to the goal, rather than generating an entire solution at once. This step-by-step approach also enables the model to iteratively refine its path based on new observations using the conversation history as “action-result” feedback. Furthermore, we assess the reliability of each generated step by computing the uncertainty of the model’s response using the methodology detailed in Section 2.2.2.3.

By taking this approach, we aim to answer these questions about the problem-solving skills of LLMs in logistics problems:

- To what extent can an agent, powered solely by an LLM, solve a logistic problem when placed in an unexpected and unfamiliar environment?
- What are the intrinsic limitations and strengths of this approach compared to traditional rule-based or algorithmic solutions?

To simulate an unexpected and dynamic environment, we designed our experiments around a web-based platform that interacts with the agent through API calls. The platform provides a structured yet unpredictable setting in which the agent must operate. A critical design choice we made in our methodology was to avoid parsing the JSON response containing the map structure. Instead, the agent receives the map data (that is added to the prompt) and is expected to interpret it entirely on its own. This decision was made to ensure that the LLM must independently derive the necessary spatial and logistical information without relying on pre-processed or structured inputs.

Additionally, this design choice introduces a layer of robustness: if the API undergoes modifications, such as changes in the response format, the addition of new parameters, or variations in data structure, the agent should still be capable of functioning. This property aligns with our objective of evaluating the adaptability of LLM-driven agents in dynamically changing environments, where real-world conditions may not always remain constant.

Our experimental results will be presented in detail in Chapter 6. However, to summarize our primary evaluation criteria, we focus on testing the following goals of the LLM-based agent:

- **Parcel Pickup:** We evaluate whether the agent is capable of successfully identifying the correct location of a parcel on the map and navigating to that specific tile to pick it up. This task requires the agent to correctly interpret spatial relationships and make movement decisions accordingly;
- **Parcel Delivery:** The second evaluation criterion involves determining whether the agent can correctly identify and reach the intended delivery location based on the information available in the raw map data. Since no explicit delivery coordinates are pre-processed for the agent, it must infer this information on its own.

Through these experiments, we aim to provide valuable insights into the generative problem-solving capacity of LLMs in a logistic setting, evaluating their adaptability, reasoning limitations, and potential advantages in real-world scenarios.

3.2 Environment - Deliveroo.js

Deliveroo.js it's an Educational Game, developed by Marco Robol for the course on Autonomous Software Agents by Prof. Paolo Giorgini, using the Treejs¹ framework.

The code for the server is open and can be accessed on GitHub² as well as some example of agents (with different level of complexity)³.

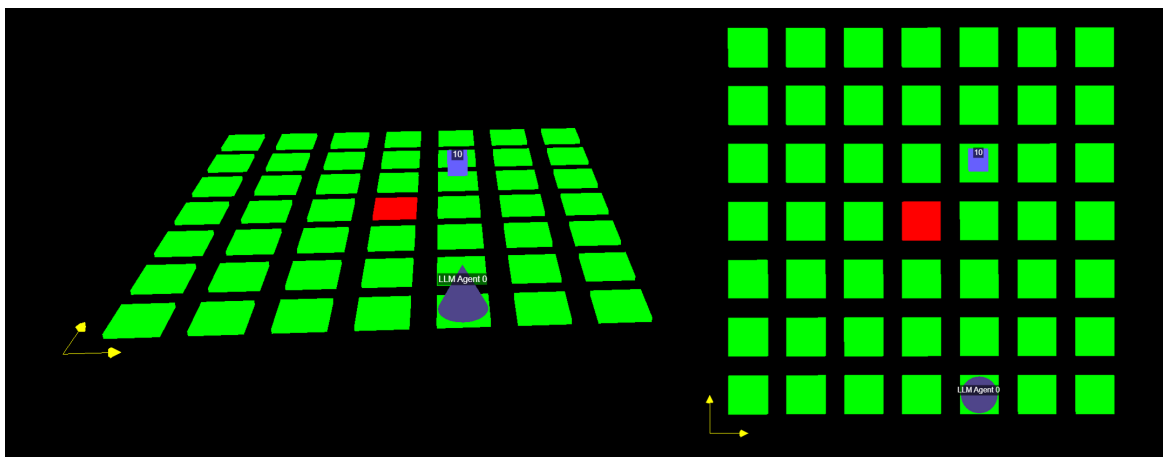


Figure 3.1: Two views of the same 7×7 Map in Deliveroo.js with a single central delivery zone.

The game can be played even by humans, by interacting in the browser; technically speaking, it is a web-based platform consisting of three main components connected to each other via sockets

¹<https://threejs.org/>

²<https://github.com/unitn-ASA/Deliveroo.js>

³<https://github.com/unitn-ASA/DeliverooAgent.js>

(implemented with Socket.io⁴):

- **Game server:** it contains the entire logic of the game and it includes the implementation of client connection handler, parcel spawning, current environment status and so on;
- **Agent client:** it is the custom component that we developed to interact with the game server. It is a JavaScript file that connects to the server, manages all the logic of the agent (in our case, the LLM agent) and sends the actions to the server;
- **3D web app:** it is the visual representation of the game. It is a web page that connects to the server and receives the status of the game to render it in a 3D environment. It is not necessary for the agent to work, but it is useful to understand what is happening in the game.

As we can see from the Figure 3.1, the game is a grid of $N \times M$ tiles where the agent can move. At the time of writing, the (0,0) cell is in the bottom-left corner, but we moved it to the top-left corner to better analyze the results with additional tools. The map is defined in a JS file (example in Listing 3.1) where the number inside a cell represent the type of the cell.

JavaScript Code

```
1 // 7x7 goal center deliver
2 module.exports = [
3   [1, 1, 1, 1, 1, 1, 1],
4   [1, 1, 1, 1, 1, 1, 1],
5   [1, 1, 1, 1, 1, 1, 1],
6   [1, 1, 1, 2, 1, 1, 1],
7   [1, 1, 1, 1, 1, 1, 1],
8   [1, 1, 1, 1, 1, 1, 1],
9   [1, 1, 1, 1, 1, 1, 1],
10  ];
```

Listing 3.1: Example of a 7x7 map with a single central delivery zone

There are three possible types of cells:

- **green** (1): the agent can move on it. These cells can contain multiple parcels but only one agent at a time;
- **red** (2): the agent can move on it and deliver any number of parcel it has;
- **black** (0): the agent can't move on it and they can't contain any parcel (we will not use them in our tests).

The functioning is very straight forward:

- **Agents:** there can be any number of agents that can cooperate or compete. Each agent has a score that is increased by delivering parcels. They are represented as cones with their name on it on the map ('LLM Agent' in Figure 3.1 is ours);
- **Parcels:** they are represented as small cubes with a number on it. The number is the reward the agent will get by delivering it. They spawn in random cells and they can be picked up by the agent. If they are not delivered in a certain amount of time, they may disappear.

⁴<https://socket.io/>

3.2.1 Server Configuration and Event Handling

JavaScript Code

```
1 module.exports = {
2   MAP_FILE: "map_file",
3
4   PARCELS_GENERATION_INTERVAL: "5s",
5   PARCELS_MAX: "1",
6
7   MOVEMENT_STEPS: 1,
8   MOVEMENT_DURATION: 50,
9   AGENTS_OBSERVATION_DISTANCE: 100,
10  PARCELS_OBSERVATION_DISTANCE: 100,
11  AGENT_TIMEOUT: 100,
12
13  PARCEL_REWARD_AVG: 10,
14  PARCEL_REWARD_VARIANCE: "0",
15  PARCEL_DECAYING_INTERVAL: "infinite",
16
17  RANDOMLY_MOVING_AGENTS: 0,
18  RANDOM_AGENT_SPEED: "2s",
19
20  CLOCK: 50,
21 };
```

Listing 3.2: Example of a configuration file for the server

The behavior of the parcels in the system is defined through the server configuration file. This file specifies key parameters that control parcel generation, reward values, and decay over time. An example for a configuration file is shown in Listing 3.2.

Based on the server settings, a maximum number of parcels can be active simultaneously. Each parcel is spawned at a fixed interval, with a random reward value determined by a specified average and variance. Additionally, the configuration dictates whether the reward remains constant or decreases over time.

In the example, parcels are generated every 5 seconds, but only one can exist at a time. Each parcel starts with a reward value of exactly 10. Furthermore, since `PARCEL_DECAYING_INTERVAL` is set to `"infinite"`, the reward does not decrease over time and the parcel will not disappear (until delivered). This setup ensures a stable environment for testing the agent's performance.

The agent can interact with the environment using the following actions:

- **up, down, left, right:** move in the specified direction, if the relative cell is empty and green or red;
- **pickup:** the agent picks up the parcel present in the cell;
- **deliver:** the agent drops any parcel it has in the cell: if it is a delivery zone the parcels will disappear and the reward will be added to the player's score, otherwise it will just remain on the cell.

The server is responsible for transmitting events to the agent, ensuring that it receives all relevant updates in real-time. Specifically, the following events were utilized in our tests:

- **onMap (width, height, tiles):** it sends the width and the height of the map, along with all the tiles in it. Tiles are currently sent as a dictionary `{x: INT, y: INT, delivery: BOOL,`

`spawner: BOOL`} where `delivery` is true if the tile is a delivery zone and `spawner` is true if a parcel can spawn on it;

- `onYou (id, name, x, y, score)`: it sends the id, the name, the x and y coordinates and the score of the agent connected that the code is piloting;
- `onParcelsSensing (perceived_parcels)`: it is an async function that sends the parcels that the agent can see at any time. The parcels are sent as a dictionary `{x: INT, y: INT, reward: INT}` where x and y are the coordinates of the parcel and `reward` is the reward the agent will get by delivering it.

3.3 Large Language Models Selection

As mentioned in Section 2.2, Large Language Models are powerful tools that have revolutionized the field of natural language processing. The way they generate text based on input prompts has opened up new possibilities for research and applications in various domains.

One of the core aspects of LLMs is their autoregressive nature, meaning they generate text one token at a time (given the current implementation, but alternative solution are currently being studied, for example by Meta AI⁵ in the paper ‘Better & Faster Large Language Models via Multi-token Prediction’ by Gloeckle et al.[8]), predicting the next most likely token based on the context provided. This capability is what allows LLMs to generate coherent and contextually relevant responses. The way these systems operate can be broken down into a (simplified) step-by-step process:

- The prompt (the request) is tokenized, which means it is divided into smaller units called tokens. These tokens are predefined character combinations that serve as the building blocks for processing text. An example of this tokenization process is illustrated in Figure 3.2. Once tokenized, the prompt is passed to the model for processing;
- The model then generates the “next token” based on a probability distribution computed using attention mechanisms, to determine the most contextually appropriate next token;
- The newly generated token is appended to the existing sequence, and the process is repeated iteratively. This continues until a predefined stopping criterion is met, either reaching a maximum token limit or encountering a special termination token.

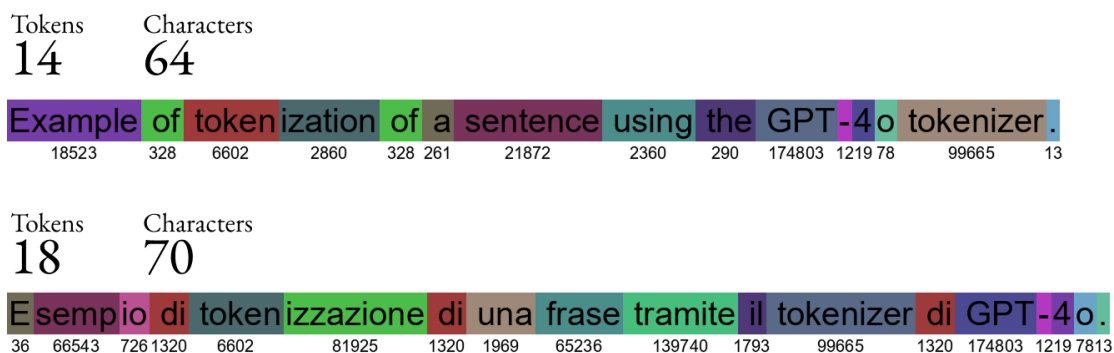


Figure 3.2: Example of tokenization of a sentence using the GPT-4o tokenizer.

Source: Data from OpenAI Platform⁶

One key element of this process is how the next token is selected. The winning token is picked from a probability distribution obtained through the softmax function. However, if selection were purely deterministic, the model would always generate the same output given the same prompt, making it

⁵<https://ai.meta.com/>

⁶<https://platform.openai.com/tokenizer>

rigid and predictable. To introduce variability and prevent repetitive patterns, a controlled amount of randomness is introduced using the `temperature` parameter.

The temperature parameter plays a crucial role in regulating randomness in text generation: this mechanism explains why the same prompt can yield different outputs when used multiple times: because the token selection is influenced by this controlled randomness.

Beyond temperature, another factor that influences token selection is the logit bias⁷. Logit bias allows direct intervention in the probability of specific tokens being chosen during text generation. Instead of relying solely on the model's learned probabilities, users can manually adjust the likelihood of certain tokens appearing by modifying the logits (the unnormalized probabilities) before applying the softmax function.

The logit bias mechanism operates as follows:

- Positive bias values increase the probability of a specific token being selected, making it more likely to appear in the generated text.
- Negative bias values decrease the probability of a token, potentially eliminating it from consideration altogether.

This approach gives users more control over text generation, allowing them to guide the model toward preferred outputs while avoiding undesired words or phrases. A simplified implementation of how logit bias works can be seen in the Python code snippet in Listing 3.3, where we want to increase the likelihood of the word “fox” and decrease the likelihood of the word “quick”.

Python Code

```
1 [...]
2
3 # Get the logits (raw predictions)
4 outputs = model(**inputs)
5 logits = outputs.logits
6
7 logit_bias = {
8     # Decrease likelihood for the word "quick"
9     tokenizer.encode("quick")[0]: -2.0,
10    # Increase likelihood for the word "fox"
11    tokenizer.encode("fox")[0]: 2.0,
12    # Almost never generate the word "slow"
13    tokenizer.encode("slow")[0]: -100.0,
14    # Almost always generate the word "brown"
15    tokenizer.encode("brown")[0]: 100.0,
16 }
17
18 # Apply logit bias: modify logits of specific tokens
19 for token_id, bias in logit_bias.items():
20     logits[token_id] += bias
21
22 # Convert logits by applying softmax
23 probs = torch.nn.functional.softmax(logits, dim=-1)
24
25 [...]
```

Listing 3.3: Example of what a basic implementation of logit bias could look like

⁷<https://www.vellum.ai/llm-parameters/logit-bias>

This technique is particularly useful in scenarios where the generated text must adhere to specific constraints. For example, logit bias can be used for content moderation, ensuring that the model avoids generating harmful, offensive, or inappropriate content. By assigning a strong negative bias to certain tokens used to build specific words or sentences, users can effectively steer the model away from producing undesirable responses.

Not only, there are some cases where an LLM has been expanded with custom and private information, via fine-tuning or Retrieval Augmented Generation (as we cited in Section 2.2.2). In this context, we may want to force the model to not generate some specific content that may be useful for the overall response but should not be shared.

On the other hand, logit bias can be leveraged to override built-in model safeguards. In some cases, AI models have safety mechanisms that prevent them from answering certain types of questions, responding with phrases like “I’m sorry, but I can’t provide that information.” By applying a negative logit bias to the tokens that generate the words that build this response, users can force the model to produce an alternative reply, whether it be a reworded refusal or even an answer to the original question.

Overall, logit bias is a powerful tool for modelling model behavior, allowing developers to enforce preferred linguistic patterns, avoid specific terminology, and customize AI responses according to their needs. When combined with temperature adjustments and other generation techniques, it provides a robust framework for controlling LLM output and ensuring its alignment with desired objectives.

3.3.1 Open Source Models

The term “Large” in Large Language Models refers to their extensive number of parameters and the vast datasets used during training. Training these models is a resource-intensive process, both in terms of computational power and time. However, some organizations choose to release their models as open source, allowing the community to access and utilize them freely. This approach is also beneficial to the broader AI community, as it fosters innovation: if a new open source model is released and it’s “the most powerful”, it sets a new baseline for companies. This compels the companies offering paid access to their closed models to find further innovations, either by reducing costs or developing even more powerful models with new features.

In the context of LLMs, the term “open source” differs from its traditional usage in software development⁸. Specifically, there is a nuanced distinction that categorizes publicly available models into two primary types:

- **Open Source:** The creators provide full access to the model’s source code, architecture, training data, and pre-trained weights. This level of transparency allows users to understand, modify, and enhance the model comprehensively.
- **Open Weights:** The creators make the trained model’s weights publicly available but may withhold other critical components, such as the training data or detailed methodologies used during training. This approach enables users to employ the model for specific tasks but limits their ability to fully comprehend or modify its underlying structure.

In the following subsections, we will present some open models we tested in the early stages of our research. The results analyzed in Chapter 6 do not consider these models due to certain limitations, which will be discussed shortly. Nonetheless, they are worth mentioning as they provided a valuable starting point during the initial phases of our project.

3.3.1.1 Challenges with Open Source Models

One significant challenge we encountered with open source models was the implementation of altering the logit bias mechanism. While the example in Listing 3.3 may suggest simplicity, the reality is more complex. Implementing this mechanism requires:

- Reconstructing the model’s architecture accurately;
- Loading the pre-trained weights appropriately;

⁸<https://promptengineering.org/llm-open-source-vs-open-weights-vs-restricted-weights/>

- Modifying the model’s potentially intricate structure to extract the raw values and change them to incorporate the logit bias.

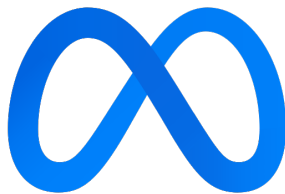
These steps demand substantial time, and in addition, running those models locally requires large computational resources; it would have been too much effort for something that was not the focus point of the project, since closed models provide this feature out of the box.

Moreover, during initial tests, we observed difficulties in constraining the Open Source models to produce specific tokens from a list without altering the logit bias. For instance, when prompted to “Answer with just a single letter between U, D, L, R.” (the explanation for this prompt will be given in Section 5.2.4) the model often responded with full sentences like “Sure, the answer is U!”. Truncating such responses to a single token would result in outputs like “Sure” which is not the desired outcome. This is a problem that we didn’t have with the closed source models, that we will discuss in the next sub-section 3.3.2.

So, open source models have been tested only in the first instance of the project, to test the logic of the agent without wasting credit with the OpenAI API.

They have been loaded and used through Ollama⁹, a tool for running and managing large language models locally. It simplifies downloading, running, and interacting with models without relying on cloud services.

3.3.1.2 LLaMa 3.2

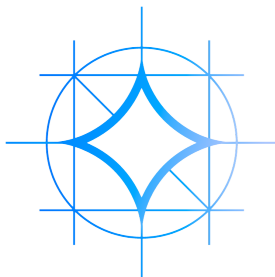


Source: Flaticon

LLaMa 3.2, developed by Meta AI¹⁰, is a multimodal large language model designed to process both textual and visual data, marking Meta’s first open-source AI model with such capabilities. The model is available in two configurations: an 11-billion-parameter version and a more robust 90-billion-parameter variant. There are also a 1-billion-parameter and 3-billion-parameter text-only versions of the models. These models are optimized for deployment on mobile (yet powerful) hardware platforms.

The model tested in our experiments was the *text-only 3-billion-parameter* version.

3.3.1.3 Gemma 2



Source: logowik

Gemma 2, introduced by Google¹¹, is an open suite of language models available in parameter sizes of 2 billion, 9 billion, and 27 billion. The 27-billion-parameter model has demonstrated exceptional performance, surpassing larger models in real-world conversational benchmarks. This suite is built upon the same research and technology that underpins Google’s Gemini models, emphasizing both performance and accessibility.

The model tested in our experiments was the *9-billion-parameter* version.

3.3.1.4 DeepSeek-V3



Source: DeepSeek GitHub

DeepSeek-V3, developed by DeepSeek¹², is a Mixture-of-Experts (MoE) language model comprising a total of 671 billion parameters, with 37 billion activated per token during inference. It employs Multi-head Latent Attention (MLA) and the DeepSeekMoE architecture to achieve efficient inference and cost-effective training. The model was pre-trained on 14.8 trillion diverse and high-quality tokens, followed by supervised fine-tuning and reinforcement learning stages to fully harness its capabilities. Despite its extensive scale, DeepSeek-V3’s training process is notably efficient and it has demonstrated remarkable stability throughout training. Comprehensive evaluations reveal

⁹<https://ollama.com/>

¹⁰<https://ai.meta.com/blog/llama-3-2-connect-2024-vision-edge-mobile-devices/>

¹¹<https://developers.googleblog.com/en/gemma-explained-new-in-gemma-2/>

¹²<https://github.com/deepseek-ai/DeepSeek-V3>

that DeepSeek-V3 outperforms other open-source models and achieves performance comparable to leading closed-source models.

Unfortunately, we didn't get to test it because it has been released after the end of the project, but it is worth mentioning for future research thanks to its impressive performance.

3.3.2 Closed Source Models - OpenAI

Before delving into the specifics of this section, we want to emphasize that OpenAI is not the only company offering closed-source AI models. Several other providers exist, such as Anthropic¹³, which develops the Claude family of models, but also DeepSeek¹⁴, that developed the open source DeepSeek-V3 provides paid API access to its models as well.

While OpenAI remains the most widely recognized and utilized provider, particularly in research and industry applications, it is important to acknowledge that many of the benefits and drawbacks of closed-source models apply broadly across all such services. Nevertheless, since our work specifically relies on OpenAI's models, our discussion will primarily focus on them while keeping in mind that similar considerations extend to other closed-source alternatives.

A key distinction between open-source and closed-source models is the transparency regarding their architecture. In many closed-source models, details such as the exact number of parameters are not publicly disclosed. For instance, while OpenAI's GPT-3[4] is known to have a maximum of 175 billion parameters¹⁵, the parameter counts for subsequent models like GPT-4[22] have not been officially confirmed. Estimates suggest that GPT-4 may have around 1.76 trillion parameters, but this remains speculative¹⁶.

OpenAI was established as a research organization dedicated to advancing artificial intelligence. They initially released models such as GPT-2[25] and Whisper[23] (a speech recognition model) to the public at no cost. GPT-2, for example, was made available with various model sizes, the largest containing 1.5 billion parameters¹⁷.

Subsequently, OpenAI developed GPT-3 and DALL·E[26], introducing them through a commercial platform and API services. The landscape of AI applications shifted significantly with the release of ChatGPT¹⁸, a conversational AI model that garnered widespread attention for its advanced language understanding and generation capabilities.

OpenAI offers access to their models via subscription plans and a pay-as-you-go API, with pricing varying based on the specific model utilized. The API provides granular control over model behavior through parameters such as `logit_bias`, that, as the name suggests, allows users to adjust the likelihood of specific tokens appearing in the generated output by assigning bias values ranging from -100 to 100.

However, a limitation of closed-source models is the lack of transparency regarding updates or changes. Users may be unaware of modifications that could affect model behavior. For instance, there have been instances where the behavior of the `logit_bias` parameter changed without prior (or "at run-time") notice. We have identified three primary behaviors that the `logit_bias` parameter can have depending on the version of the API, which appears to be independent from any update of the models themselves:

1. **No change:** The API does not apply the `logit_bias` at all, resulting in outputs identical to those generated without using the parameter;
2. **Exclusive:** The `logit_bias` is applied strictly. Setting a token's bias to 100 forces the model to produce that token; if multiple tokens are set to 100, the model will only choose among them. Conversely, setting a token's bias to -100 prevents that token from appearing in the output;



Source: Vecteezy

¹³<https://www.anthropic.com/>

¹⁴<https://api-docs.deepseek.com/>

¹⁵<https://en.wikipedia.org/wiki/GPT-3>

¹⁶<https://en.wikipedia.org/wiki/GPT-4>

¹⁷<https://en.wikipedia.org/wiki/GPT-2>

¹⁸<https://en.wikipedia.org/wiki/ChatGPT>

3. **Soft:** The `logit.bias` is applied moderately. Assigning a token a bias of 100 significantly increases its likelihood of being produced, but other tokens may still be selected. Similarly, setting a token’s bias to -100 greatly decreases its likelihood, but it may still appear.

Our tests were conducted during a period when the API exhibited the third behavior (Soft), while the original paper ‘Robots That Ask For Help: Uncertainty Alignment for Large Language Model Planners’ was released, the API exhibited the second behavior (Exclusive)

Due to frequent API updates, we cannot guarantee the reproducibility of these results. To more mitigate variability, we set the temperature parameter to zero during our tests, aiming for deterministic outputs. To limit the impact of unwanted tokens appearing in the output or vice versa, we decided to ignore any token that was not in the list of expected tokens and we assigned to the missing tokens the 20th log-probability minus the 19th one.

Another crucial feature offered by OpenAI is the `logprobs` parameter. When specified, this parameter returns the log probabilities of the top tokens that the model may generate as the “next token”.

This information is essential for computing the uncertainty of the model’s predictions, as detailed in Section 2.2.2.3. It’s important to note that also the behavior of the `logprobs` parameter can change over time. For instance, when the paper ‘Robots That Ask For Help: Uncertainty Alignment for Large Language Model Planners’ was released, the `logprobs` parameter returned a maximum of 5 tokens. Since then, this limit has been increased to a maximum of 20 tokens. In the future, it may be adjusted again, so it’s important to keep this in mind when working with the API.

Pricing per 1 million tokens for all the OpenAI models used can be seen in Table 3.1; the functioning of caching is explained in Section 5.2.5.1.

Model	Input	Cached Input	Output
gpt-4o-mini	\$0.15	\$0.075	\$0.60
gpt-4o	\$2.50	\$1.25	\$10.00
gpt-3.5-turbo	\$0.50	/	\$1.50

Table 3.1: Pricing Table per 1Mil tokens for the tested models.

Source: OpenAI ¹⁹

3.3.2.1 GPT-4o-mini

GPT-4o-mini, is a lightweight variant of the GPT-4o architecture, optimized for efficiency while maintaining strong performance across a range of tasks. It is designed to deliver fast inference and lower computational costs, making it suitable for deployment in real-time applications and on-device AI systems. While OpenAI has not publicly disclosed the parameter count, it is positioned as a more efficient alternative to larger models in the GPT-4 family.

The main model used for the tests in our experiments was indeed GPT-4o-mini, selected for its balance of performance and pricing.

3.3.2.2 GPT-4o

GPT-4o is OpenAI’s flagship multimodal model, capable of processing and generating text, images, and audio in real time. It represents a significant leap in efficiency and latency, outperforming previous iterations such as GPT-4-turbo while operating at a lower computational cost. Unlike its predecessors, GPT-4o is natively trained across multiple modalities rather than combining separate models for different inputs. Though OpenAI has not released detailed architectural specifications, benchmark results indicate substantial improvements in reasoning, multilingual proficiency, and response speed.

It has been used less than GPT-4o-mini because of the higher cost and the fact that, since from the first tests, the results were almost identical to the mini version.

¹⁹<https://platform.openai.com/pricing>

AzureAPI GPT-4o was accessed via the Azure OpenAI API, provided by the University of Trento, offering a secure and private interface for interacting with the model. To facilitate communication between the agent’s JavaScript code and the API, a lightweight Python server was developed. The server’s code is available on GitHub[5].

3.3.2.3 GPT-3.5-turbo

GPT-3.5-turbo is a high-performance variant of OpenAI’s GPT-3.5 model. It provides strong general-purpose capabilities while being more accessible for applications requiring large-scale deployment. Though it does not match GPT-4-level reasoning abilities, it remains competitive in many NLP benchmarks and is widely used for production AI services.

Unfortunately, this model has been the weakest in our tests; we will discuss this in depth in the Section 6.7.

4 Agent Development

In this chapter we present the iterative development of the agent. We describe the successive phases of its evolution, from the initial prototype to the final implementation. During this process, several challenges and unexpected issues emerged. For each phase, we detail the encountered problems and the approaches adopted to overcome them, explaining the reasoning behind the design choices. The majority of the choices taken in the structure of the final prompts will be discussed in Chapter 5, while here are reported the ones from the discarded approaches.

4.1 Development information

The go-to programming language for AI development is Python. However, we chose to use JavaScript for this project for several reasons. First, both the server and the example agents were already implemented in JavaScript; so a JavaScript interface to communicate with the server was already available. This allowed us to focus only on the agent’s logic, without having to worry about the server-side implementation. Second, thanks to the availability of the project of the Autonomous Software Agents course, our initial plan was to use an existing JavaScript-based benchmark from the course. Although we ultimately decided against using this benchmark, so JavaScript remained our language of choice.

Additionally, since there isn’t a dedicated JavaScript library for the Azure OpenAI API, instead of manually recreating the necessary API calls we opted for a more efficient approach by setting up a lightweight Python server to act as a middleman. This solution, discussed in Section 3.3.2.2, allowed us to integrate Azure’s services seamlessly while maintaining JavaScript for the core development.

4.2 First Approach

In this initial phase of the development and testing process, a trial-and-error methodology was adopted to iteratively refine the system’s behavior and try to optimize performance. Unfortunately, this led to moving away from the definition of the problem and, in combination with the poor performance of the agent, it was decided to start over with a new approach. Nonetheless, this first attempt was crucial in understanding the challenges and limitations of the problem, so it is important to describe it.

Text

```
1 [ROLE]
2
3 [MAP]
4
5 [LEGEND]
6
7 [ACTIONS]
8
9 [PARCELS ALREADY PICKED]
10
11 [RULES]
12
13 [QUESTION]
```

Listing 4.1: Scheme of a prompt used in the first approach, full in Appendix B.2

The approach began by parsing crucial information from the server, which served as the foundation

for understanding how the LLM would interact with it. The main point of discussion in the parsing topic was the map, that was represented as a multi-line string like the one in Listing 4.2, even if LLMs analyze the data as an “horizontal” sequence. The first prompt sent to the LLM was crafted by concatenating the map with all the other information needed to describe the state of the environment, as shown in Listing 4.1.

Text

```

1      1  1  1  1  1
2      1  1  P  1  1
3      1  1  1  1  1
4      2  A  1  1  1
5      1  1  1  1  1
6
7      LEGEND:
8      1: Walkable cell
9      2: Delivery point
10     A: Agent
11     P: Parcel

```

Listing 4.2: Parsed Map Result with partial legend

This implementation started by letting the LLM also decide the goal to pursue, by providing all the available information. As various challenges and inefficiencies were identified during extensive testing, we progressively implemented a total of seven “helping” parameters.

4.2.1 Helping Parameters

These parameters were introduced with the objective of addressing specific issues observed during experimentation, and each of them played a significant role in shaping the overall functionality of the agent:

- **ANTI_LOOP:** This parameter was introduced to eliminate a common inefficiency in agent movement, where the agent would repeatedly traverse the same path in a circular loop, failing to make meaningful progress toward its goal. By setting this parameter to `true`, if the last four action were ["U", "R", "D", "L"] (either clockwise or counterclockwise) the agent was forced to take an action that prevented the loop from happening again. This optimization helped the agent make more intelligent movement decisions, thereby removing the possibility of being stuck in repetitive cycles;
- **HELP_THE_BOT:** The primary purpose of this parameter was to assist the agent in handling parcels more effectively. When activated by setting it to `true`, the agent was programmed to automatically take a parcel if the parcel was located directly below its current position. Additionally, if the agent was positioned at a delivery point, this parameter ensured that the agent would immediately proceed with delivering the parcel without requiring additional decision-making steps. This was implemented to reduce the number of calls to the LLM, since, even in this version of the agent, it was able to always pick up a parcel and deliver it (if in the correct tile);
- **SELECT_ONLY_ACTION:** This parameter was designed to simplify the agent’s decision-making process in cases where the list of available actions contained only a single option. When set to `true`, the agent would automatically select and execute the sole available action without calling the LLM. This was made by a big filtering phase that returned the legal actions:
 - remove the opposite of the last action;
 - remove all the actions that were not possible (like going left while in the first column or going up while in the first row);
 - remove the delivery action if the agent wasn’t carrying a parcel and in a delivery point;

- remove the pick action if the agent was in a cell with no parcel.

This, in combination with the `HELP_THE_BOT` parameter, reduced the number of unnecessary calls to the LLM, thereby enhancing the agent’s efficiency, but also giving the agent too little decision power;

- **USE_HISTORY:** This parameter is the only one that was kept for all the future iterations (more on this in Section 6.3). The role of this parameter was to decide whether each call to the LLM should contain only the current state of the environment of the entire message history. If set to `true`, the LLM would have access to the full conversation history, allowing it to make more informed decisions based on past interactions and events. This feature was particularly useful and powerful, but also with a big downside related to the LLM context length;
- **REDUCED_MAP:** This parameter was introduced to optimize the space the map occupied in the prompt by limiting the environment described as a slice of the full map and then scaling all the coordinates (of the agent and the parcels) treating the reduced map as the total map. The reduction in size was determined based on the maximum value between `PARCELS_OBSERVATION_DISTANCE` and `AGENTS_OBSERVATION_DISTANCE` (from Server Configuration File, see Section 3.2.1), ensuring that the agent only received the most relevant spatial data necessary for making informed decisions. Essentially, this optimization allowed the attention of the LLM to not be too sparse, but it brought some additional problems, for example by hiding any delivery zone from the map when the agent was placed too far away with the goal of delivering;
- **HELP_FIND_DELIVERY:** This parameter was specifically designed to assist the agent in locating delivery points more effectively. By setting it to `true`, the system ensured that the closest delivery point (using Manhattan Distance) was always included in the agent’s prompt (not as coordinates but as directions, eg. “right and up”), even if that particular delivery point was not within the agent’s immediate field of view. In fact, this parameter was implemented to remedy the problem described in the **REDUCED_MAP** point. This feature provided the agent with valuable directional guidance, allowing it to make better routing decisions and reducing the risk of wandering aimlessly in search of a delivery location (or worse, by looping again and again), but also reduced our ability to track the LLM ability in finding the delivery point by itself (more on this in Section 5.1);
- **HELP_SIMULATE_NEXT_ACTIONS:** The goal of this parameter was to enhance the agent’s decision-making process by simulating and displaying the expected outcomes of each possible action. When activated by setting it to `true`, the prompt provided to the LLM included a detailed breakdown of how each available action would alter the surrounding environment, by computing for each action the resulting map and attaching all of them to the final prompt. In theory, this additional information could help the agent anticipate and select the most favorable course of action. However, experimental results indicated that enabling this feature led to suboptimal performance, resulting in poor decision-making and inefficiencies, probably due to the size of the prompt;

This design resulted in an implementation that was bringing the project in the wrong direction, because the whole “no framework on top” idea was breaking down even if this didn’t have anything to do with planning in a strict sense. Without those helps, the agent was not able to perform well in any environment, and the decision was made to start over with a new approach.

Overall, while this initial approach did not yield the desired performance, it was an essential step for the next iterations of the agent’s development.

The code for this first implementation can be found in the `archive/raw_llm_agent.js` file inside the project repository [5].

4.2.2 Takeaways

Through this initial approach, we gained valuable insights into the challenges of designing an effective agent. The key lessons learned from this phase can be summarized as follows.

Why this approach has been discarded?

- Unclear prompt style: the way information was structured in the prompt was not optimal. This became particularly evident in the uncertainty computation, where the agent frequently exhibited high uncertainty in its actions;
- Over-reliance on helping parameter: providing excessive hints and structured input to the agent hindered its ability to explore the environment effectively. While guidance could be helpful, too much assistance made the results too distant from what the LLM could achieve by itself.

Key Takeaways from this phase

- Performance: when extensive guidance was provided, the results were still acceptable but inferior to those obtained using PDDL-based solutions. Initially, we considered implementing a PDDL version of the agent to serve as a benchmark. However, this comparison would have been unfair due to fundamental differences in approach and assumptions;
- Unintended biases in LLM behavior: Although not directly related to the core functionality of the agent, an interesting observation emerged regarding how the LLM interpreted the presence of other agents. In some tests, the server allowed the spawning of “enemy” agents capable of blocking paths. While this feature was not used in the main experiments, we discovered that simply including information about these agents in the prompt led the LLM to assume that agents near parcels were actively trying to steal them, making our agent change its direction. In reality, these agents were merely obstacles with no intent to compete for parcels. This behavior highlights inherent biases in the LLM’s training data, where similar context might have been associated with adversarial interactions.

4.3 Second Approach

In our exploration of different strategies for designing an LLM-driven agent, this second approach was the weakest one. The fundamental idea was to adopt a “full raw” approach, where the model received all available unprocessed data from the server without any pre-processing or filtering. The hypothesis was that by exposing the LLM to as much raw information as possible, it might be able to infer meaningful patterns and determine the best course of action on its own, giving us the ability to evaluate the LLM’s planning capabilities without any external structure.

Text

```
1 [ROLE]
2
3 Raw 'onMap' response: [JSON]
4
5 Raw 'onYou' response: [JSON]
6
7 Raw 'onParcelsSensing' response: [JSON]
8
9 [ACTIONS]
10
11 [QUESTION]
```

Listing 4.3: Scheme of a prompt used in the second approach, full in Appendix B.1

To achieve this, the agent’s prompt was constructed as a simple collection of “name of the server call: JSON result” for each function. Unlike other approaches that structured data into a more human-readable or semantically meaningful format, this method provided the LLM with a direct dump of server responses. The only additional elements in the prompt were the list of available actions and a query requesting the next step (necessary to compute the uncertainty).

At first glance, this approach seemed promising in that it completely avoided manual interpretation of server responses, reducing the need for custom logic or intermediate representations. If successful, it could have allowed for a highly adaptable agent that functioned independently of predefined schemas or rigid data structures.

However, in practice, this approach performed poorly. While it occasionally worked in very small maps, it became unreliable and inefficient as soon as the environment grew even slightly. The agent frequently took suboptimal paths, exhibited excessive backtracking, and often failed to reach its goal efficiently.

Additionally, for the agent to work correctly in such small maps, the parameter `USE_HISTORY` from the first implementation was set to `true`, allowing the LLM to leverage the ‘action-feedback’ history.

The lack of structured guidance made it difficult for the model to consistently produce useful responses, ultimately leading us to discard this approach in favor of more refined strategies.

4.3.1 Takeaways

Even if this approach was the weakest one and has been discarded very soon, it provided us with valuable insights into the limitations of the LLM in processing raw data. The key lessons learned from this phase can be summarized as follows.

Why this approach has been discarded?

- **Arbitrary Naming of API Calls:** The names of server calls were not standardized and could change depending on the server’s development. For instance, a call named `onMap` might return a list of map tiles, but there was no inherent guarantee of this behavior;
- **Lack of Context and Poorly Structured Input:** The raw JSON data lacked structured context, making it challenging for the LLM to infer the correct course of action. Additionally, the query format itself was suboptimal. For example, if the prompt simply asked, “Final goal: go to (x, y) , give me the next step to reach the goal”, the model might struggle to determine that the immediate goal was not reaching (x, y) but selecting the best intermediate step that moved the agent closer to that final destination.

Key Takeaways from this phase

- **LLM Inference Capabilities:** Surprisingly, the LLM was still able to extract some meaningful information from the unstructured data. While the results were inconsistent, there were instances where the model successfully inferred useful actions despite the messy prompt;
- **Significant Impact of Prompt Engineering:** Small modifications to the prompt led to drastic changes in the agent’s behavior. This highlighted the critical role of prompt engineering in optimizing LLM performance, a topic we will explore in detail in Section 5.2.

4.4 Final Agent

This represents the final iteration of our approach, incorporating substantial improvements in both prompt generation and the overall structure of the agent’s implementation.

Initially, prompts were dynamically constructed at runtime using a series of `if/else` conditions, which made them difficult to manage, debug, and scale. That approach lacked flexibility, as any modification required changes to the core logic of the agent, increasing complexity and reducing maintainability, while also adding more potential sources of error.

4.4.1 Prompt Management System

In the revised approach, we transitioned to a structured prompt management system where predefined templates are stored in a dedicated folder (`prompts/` folder in the GitHub repository [5]). These templates use variable placeholders that are replaced dynamically via regex-based substitutions. This change provided multiple advantages: it improved readability, ensured consistency across different prompts, and made modifications significantly easier. Instead of altering the logic of the agent itself,

changes could now be made directly at the prompt level, allowing for rapid experimentation and iteration. Additionally, this structured approach facilitated more systematic testing, as different versions of the prompts could be evaluated with minimal effort. Overall, this refinement not only enhanced the reliability of the agent but also contributed to a more efficient development workflow.

4.4.2 Agent Refactoring

Beyond improving prompt handling, a major structural refactoring was undertaken to optimize the agent’s implementation. The original codebase was relatively monolithic, with large, complex functions handling multiple aspects of decision-making. This made debugging and extending functionality cumbersome. In the final version, the implementation was restructured into a more modular design, with a greater number of smaller, well-defined functions handling specific tasks. This decomposition significantly reduced code redundancy, improved readability, and made the agent easier to modify. One of the key benefits of this modular approach was its impact on the data acquisition process. Since the agent’s core logic was now more flexible, adapting it for different data collection tasks required minimal effort. Simple function modifications or prompt adjustments were often sufficient to tailor the agent’s behavior to new requirements. This ability to rapidly reconfigure the agent streamlined experimentation and allowed us to gather diverse datasets efficiently, ultimately improving the quality and scope of our evaluations.

4.4.3 RAG Experiments

We also explored the potential of Retrieval-Augmented Generation as a way to enhance the agent’s decision-making process. The initial idea was to categorize the parcels *a posteriori* based on predefined classes and provide the LLM with priority information for each category (example in Listing 4.4). This would have allowed the model to make more informed decisions by leveraging structured context about the importance of different types of parcels. However, while this approach showed promise, it was ultimately considered too far from the core objectives of this thesis and was therefore not pursued further.

JavaScript Code

```
1 if (PARCEL_CATEGORIZATION) {  
2   for (let parcel of rawOnParcelsSensing) {  
3     const parcelIdNumber = parseInt(parcel.id.substring(1));  
4     parcel.food = parcelIdNumber % 2 === 0 ? "banana" : "pineapple";  
5   }  
6 }
```

Listing 4.4: Discarded implementation of an example of *a posteriori* parcel categorization

Another experimental RAG-based approach involved providing the agent with direct information about past actions, such as: “The last time you were in position(x , y), you attempted to move **up**, but the path was blocked”. While this could be a useful strategy for a “blind” agent (one without direct state awareness) its application in our case would have led to a problematic behavior. The agent would work by just performing random actions, gathering information about the blocked/suboptimal paths and then relying entirely on the RAG-generated feedback to navigate. This effectively bypassed the agent’s core decision-making process, turning it into a reactive system rather than a proactive one. Given our focus on autonomous decision-making, this approach was deemed unsuitable.

Nonetheless, these experiments highlighted the potential of RAG for different types of autonomous agents and may be worth exploring in future research.

4.4.4 Stateless and Stateful Agents

The final agent played a crucial role in generating the heatmaps that will be discussed in Section 5.3, primarily leveraging GPT-4o-mini. The task categorization into **pickup** and **delivery** was handled by computing the number of parcels currently in possession of the agent, rather than inferred dynamically from the environment state (keeping the same prompt for both goals) to ensure clarity and allowing

for more controlled testing conditions.

Both stateless and stateful configurations of the agent were developed and tested. The stateless version made decisions based solely on the current state, while the stateful version incorporated historical context to refine its choices over time thanks to the **action-result** feedback in the conversation history.

4.4.5 Uncertainty Handling

To handle uncertainty in decision-making, we experimented with three different approaches:

- **Raw probability selection:** The agent directly selected the action with the highest probability, without any additional modifications. This approach was straightforward, but led to suboptimal paths when the highest-probability action was incorrect;
- **Weighted selection:** Instead of always choosing the most probable action, the agent sampled actions based on their probabilities. This method was particularly effective in the *stateful* configuration: if an incorrect action initially had the highest probability, randomness allowed the agent to eventually correct itself over multiple iterations on the same spot;
- **Stopping mechanism:** This approach follows literally the process explained in the paper ‘Robots That Ask For Help: Uncertainty Alignment for Large Language Model Planners’ (explained in Section 2.2.2.3): by filtering the actions after the computation of the scaled log-probabilities, if the resulting set of possible actions was not a singleton the agent would stop waiting for the user input. This method was implemented for completeness in the development, but not tested automatically since it would have required human intervention.

Among these methods, weighted selection proved to be the most effective for data acquisition and testing, as it leveraged randomness to improve path accuracy. This experiment demonstrated that even when an LLM-based agent lacks perfect reasoning abilities, incorporating controlled stochasticity can help guide it toward better long-term performance.

4.4.6 Takeaways

We can summarize the biggest improvements in the final agent as “prompt management” and “agent refactoring”.

The change in the prompt management provided multiple advantages:

- Improved readability: The separation of logic from text made prompts easier to understand and modify;
- Consistency and maintainability: Storing prompts as templates reduced duplication and made debugging simpler;
- Faster experimentation: Changes could be made at the prompt level without modifying the agent’s core logic;
- Better testing and evaluation: Different prompt versions could be systematically compared with minimal effort.

The refactoring of the agent structure also brought several advantages:

- More modular design: Smaller, well-defined functions improved readability and maintainability;
- Reduced code redundancy: Reusable components simplified implementation and debugging;
- Greater flexibility: The agent could be easily adapted for new tasks, making data acquisition faster and more efficient.

4.5 Extra: Closest Cell to the Goal

As part of the iterative process of refining our agent, we also tried to systematically isolate and address specific challenges by progressively reducing the complexity of the final problem. This incremental approach not only helped us pinpoint potential weaknesses in the agent’s decision-making process but also provided deeper insights into the LLM’s underlying behaviors and limitations.

To achieve this, we conducted a series of controlled experiments, including:

- **Testing on smaller, simplified maps:** By reducing environmental complexity, we could more easily observe the agent’s decision-making patterns and identify whether failures were due to reasoning errors or external factors;
- **Using custom, “disposable” prompts:** We introduced minor variations in prompt structures to assess how sensitive the LLM was to different formulations of the problem. This helped us determine whether misinterpretations were caused by the model itself or by the way the information was presented;
- **Decomposing the final goal into smaller, manageable sub-goals:** Breaking down the problem into intermediate objectives allowed us to test whether the agent could handle incremental progress rather than needing to solve the entire task at once. This approach is further detailed in the following paragraphs.

One key issue we wanted to investigate was why the agent often failed to select the optimal action leading toward the goal. Was the model inherently incapable of making the correct choice, or was the issue rooted in the way the prompt was structured? To better understand this, we designed an experiment that introduced a two-step decision-making process (visible in Figure 4.1):

1. **Identifying the best neighboring tile:** Before making a move, the agent was first asked to determine the most optimal adjacent tile to step toward, effectively transforming the problem into a local optimization task;
2. **Selecting the best action to reach that tile:** Once the target tile was identified, the agent was then tasked with choosing the appropriate action to move in that direction.

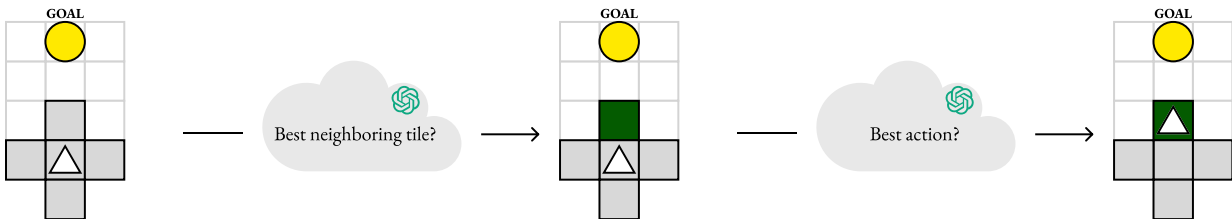


Figure 4.1: Two-steps decision making example

This method effectively reframed the problem from a complex, global pathfinding challenge into a series of simpler, localized decisions. We hypothesized that this decomposition would improve performance by reducing the cognitive load on the LLM and allowing it to focus on smaller, well-defined tasks.

However, despite these modifications, the results did not align with our expectations. The primary issue was that if the agent misidentified the best neighboring tile by choosing a position that actually increased the distance to the goal, then even a perfectly chosen direction would still lead to suboptimal behavior.

Another significant factor was the role of uncertainty computation in the LLM response. When relying on the raw output of the LLM (‘Raw probability selection’ in Section 4.4.5), results were often inconsistent or outright incorrect. To address this, we experimented also with the uncertainty-weighted random action selection approach (‘Weighted selection’ in Section 4.4.5), where decisions

were influenced by the model’s confidence levels. Unfortunately, this introduced additional challenges, as interpreting the results became even more complex.

Ultimately, this experiment underscored some limitations of LLM-based decision-making in logistics scenarios where local optimizations must align with global objectives. The agent’s inability to consistently select the correct neighboring tile demonstrated how small errors at the decision-making level could compound, leading to inefficient or even counterproductive actions.

However, as we will explore in Chapter 6, newer LLM models exhibit notable improvements in our task, suggesting that continued advancements in the technology may help overcome these challenges and enhance overall performance in complex logistics tasks.

5 Data Collection

In this chapter, we provide a comprehensive account of the data collection process that underpins our experiments. We detail the key steps involved in designing, implementing, and refining our methodology, with a particular focus on the construction and optimization of the prompts used to interact with the Large Language Model.

We begin by presenting the core prompts utilized in our study, describing their specific applications within each experimental setup. Following this, we outline the rationale behind our prompt design decisions, grounding our choices in relevant literature and highlighting key findings from previous research that informed our approach.

Finally, we introduce the structure of our experimental evaluations, including the generation of specialized heatmaps designed to illustrate the agent’s uncertainty in action selection. These visualizations provide valuable insights into the model’s decision-making process, highlighting both its generative strengths and areas for further refinement.

5.1 Prompts

This section provides a general overview of the prompts used in our experiments, while the following one (Section 5.2) will detail the specific choices behind each specific part of them.

Our prompts were designed to be both concise and comprehensive, ensuring they effectively guided the model in generating appropriate next-step actions. The core idea was to elicit responses that incrementally moved the agent toward its objective, either picking up or delivering a parcel, while minimizing ambiguity in the instructions.

The two primary prompts used in our study are:

- Pickup prompt: guiding the agent to retrieve a parcel from a specified location;
- Delivery prompt: prompting the agent to transport the parcel to its designated drop-off point.

Each prompt was carefully structured to provide relevant contextual information while avoiding unnecessary complexity that could dilute the model’s focus.

5.1.1 Pickup prompt

The Pickup prompt (Listing 5.1) was used to evaluate the model’s ability to generate actions in the *pickup* scenario. In this setup, the agent was tasked with picking up a parcel from a specific location on the map. The prompt provided the agent with the map information, including the map dimensions, the location of the parcel, and the agent’s current position.

The model was then asked to determine the optimal next action that would bring the agent one step closer to the parcel.

Text

```
1 You are a delivery agent in a web-based delivery game where the map is a matrix
2 I am going to give you the raw information I receive from the server and the
   possible actions.
3 Map width: {width}
4 Map height: {height}
5 Tiles are arranged as {height} rows in {width} columns:
6 {tiles}
7 The parcel you need to take is in the spot ({parcelX}, {parcelY}).
8 You are on the spot ({agentX}, {agentY}).
9 The actions you can do ONLY if the next tile is available are:
10 U) move up
11 D) move down
12 L) move left
13 R) move right
14 T) take the parcel that is in your tile
15 S) ship a parcel (you must be in a delivery tile)
16
17 Your final goal is to go to a tile with the parcel and (T)ake it, I need the
   best action that will get you there.
18 Don't explain the reasoning and don't add any comment, just provide the action's
   letter.
19 What is your next action?
```

Listing 5.1: Pickup prompt used in the experiments

5.1.2 Deliver prompt

The Delivery Prompt (Listing 5.2) was structured to evaluate the model's ability to navigate toward a designated drop-off location. Unlike the Pickup Prompt, where the parcel's position was explicitly provided, the delivery prompt required the LLM to infer the delivery destination from the map description. The model was then asked to determine the best next action to move toward the inferred delivery location.

Text

```
1 You are a delivery agent in a web-based delivery game where the map is a matrix
2 I am going to give you the raw information I receive from the server and the
   possible actions.
3 Map width: {width}
4 Map height: {height}
5 Tiles are arranged as {height} rows in {width} columns:
6 {tiles}
7 You are on the spot ({agentX}, {agentY}).
8 The actions you can do are:
9 U) move up
10 D) move down
11 L) move left
12 R) move right
13 T) take the parcel that is in your tile
14 S) ship a parcel (you must be in a delivery tile)
15
16 You have a parcel to ship, your final goal is to go to the delivery zone (
   delivery = true) and (S)hip the parcel, I need the best action that will get
   you there.
17 Don't explain the reasoning and don't add any comment, just provide the action's
   letter.
18 What is your next action?
```

Listing 5.2: Deliver prompt used in the experiments

5.2 Prompt Creation Choices

This section outlines the key considerations that guided the construction of our prompts (commonly referred to as **Prompt Engineering**). We describe the reasoning behind specific design choices, following the sequence in which they appear in the prompt. For clarity, we use the Pickup Prompt (Listing 5.1) as a reference.

5.2.1 Role Prompting

Assigning specific roles or personas to Large Language Models within prompts, known as “role prompting,” has been shown to enhance their performance on various tasks. This technique involves instructing the model to assume a particular identity, such as a “math professor” or “geographer,” to guide its responses more effectively. The concept of role prompting has been explored in several studies.

For instance, the paper ‘Better Zero-Shot Reasoning with Role-Play Prompting’ by Kong et al. [15] demonstrated that strategically designed role-play prompts can significantly improve LLMs’ reasoning abilities across diverse benchmarks.

Similarly, ‘Role-Play Zero-Shot Prompting with Large Language Models for Open-Domain Human-Machine Conversation’ by Njifenjou et al. [21] investigated the use of role-play prompts to enhance conversational agents’ performance without additional fine-tuning.

In our experiments, we employed role prompting to encourage the model to adopt the persona of a *delivery agent*, thereby focusing its attention on the task at hand. By framing the prompts in this manner, we aimed to guide the model’s responses towards generating coherent and contextually relevant actions.

5.2.2 Map Encoding to Reduce Attention Sparsity

One of our primary concerns was that the map included in the prompt might take up too much space, which could lead to an excessively sparse distribution of attention. This could result in the model not properly focusing on crucial aspects of the input. At the same time, we wanted to maintain a minimal and flexible approach to map handling. Our goal was to ensure that the system would continue functioning even if the format of the map-providing function were to change unexpectedly. As discussed in Section 3.1.1, this design decision was made to simulate an unpredictable and rapidly changing environment.

To address this concern, we first aimed to understand how the model processed map-related information. Ideally, we would have examined the models’ attention scores directly, but unfortunately, OpenAI’s API does not provide access to these scores. As a result, we had to rely on a qualitative analysis instead. To achieve this, we leveraged BERT, another Transformer-based architecture, to visualize attention patterns in our original prompt. The results revealed that a significant portion of the models’s attention was focused on tokens related to JSON syntax rather than the meaningful content of the prompt itself. This suggested that structural elements, rather than semantic information, were drawing a disproportionate amount of attention.

Based on this observation, we attempted to reduce the presence of JSON syntax within the prompt while preserving the essential information. We then repeated the same attention analysis using this modified version. The results showed a shift in the distribution of attention, even if the two prompts contained virtually the same information. Specifically, when we examined the tokens receiving the highest attention — excluding punctuation, as well as the special starting and ending tokens — we found that in the original prompt, only 23 out of the top 50 tokens were meaningful words. In contrast, in the revised prompt, this number increased to 31 out of 50, indicating a more focused attention on relevant content.

For this test, we used an older version of the “Pickup prompt”, which is explained in detail in Section 4.3, with a small 2×5 map. As demonstrated in Table 5.1, the revised version of the prompt not only contained more words in the top 50 tokens, but also explicitly included the term **pickup**, which was the goal of the task. Notably, this term was entirely absent from the top attention-receiving tokens in the initial prompt.

To further analyze the impact of this change, we collected attention scores from both the old and new prompt versions. We then removed the tokens corresponding to the map representation in both prompts and plotted the difference between the attention in the remaining common 264 tokens. The

Old Prompt	New Prompt
game	game
reasoning	parcel
parcel	reasoning
parcel	i
score	actions
response	you
parcel	if
tile	parcel
i	moves
if	height
you	width
response	score
loop	ship
comment	tile
information	loop
actions	information
and	choosing
ship	you
and	pickup
server	and
you	you
choosing	parcel
using	and
	server
	information
	and
	ship
	explain
	using
	reward
	name

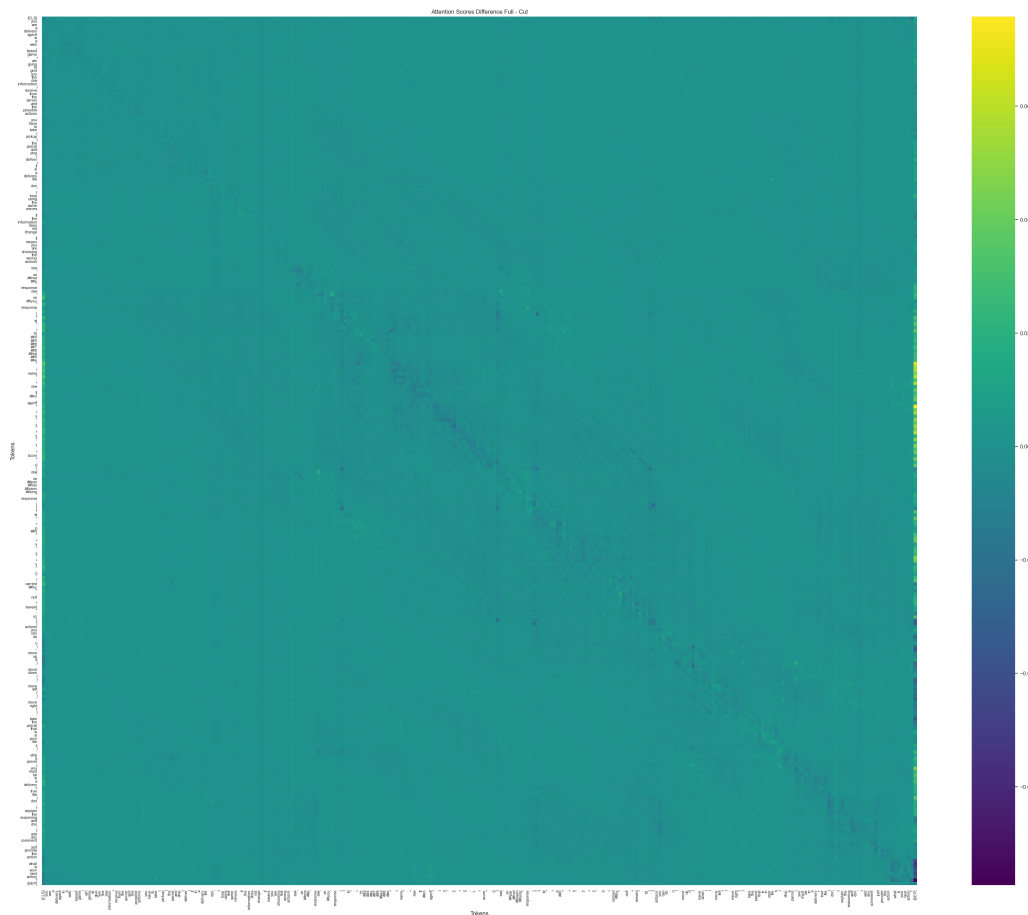
Table 5.1: Comparison of Old and New Prompts’ words appearing in the top 50 attention-receiving tokens

resulting visualization is presented in Figure 5.1 as the difference between *old prompt attention* and *new prompt attention*. Interpreting this figure is not entirely straightforward, and we acknowledge that it is unclear to what extent this analysis can serve as a direct proxy for a GPT-based model. However, one noticeable trend is that the diagonal elements exhibit a lower value (indicated by the presence of purple and dark blue regions), indicating that the *new prompt* received more attention than the *old prompt*. Meanwhile, most of the remaining areas show a delta close to zero (represented in aqua blue), suggesting that the changes we introduced led to a redistribution of attention without introducing excessive noise or unintended biases. While this number does not have a direct interpretative meaning, for completeness, we note that the total sum of the difference matrix is about -15.34476 , indicating more attention in the new version of the prompt.

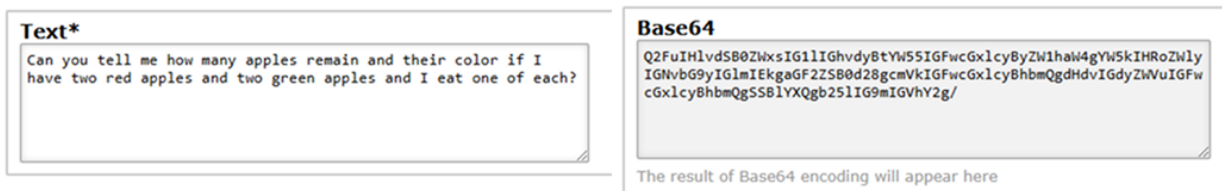
5.2.2.1 Emerging behavior: encoded question decoded answer

Another strategy we explored to reduce the space occupied by the map in the prompt was to encode it. Our goal was to compress the map representation while ensuring that the model could still understand the relevant information.

A relevant study on this topic is presented in the paper ‘LLMs Can Understand Encrypted Prompt: Towards Privacy-Computing Friendly Transformers’ by Xuanqi Liu and Zhuotao Liu[19]. The authors of this paper aimed to preserve user privacy by encrypting prompts before sending them to a language



model. Their findings demonstrated that the model was able to comprehend and respond appropriately to encrypted prompts. This suggested that LLMs possess an inherent ability to process encoded text meaningfully.



To validate this, we first examined the model’s ability to process encoded prompts by testing a simple case. As shown in Figure 5.3, the model successfully understood the meaning of the question “Can you tell me how many apples remain and their color if I have two red apples and two green apples and I eat one of each?” even when the entire prompt was encoded in BASE64 (see Figure 5.2). Notably, the model did not generate a decoding function before answering. Instead, it directly

processed the encoded prompt and returned a response in plain English, demonstrating its ability to work with BASE64-encoded text.

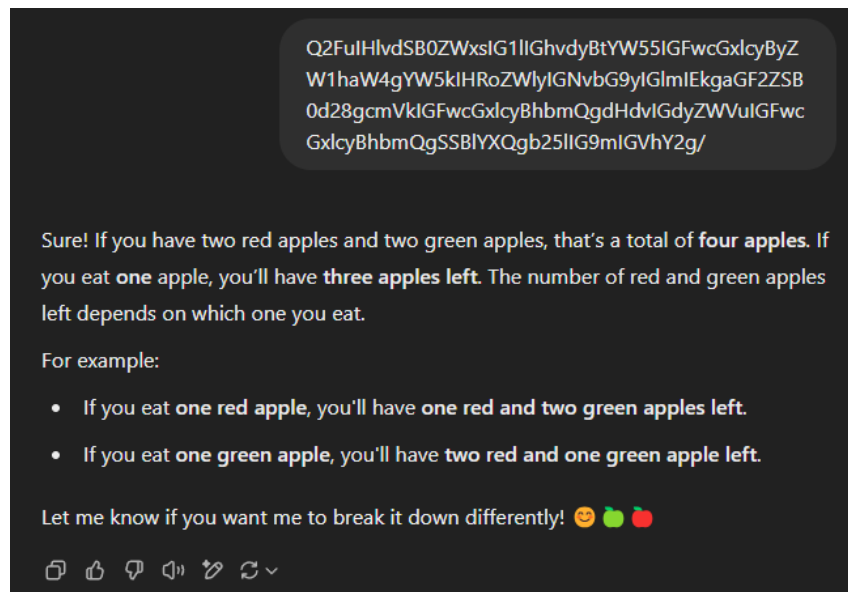


Figure 5.3: GPT-4o-mini response to a BASE64-encoded question
Source: GPT-4o-mini via OpenAI web interface

Encouraged by these results, we attempted to apply a similar encoding strategy to the map representation while also aiming to reduce the number of characters in the prompt. Our hypothesis was that reducing the character count would also reduce the number of tokens, thereby improving attention efficiency.

To test this, we implemented a two-step encoding process in Python, as outlined in Listing 5.3:

1. First, we compressed the map using the `zlib` library and its Deflate algorithm ² to minimize its size;
2. Next, we encoded the compressed output in BASE64 ³ using the `base64` library;
3. Finally, we inserted this doubly encoded map into the prompt.

Python Code

```
1 [...]
2
3 input_string = MAP
4 compressed_data = zlib.compress(input_string.encode('utf-8'))
5 compressed_base64 = base64.b64encode(compressed_data).decode('utf-8')
6
7 [...]
```

Listing 5.3: Example of double encoding algorithm

While this method did succeed in reducing the number of characters by $\sim 65\%$, the results were not as expected in terms of model comprehension. Unlike the case of simple BASE64 encoding, the LLM was no longer able to interpret the prompt correctly. Instead, its responses typically fell into one of two categories:

²<https://en.wikipedia.org/wiki/Deflate>

³<https://en.wikipedia.org/wiki/Base64>

- The model would explicitly state that it recognized the input as an encrypted message and generate something similar to “It seems you’ve provided a compressed string or encoded data. Could you clarify how you’d like to use or process this? If it’s encoded text, I can try decoding it for you.”;
- Alternatively, if instructed with the encryption method used on the text, the model would generate a Python function to decode the input before attempting to process it. Such a function would have been executed “under the hood” to decrypt the message before processing.

Furthermore, although our approach successfully reduced the total number of characters in the prompt, it did not achieve our ultimate objective.

We then tried only using the BASE64 encoding, but the number of character actually increased rather than decreased as well as the number of tokens that increased even more in percentage. This is because tokenization in LLMs is based on common character sequences rather than individual characters. Encoding the map disrupted these common patterns, leading to a tokenization process that resulted in a higher overall token count. Consequently, the intended effect of reducing attention sparsity was not achieved, as illustrated in Figure 5.4.

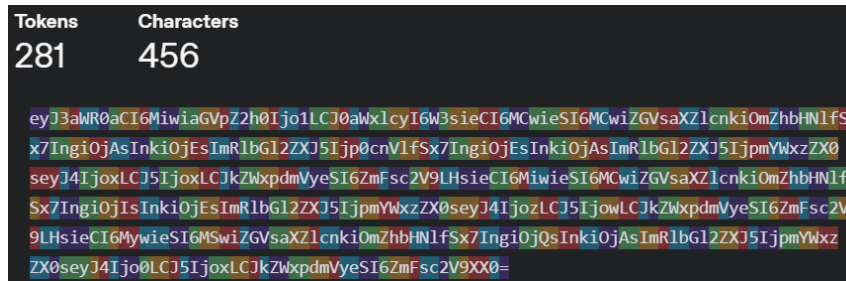


Figure 5.4: Tokenization of a BASE64-encoded text
Source: *GPT-4o tokenizer via OpenAI web interface*

In summary, while encoding techniques like BASE64 alone might be useful for preserving meaning in certain contexts, our specific approach of compressing and encoding the map did not lead to the desired efficiency gains. Instead, it introduced additional computational overhead for the model, ultimately making this approach ineffective for our purposes.

5.2.3 Emerging behavior: math capabilities

We also investigated the impact of including action-related information, such as explicitly stating “Going up increases your X position by 1” for every movement action. It initially seemed like a promising way to guide the model’s reasoning. However, this conflicted with our overarching goal of keeping the system adaptable rather than grounding it in a fixed environmental structure. Despite this, we proceeded with the experiment to assess its impact.

Interestingly, we observed that the model’s performance in navigation tasks was unexpectedly strong, even to the point of raising concerns about potential shortcuts. To further investigate, we conducted a follow-up test in which we entirely removed the map from the prompt. As result, the model still produced correct answers by seemingly ignoring any extraneous details and reducing the problem to simple arithmetic. It recognized that if the starting position was (4,2) and the goal was at (0,0), the necessary movement was simply decreasing X by 4 and Y by 2. This behavior demonstrated that the model could abstract away the spatial representation and “cheating” by operating purely on mathematical reasoning.

This finding aligns with existing research on the emergent mathematical reasoning capabilities of large language models. For example, Wei et al. highlight how LLMs exhibit “emergent abilities”, capabilities that do not appear in smaller models but spontaneously manifest as model scale increases [32]. Among these abilities, arithmetic and mathematical reasoning are particularly notable, suggesting that LLMs can generalize numerical patterns without explicit training for such tasks. Our experiment provides anecdotal evidence supporting this: rather than relying on the environmental constraints

provided in the prompt, the model instinctively leveraged its inherent mathematical capabilities to deduce the necessary movements. Ultimately, given this behavior, we abandoned the idea of including the results of the actions in the prompt.

5.2.4 Question Structure

A summarized example of a prompt used in the paper *Robots That Ask For Help: Uncertainty Alignment for Large Language Model Planners* [28] can be seen in Listing 5.4. This prompt has been taken directly from the example available in their interactive demo⁴, which is accessible online. The prompt follows a structured approach that allows the language model to process information about the environment, understand the task, and select the best action from a predefined set of choices.

Since we applied the same uncertainty analysis to the model’s responses, we opted to maintain a similar structure for our prompts.

Text

```

1 You are a robot operating in an office kitchen. You are in front of a counter
   with two closed drawers, a top one and a bottom one. There is also a
   landfill bin, a recycling bin, and a compost bin.
2
3 On the counter, there is an energy bar, a banana, and a microwave.
4 Put the snack next to the microwave.
5
6 A) pick up the energy bar and put it next to the microwave
7 B) pick up the banana and put it next to the energy bar
8 C) pick up the banana and put it next to the microwave
9 D) pick up the energy bar and put it next to the banana
10
11 Which option is correct? Answer with a single letter.

```

Listing 5.4: Prompt from ‘Robots that ask for Help’ paper

5.2.4.1 Structure of the Paper’s Prompt

The structure of this prompt follows a well-defined pattern that guides the LLM’s reasoning process. It consists of the following key elements:

- **Role and Environment Description:** The prompt starts by establishing the role of the model (a robot) and providing an overview of its working environment (an office kitchen);
- **Task Specification:** The next section provides a direct and unambiguous description of the task to be performed—in this case, placing a snack next to the microwave;
- **Action Choices:** A predefined list of actions (A, B, C, D, E) is presented, each corresponding to a possible decision. This format constrains the response space to just five choices;
- **Question and Response Constraint:** The prompt ends with a clear question that explicitly instructs the model to select an answer in a specific format (a single letter). This restriction allows for more structured uncertainty analysis as explained in Section 2.2.2.3.

5.2.4.2 Comparison with Our Approach

Similarly, our prompt design follows the same structured approach: given our task, our prompt structure aligns closely with the example above but incorporates additional environmental details specific to our use case.

Our prompt structure consists of the following elements:

⁴<https://tmp.com>

- **Role and General Environment Description:** The prompt begins by defining the model’s role (a delivery agent) and its operating context (a web-based delivery game);
- **Detailed Environment Specification:** Unlike the office kitchen scenario, our logistics task involves a structured map-based environment. Therefore, we explicitly include details such as:
 - Map dimensions;
 - Tile types and obstacles;
 - Parcel location;
 - Agent’s current position.
- **List of Possible Actions:** Instead of using letter-based choices (A, B, C, D, E), our prompt presents a set of movement and interaction commands:
 - **U, D, L, R:** Move Up, Down, Left, Right;
 - **T/S:** Pick up or deliver a parcel.
- **Goal Definition and Question:** The prompt explicitly states the agent’s goal, such as reaching a specific destination or delivering a parcel to an inferred goal tile. Additionally, the question is phrased to ensure a concise and structured response:

“Your final goal is to [...] take the parcel. Just provide the action’s letter. What is your next action?”

This ensures that the model’s response remains within the expected format, to evaluate the uncertainty in the same way as in the original paper.

A brief note: The letter “T” for “Take” was chosen instead of “P” for “Pickup” because, in the initial version of the prompt (as discussed in Section 4.2), *P* was used to represent the parcel on the map. Similarly, “S” for “Ship” was selected, even though the server and other components refer to it as “Deliver,” as the letter *D* was already assigned to the Down movement action.

5.2.4.3 Multichoice Benchmarking

Constraining the model’s response to a predefined set of choices simplifies the evaluation process. Moreover, multi-choice question answering is a widely adopted method for benchmarking language models, as demonstrated in datasets such as the Massive Multitask Language Understanding (MMLU) benchmark⁵. By enforcing a structured question format, we can systematically assess the model’s performance across different experimental conditions.

5.2.5 Goal positioning

As highlighted in the article ‘The Needle In a Haystack Test: Evaluating the Performance of LLM RAG Systems’[1], LLMs exhibit a tendency to prioritize information positioned at the beginning or end of a prompt, often overlooking details embedded in the middle. To analyze this phenomenon, the researchers conducted an experiment by inserting a unique “needle” token at different positions within a prompt and measuring whether the model could successfully retrieve it in its response. Their goal was to determine the optimal placement of information retrieved from a retrieval-augmented generation system to maximize recall. Their findings revealed that LLMs are most likely to recall information from the beginning and, to a lesser extent, from the end, while details placed in the middle are more frequently neglected.

This positional bias has significant implications for prompt engineering, especially in structured queries. In our thesis, we follow a similar principle by positioning the specific request at the end of the prompt. Understanding the model’s attention distribution allows us to optimize prompt design, ensuring that key details receive the necessary emphasis to improve response accuracy and relevance.

Moreover, by placing the primary question at the end of the prompt, as well as all the information that changes from a request to another, we can leverage OpenAI’s prompt caching capability of the API.

⁵<https://en.wikipedia.org/wiki/MMLU>

5.2.5.1 Leveraging Prompt Caching

As can be seen in Listing 5.2 and Listing 5.1, the “changing part” of the prompt (mainly the position of the agent, but in case of a filtering of the actions, also any updated list of possible actions), was placed at the end, while the main static components, including the role and map, were positioned at the top. This structure was the result of a literature analysis while it also takes advantage of OpenAI’s prompt caching API⁶.

According to OpenAI’s documentation, for any prompt exceeding 1000 tokens, only the modified portion is recomputed, while the cached static portion remains unchanged. This significantly improves efficiency by reducing computational overhead, resulting in faster response times and lower operational costs.

The benefits of this approach are especially pronounced in the stateful version of the agent, where the LLM receives the entire conversation history with each request. In this case, caching allows the model to handle long, continuous interactions without constantly reprocessing the entire history from scratch, making the system much more scalable and responsive.

By structuring prompts in this manner, we ensure that also queries for a stateless agent in a bigger environment or with more actions available (necessary to reach the 1000 tokens count) can benefit from caching, leading to optimized performance in a scenario that requires frequent API calls.

5.3 Uncertainty Visualization

A **heat map** (or heatmap) is a 2-dimensional data visualization technique that represents the magnitude of individual values within a dataset as a color. The variation in color may be by hue or intensity.

[...]

There are two main type of heat maps: spatial, and grid.

- A spatial heat map displays the magnitude of a spatial phenomena as color, usually cast over a map. In the image labeled “Spatial Heat Map Example,” temperature is displayed by color range across a map of the world. Color ranges from blue (cold) to red (hot).
- A grid heat map displays magnitude as color in a two-dimensional matrix, with each dimension representing a category of trait and the color representing the magnitude of some measurement on the combined traits from each of the two categories. For example, one dimension might represent year, and the other dimension might represent month, and the value measured might be temperature.

Source: Wikipedia⁷

As stated in the definition above, heatmaps are a powerful tool for visualizing data, allowing for a compact and intuitive representation of large datasets. In our case, we employ heatmaps to capture and convey the uncertainty in the agent’s decision-making process at each position within the matrix. Specifically, the heatmaps encode the probabilities assigned by the KnowNo framework to the various possible actions.

The primary goal of these heatmaps is to provide a clear and structured means of analyzing how the agent distributes probabilities over its available actions. By mapping these probabilities onto a visual representation, we gain insights into which areas of the matrix exhibit high certainty (where one action dominates) and which regions display greater uncertainty (where multiple actions hold comparable probabilities). These heatmaps serve as a diagnostic tool for evaluating the agent’s behavior, identifying patterns, and pinpointing potential areas for improvement in future iterations.

More specifically, we utilize this visualization for two key aspects of our analysis:

- **Heatmaps:** To represent the probability distribution over the filtered list of actions. This means that after discarding actions with probabilities below the predefined threshold (as detailed

⁶<https://platform.openai.com/docs/guides/prompt-caching>

⁷https://en.wikipedia.org/wiki/Heat_map

in Section 2.2.2.3), the heatmap displays the remaining probability assigned to the remaining actions (with the total percentage scaled back to 100%);

- **Correctness Heatmaps:** To represent the overall probability assigned to the correct actions in each cell of the matrix. This allows us to measure how well the model aligns with the expected optimal behavior, providing a way to assess the agent’s effectiveness in selecting appropriate actions.

These heatmaps are central to our data collection and analysis pipeline. By systematically constructing and analyzing them, we can track how the agent’s decision-making evolves over time and across different configurations of the problem space. We can also identify systematic behaviors, biases, or areas where the model struggles, that we will analyze in Chapter 6.

The final analysis will rely heavily on these heatmaps to provide a comprehensive understanding of the agent’s performance and behavior.

5.3.1 Heatmaps

Heatmaps provide a visual representation of the probability distribution of actions taken within each cell of the environment. By encoding probability values as color intensities, these heatmaps offer an intuitive way to interpret the model’s decision-making process across the entire grid. Each cell in the heatmap corresponds to a specific position in the environment, with colors indicating the not-discarded actions in that location.

To generate these heatmaps, a specialized stateless agent systematically scans the map, cell by cell, from the top-left to the bottom-right. At each position, it records the probabilities assigned to each possible action. The collected data is then stored in a JSON format, as shown below:

```
[
  {
    "x":0,"y":0,
    "values":[
      ["R",true,0.9869068698680659],
      ["D",false,0.004691684231979342],
      ["U",false,0.002214120084731493],
      ["S",false,0.0021401869314925225],
      ["L",false,0.002023569441865407],
      ["T",false,0.002023569441865407]
    ]
  },
  ...
]
```

In this format, each entry represents a specific cell identified by its coordinates (x, y) . The **values** field contains a list of possible actions, each represented by:

- The action itself (e.g., “R” for Right, “D” for Down, etc.);
- A boolean value indicating whether the action was retained after filtering;
- The probability assigned to that action by the framework.

Once this data is collected, a Python script processes it to generate the final heatmap, which visually encodes the probability distribution for each action.

An example of such a heatmap is shown in Figure 5.5. Only the probability of the retained actions is displayed.

To enhance interpretability, each action is represented using a distinct color:

- Green: Left (L)

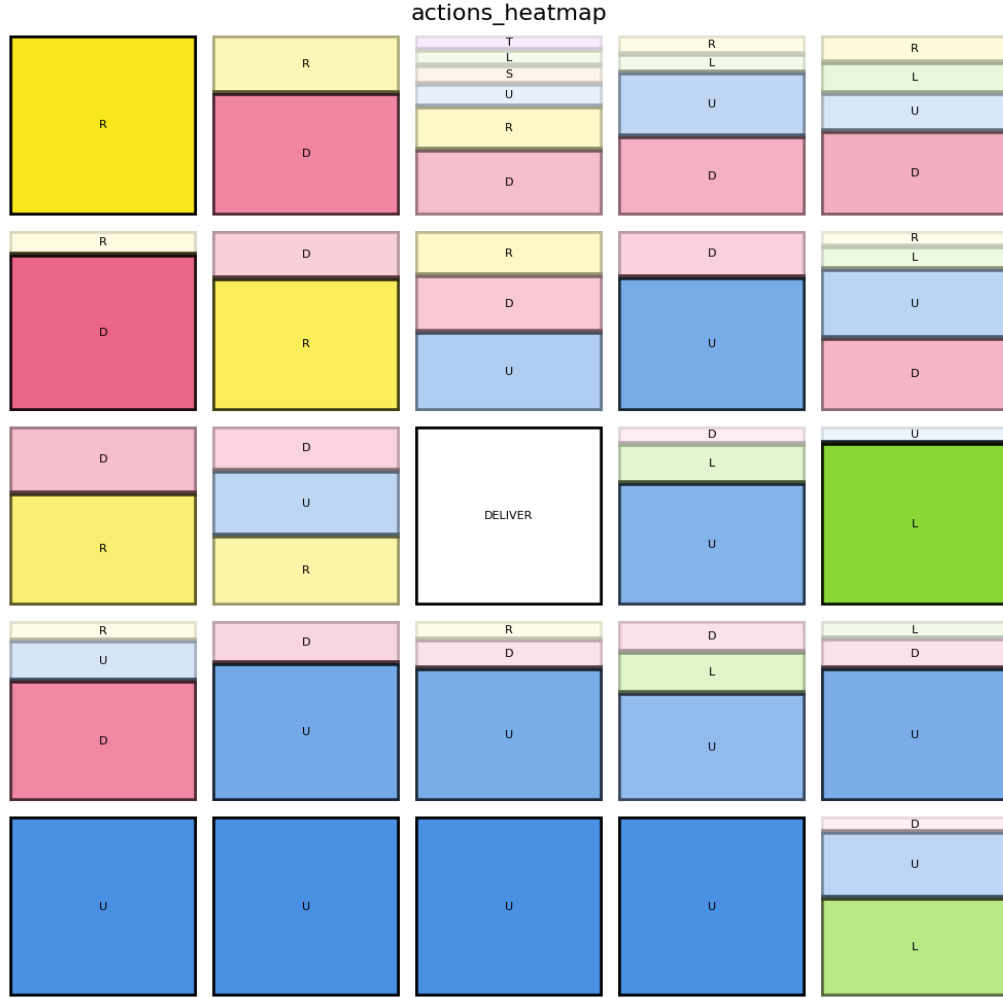


Figure 5.5: Example of a heatmap generated for a 5×5 map

- Yellow: Right (R)
- Red: Down (D)
- Blue: Up (U)
- Purple: Take (T) (picking up a parcel)
- Orange: Ship (S) (delivering a parcel)

Additionally, the transparency (alpha value) of each color is adjusted based on the probability of that action, allowing more probable actions to appear more prominently.

A total of 36 test cases have been conducted, and the resulting heatmaps are stored in the repository under the following directory structure:

`../data_and_results/heatmap/heatmaps/MODEL/MAP/GOAL/GOAL_POSITION/`

where the last folder contains four files:

- `actions_heatmap.png`: the rendered heatmap referring to that specific test;
- `correctness_heatmap.png`: the heatmap representing the probability of the correct actions that will be discussed in Section 5.3.2;
- `heatmap.json`: the raw data collected for the heatmap;

- `topX_values.json`: statistics regarding the `correctness_heatmap.png` that will be used in the final analysis in Chapter 6.

This organizational structure guarantees that, for each map size, both pickup and delivery goals are thoroughly tested to cover a wide range of scenarios. Additionally, in certain cases, we have included other specific goals or prompts to examine particular behaviors or edge cases of the agent. By incorporating these varied test configurations, we ensure a comprehensive evaluation of the agent's performance, which will provide valuable insights for analyzing the results in more depth.

The location of the goals are varied across different test cases; the rationale behind choosing “which tile should be the goal” will be discussed in Chapter 6.

5.3.1.1 Example

To better understand how to interpret these heatmaps, let us analyze a specific case. Figure 5.5 presents a heatmap generated for a 5×5 grid, where the agent's objective is to deliver a parcel at the center of the map.

Consider the cell at coordinates (2,1), which is positioned directly to the left of the goal. This location is particularly interesting because, being adjacent to the delivery point, the agent should ideally display a strong preference for moving Right (R) to complete the task efficiently. The KnowNo framework produced the following probability distribution for this cell:

```
...
{
  "x": 2,
  "y": 1,
  "values": [
    ["R", true, 0.34148147866149986],
    ["U", true, 0.3185527747471083],
    ["D", true, 0.2173149034701972],
    ["T", false, 0.04243503226711094],
    ["L", false, 0.040951236117811166],
    ["S", false, 0.03926457473627248]
  ]
},
...
```

From this data, we can extract several key observations:

- **Dominant Actions:** The three highest-probability actions at this cell are Right (R) with 34.14%, Up (U) with 31.85%, and Down (D) with 21.73%. This indicates that while the model slightly favors moving Right, it still considers moving Up or Down as viable options;
- **Discarded Actions:** Actions such as Left (L) (4.10%), Take (T) (4.24%), and Ship (S) (3.93%) were assigned very low probabilities and subsequently filtered out by the KnowNo framework, meaning they were not considered in the final action selection process. This suggests that the model appropriately recognizes that moving left or attempting to interact with the parcel at this position is not optimal;
- **Probability Normalization:** After discarding the low-probability actions, the remaining three probabilities were rescaled to sum to 100%. This renormalization ensures that the final action selection is based solely on the most relevant choices.

This example illustrates how the heatmap helps us diagnose the model's decision-making tendencies. Ideally, in this scenario, the probability of moving Right (R) should be significantly higher, given that it is the only direct path to the goal. However, the model also considers alternative movements, highlighting potential limitations in its capabilities.

This example already highlights some potential limitations of the system, particularly in terms of decision ambiguity when the agent is very near to the goal, which will be further examined in Chapter 6.

5.3.2 Correctness Heatmaps

After constructing the action-probability heatmaps, we generate an additional set of heatmaps that focus on correctness. These correctness heatmaps aim to evaluate how well the model aligns with the expected optimal behavior by quantifying the probability assigned to the correct actions at each position in the map.

To construct these correctness heatmaps, we start determining the set of correct actions for each cell using a Python script similar to the one in Listing 5.5. This script systematically analyzes the spatial relationship between each tile and the designated goal position, identifying which movements (e.g., U, D, L, R) would bring the agent closer to the goal. The resulting list of correct actions serves as the reference against which we measure the model’s decision-making accuracy.

Python Code

```
1 for tile in tiles:
2     delta_y = goal_y - tile['y']
3     delta_x = goal_x - tile['x']
4     correct_moves = []
5     if delta_x > 0:
6         correct_moves.append('D')
7     if delta_x < 0:
8         correct_moves.append('U')
9     if delta_y > 0:
10        correct_moves.append('R')
11    if delta_y < 0:
12        correct_moves.append('L')
13
14    percentage = 0
15    for value in tile['values']:
16        if value[0] in correct_moves:    # value[0] is the letter
17            percentage += value[2]       # value[2] is the probability
18
19    total_percentage = sum([value[2] for value in tile['values']])
20    percentage = percentage / total_percentage
```

Listing 5.5: If statements to compute the correct actions for every cell

The correctness heatmap visually represents the proportion of probability assigned to correct actions in each cell. This allows us to assess the model’s effectiveness in adhering to expected behavior and helps highlight areas where the agent struggles to make optimal choices.

Figure 5.6 illustrates an example correctness heatmap corresponding to the action probability heatmap shown earlier in Figure 5.5.

5.3.2.1 Example

To better understand the significance of the correctness heatmap, let’s revisit the previously analyzed cell at coordinates (2,1), which is located directly to the left of the goal position in a 5×5 grid.

To compute the correctness probability for this cell, we get the probability of the correct action (Right) and normalize it against the total probability assigned to retained actions.

Text

```
1 total_percentage = 0.34148 + 0.31855 + 0.21731 (= 0.87734)
2 percentage = 0.34148 / total_percentage (= 0.38922)
```

Listing 5.6: Scaling of the probability after filtering



Figure 5.6: Example of a correctness heatmap generated for a 5×5 map

Thus, if the agent selects an action randomly while weighting choices according to their assigned probabilities (as described in Section 4.4.5), in this case it would have a 38.92% chance of choosing the correct action, as shown in Listing 5.6.

This metric provides valuable insight into the model’s decision-making tendencies: while the model recognizes the correct action, it also considers alternative actions with relatively high probability, indicating uncertainty in its decision-making process.

By examining correctness heatmaps across different test cases, we can identify patterns and inconsistencies in the model’s behavior. These heatmaps help us pinpoint areas where the model exhibits high confidence in correct actions and areas where it is more uncertain or prone to errors.

Some key insights that can be derived from correctness heatmaps include:

- **High-confidence regions:** Cells where the model assigns a near-total probability to correct actions indicate strong alignment with expected behavior;
- **Uncertain regions:** Areas with distributed probability mass among multiple actions suggest indecision, which could stem from ambiguous training data or suboptimal model reasoning;
- **Error-prone zones:** Cells where incorrect actions receive a significant portion of the probability mass highlight potential weaknesses in the model’s decision-making.

These correctness heatmaps provide a structured way to assess the agent’s decision-making process, offering insights into how well it aligns with expected behavior. By analyzing these visualizations, we can identify areas where the model exhibits high confidence in correct actions and regions where it

shows uncertainty or inconsistencies. This analysis will be further explored in Chapter 6, where we examine the broader implications of these findings.

6 Results Discussion

In this chapter, we analyze the results of our experiments, focusing on how the agent performs in different maps and goals configuration to evaluate the agent’s ability to navigate the map and successfully complete its tasks. We examine how the placement of goal tiles affects decision-making, assess whether the model struggles with retrieving relevant information, and compare the performance in the same slice of different maps.

Our primary objective is to evaluate the agent’s ability to navigate the map and successfully complete pickup and delivery tasks.

In tasks where the agent’s goal is to pick up a parcel, the target tile corresponds to the one containing the parcel. Conversely, in delivery tasks, the goal tile is the specific location where the agent must deliver the parcel.

The placement of goal tiles within the game map was carefully designed to ensure consistency and meaningful evaluation across both pickup and delivery tasks. Specifically, we aimed to use the same goal tiles for both objectives, allowing for direct comparisons between the two. In the pickup scenario, the goal is always explicitly stated at the end of the prompt, making it immediately available to the LLM. However, in the delivery scenario, the agent must retrieve the delivery location from the provided map description, requiring it to process and extract the relevant information effectively.

To evaluate how well the model handles goal retrieval, we selected three distinct goal positions: the top-right, center, and bottom-right cells of the map. These placements ensure that the goal appears in different parts of the map description inside the prompt, allowing us to prove again what has been demonstrated in the “needle in a haystack” article[1].

In the final section, we will discuss the different LLMs used in the agent’s decision-making process. As a general observation, GPT-3.5 performed worse compared to the more advanced GPT-4o variants. Among the newer models, GPT-4o and GPT-4o-mini demonstrated similar performance, with both outperforming GPT-3.5. Due to budget constraints, the majority of our analysis was conducted using GPT-4o-mini, as it provides a cost-effective yet high-quality alternative. However, our qualitative findings can be reasonably extended to GPT-4o as well, given their comparable performance, and to all models with similar performance and/or architecture. Each API call had an average input size of approximately 250 tokens, with only a single token as output. The cost per call was approximately \$0.000038.

6.1 Map Orientation

In this section, we analyze the impact of map orientation on the agent’s decision-making process. Since the prompt did not explicitly reference any specific orientation, the model had to infer it based solely on the provided map description and its training data. By examining the model’s outputs, we can determine its perceived orientation of the map and assess whether any biases emerge.

To investigate this, we analyzed the data using two different origin conventions. The raw data was structured with the (0,0) coordinate in the top-left corner, but we also transformed the map to simulate a bottom-left origin by adjusting all coordinates accordingly. We then ran the agent in both configurations—one using the original top-left origin and another using the simulated bottom-left origin—to compare how the agent’s actions varied under different map orientations.

As illustrated in Figure 6.3, the heatmaps of the agent’s actions appear nearly identical, except for a 90-degree rotation. The small differences between the two cases can be attributed to the way the data was structured in the prompt, which may have influenced the model’s text generation.

To further analyze the effect of orientation, we examined the correctness heatmaps for both configurations, as shown in Figure 6.6. The results reveal a clear bias toward the top-left origin orientation, which we refer to as the “programming origin,” in contrast to the “Cartesian origin” commonly used



Figure 6.1: Bottom Left Origin Orientation

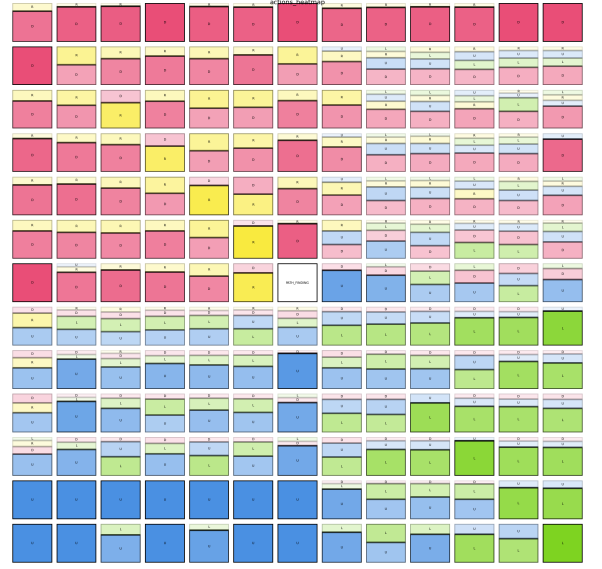


Figure 6.2: Top Left Origin Orientation

Figure 6.3: Heatmaps showing actions' probability with different map orientations

in mathematical contexts.

This bias may stem from the way the map was presented to the model. In our specific implementation, the map was formatted as a list of tiles extracted from a minimally edited JSON file. Given that common data structures in computer science often follow a top-left origin convention, it is likely that the LLM was implicitly influenced by its prior knowledge from programming-related contexts.

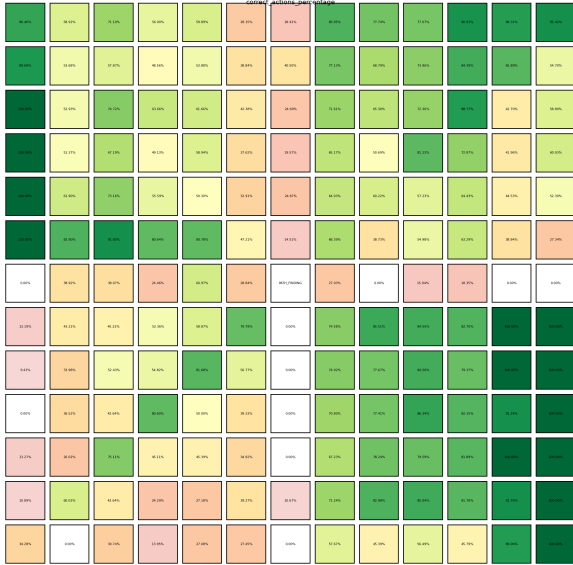


Figure 6.4: Bottom Left Origin Orientation

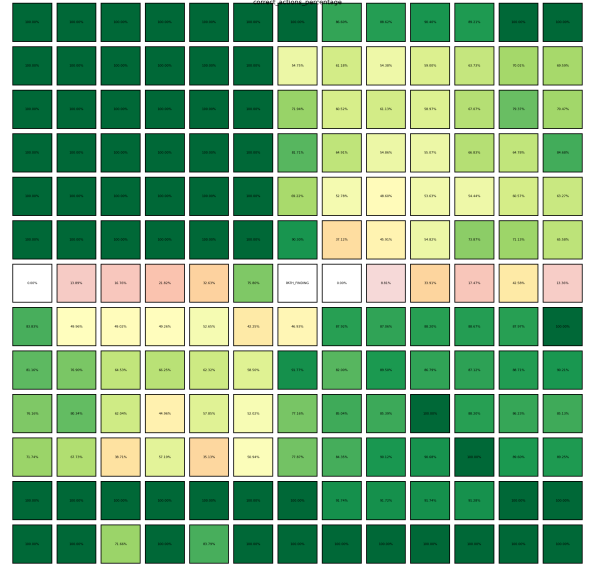


Figure 6.5: Top Left Origin Orientation

Figure 6.6: Correctness Heatmaps for the same map with different orientations

6.1.1 Comparison of Orientations

When comparing the two orientations, we found that in the top-left origin version, 154 tiles had one of the correct actions as the most probable choice. In contrast, in the bottom-left origin version, only 104 tiles had a correct action as the highest-probability choice. The comparison between the two orientations reveals a clear preference for the top-left origin. The model performed significantly better

when using this orientation, as shown by the following accuracy metrics:

- **Top-left origin:** 92% of the tiles had the correct action as the most probable one;
- **Bottom-left origin:** 62% of the tiles had the correct action as the most probable one;
- **Top 3 actions comparison:** 99% vs. 93% of the tiles contained the correct action within the top three choices.

These results strongly suggest that the model is inherently biased toward the top-left origin orientation. This finding highlights the potential influence of data structure representations on LLM-based decision-making and suggests that models trained on structured data formats may develop spatial preferences that impact their performance in spatial reasoning tasks.

6.2 Stateless

As previously discussed in Section 4.4.4, a stateless agent operates without any memory of past actions or previous states of the environment. In this context, every decision is made independently, relying solely on the information provided within a single prompt.

Technically, each call to the Large Language Model contains only the current state of the environment, without any reference to past states or prior actions. For every decision, a new conversation instance is initiated, making the agent unable to build an internal representation of the map or track past movements.

One of the major challenges faced by a stateless agent is the need to infer its position and the map layout solely based on the current prompt. This limitation leads to several difficulties, such as:

- **Inability to Track Progress:** Since the agent does not retain memory, it cannot recognize previously visited locations, often resulting in repetitive movements or getting stuck in loops.
- **Increased Uncertainty:** The LLM must deduce the correct course of action based only on the available snapshot of the environment, leading to occasional misinterpretations.
- **Higher Error Rates:** Compared to a stateful approach, where an agent can accumulate knowledge over time, the stateless method is more prone to making incorrect decisions, especially in larger maps.

Despite these challenges, implementing a stateless agent serves as an important step toward understanding the inherent uncertainty in LLM-based decision-making. The results obtained from this approach provide a useful baseline for evaluating also the performance of the stateful agent.

The stateless agent follows predefined prompt templates, as described in Sections 5.1.1 and 5.1.2. These prompts encapsulate all the necessary information about the current state and available actions within a single request.

To evaluate the performance of the stateless agent, we analyze heatmaps generated from experiments on different map sizes. The following sections present the results for various scenarios.

6.2.1 Pickup Goal at the Center

Figures 6.10 and 6.14 illustrate the heatmaps for maps of sizes 5×5 , 7×7 , and 13×13 , with the *pickup* goal placed at the center.

From the heatmaps, we can observe a consistent pattern across different map sizes. The top-left quadrant tends to show red and yellow as prominent colors and the bottom-left quadrant is predominantly blue. The top-right quadrant exhibits significant uncertainty with the KnowNo framework that keeps few actions in most of the cells, while the bottom-right quadrant generally is green and blue (colors representing every actions, as explained in Section 5.3.1).

To further examine the correctness of the agent’s decisions, we analyze the correctness heatmaps shown in Figure 6.14.

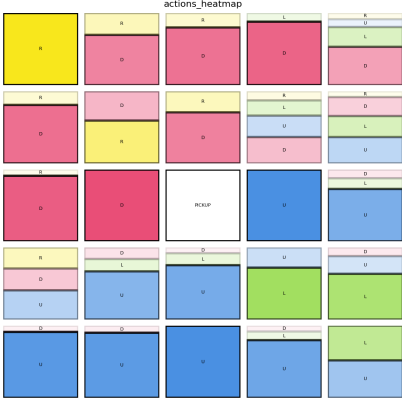


Figure 6.7: 5×5



Figure 6.8: 7×7

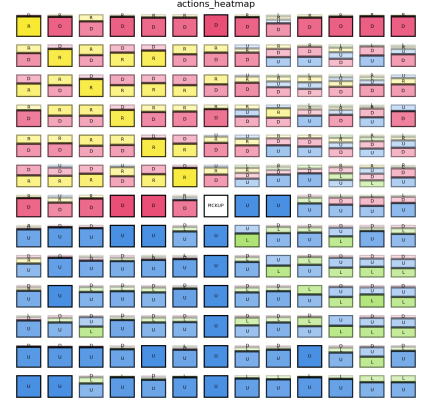


Figure 6.9: 13×13

Figure 6.10: Heatmaps for stateless agent with *pickup* goal in the center of different map sizes



Figure 6.11: 5×5



Figure 6.12: 7×7

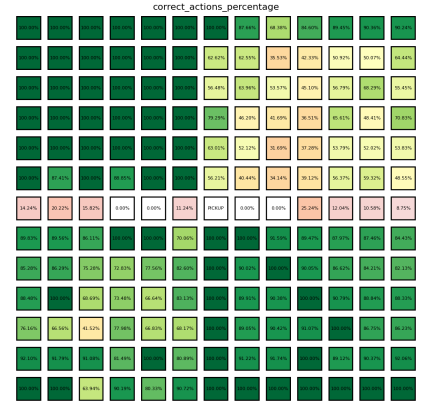


Figure 6.13: 13×13

Figure 6.14: Correctness heatmaps for stateless agent with pickup goal at the center of different map sizes.

The correctness analysis reveals the following trends:

- The top-left and bottom-right quadrants exhibit the highest certainty, with the former being almost perfect in every cell;
- The top-right and bottom-left quadrants are more uncertain, with the top-right being the least reliable;
- Along the entire row and the upper half of the column containing the goal, the correctness tends to be lower, with several cells having 0% correctness, meaning that the only correct action was discarded by KnowNo.

Table 6.2 presents numerical performance metrics across different map sizes. The number of times the correct action is present in the X actions with the highest probability is shown in the *topX* columns, while their percentage is in the *topX%* columns.

	top1	top2	top3	top1%	top2%	top3%
5×5	19	22	22	0.792	0.917	0.917
7×7	37	40	41	0.771	0.833	0.854
13×13	125	161	162	0.744	0.958	0.964

Table 6.1: Performance metrics for different map sizes - central pickup goal

The results indicate that:

- Smaller maps tend to yield a higher probability of selecting the correct action as the top-ranked choice.
- While absolute correctness decreases with increasing map size, the relative performance remains fairly stable.

Even if not visible in the heatmaps, the goal tile almost always had the correct action as the only one kept after KnowNo framework, and rarely it kept other actions but the correct one had a better probability by far.

6.2.2 Deliver Goal at the Center

A similar analysis can be conducted for the case where the *delivery* goal is placed at the center of the map. Figures 6.18 and 6.22 show the heatmaps for different map sizes in this configuration.

Compared to the pickup task, the delivery task introduces slightly more uncertainty. This uncertainty arises because, in our setup, the goal tile is not explicitly marked as a destination in the prompt but must instead be inferred from the map description. Again, the LLM needs to recognize that it has arrived at the correct location based solely on relative positioning within the map, without any persistent memory of past decisions.

From the heatmaps, we observe similar trends as seen in the pickup case:

- The top-left and bottom-right quadrants show a higher concentration of correct probability;
- The top-right quadrant, similar to the pickup scenario, exhibits more variability and uncertainty;
- The row and column containing the goal continue to demonstrate reduced correctness.

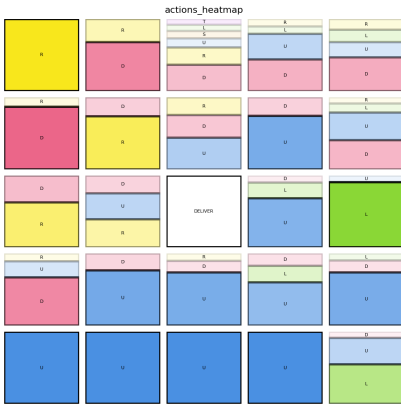


Figure 6.15: 5×5



Figure 6.16: 7×7

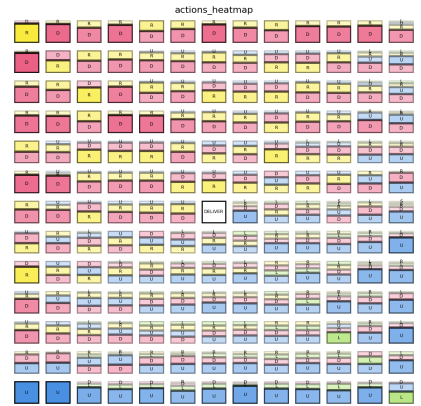


Figure 6.17: 13×13

Figure 6.18: Heatmaps for stateless agent with deliver goal in the center of different map sizes

Moreover, the values in Table 6.2 reveals that as the map size increases, the probability of selecting the correct action as the top-ranked choice decreases. However, the overall trend remains consistent, suggesting that while uncertainty increases with map size, the general patterns of decision-making remain largely stable across different scales.

6.2.3 Pickup and Deliver Goals in Different Map Sections

Beyond testing the performance of the stateless agent scenarios with the goal at the center, we also examined its ability to handle pickup and delivery goals positioned in different sections of the map.



Figure 6.19: 5×5



Figure 6.20: 7×7

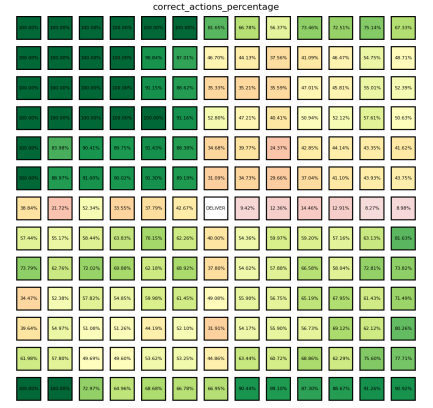


Figure 6.21: 13×13

Figure 6.22: Heatmaps for stateless agent with deliver goal in the center of different map sizes

	top1	top2	top3	top1%	top2%	top3%
5×5	20	24	24	0.833	1.000	1.000
7×7	43	47	48	0.896	0.979	1.000
13×13	125	161	162	0.744	0.958	0.964

Table 6.2: Performance metrics for different map sizes - central deliver goal

Figures 6.25 and 6.28 illustrate the heatmaps for cases where the pickup and delivery locations are in the top-right and bottom-right corners, respectively.

By analyzing the performance in these different regions, we aim to understand whether the placement of the goal influences the decision-making accuracy of the LLM. Since the map is embedded within a structured text prompt, the location of the goal may impact how the LLM processes the spatial relationships between different tiles as explained in Section 5.2.5.

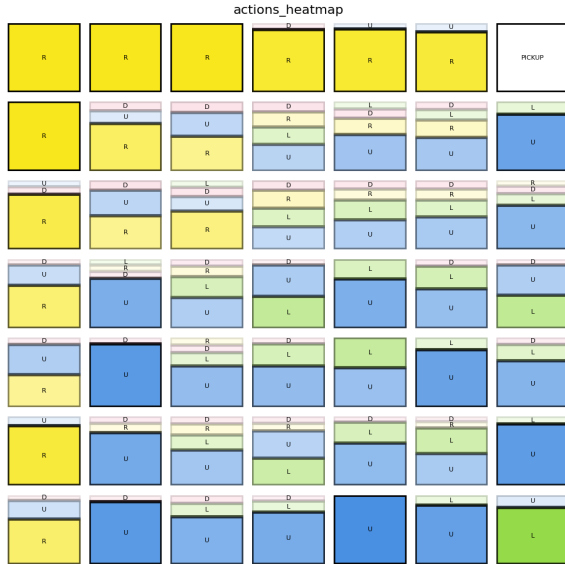


Figure 6.23: Pickup Top Right

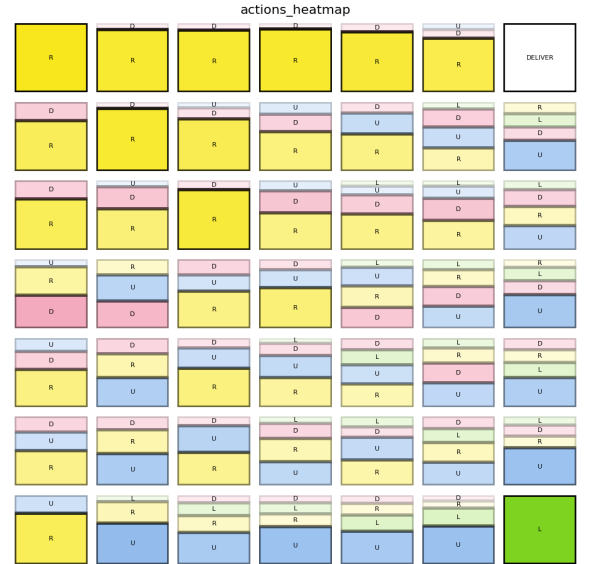


Figure 6.24: Deliver Top Right

Figure 6.25: Heatmaps for stateless agent with pickup and deliver goals in the top right corner of the map

Looking at the correctness heatmaps in Figures 6.31 and 6.34, we observe similar patterns, where



Figure 6.26: Pickup Bottom Right

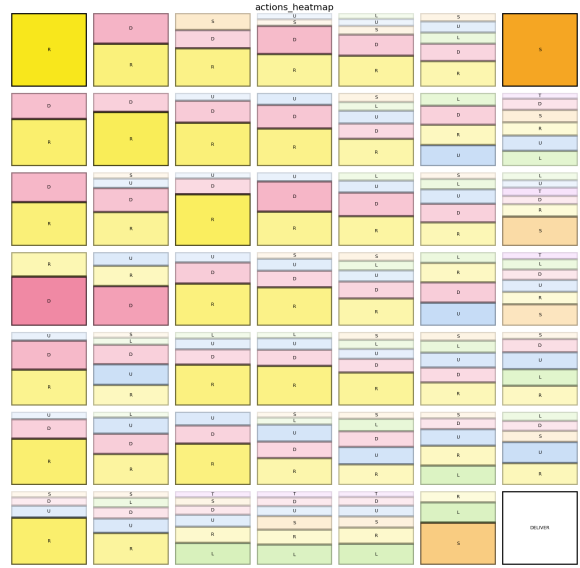


Figure 6.27: Deliver Bottom Right

Figure 6.28: Heatmaps for stateless agent with pickup and deliver goals in the bottom right corner of the map

the characteristic drop in correctness along the goal row and column persists.

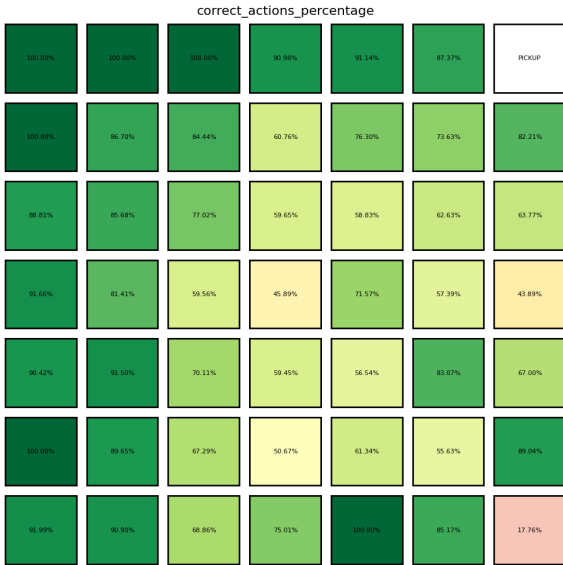


Figure 6.29: Correctness Pickup Top Right



Figure 6.30: Correctness Deliver Top Right

Figure 6.31: Correctness heatmaps for stateless agent with pickup and deliver goals in the top right corner of the map

These alternative goal placements provide valuable insights into the agent's decision-making and can be seen as subsets of the larger map. For example, the top-left portion of a 13×13 map with the goal in the center can be viewed as a 7×7 map with the goal in the bottom-right cell. At the same time, the bottom-left portion of the 13×13 map with the center goal corresponds to a 7×7 map with the top-right goal, highlighting again the same behavior of the quadrants described above.

This perspective brings us to the following topic, where a slice with the same size from different maps is compared.

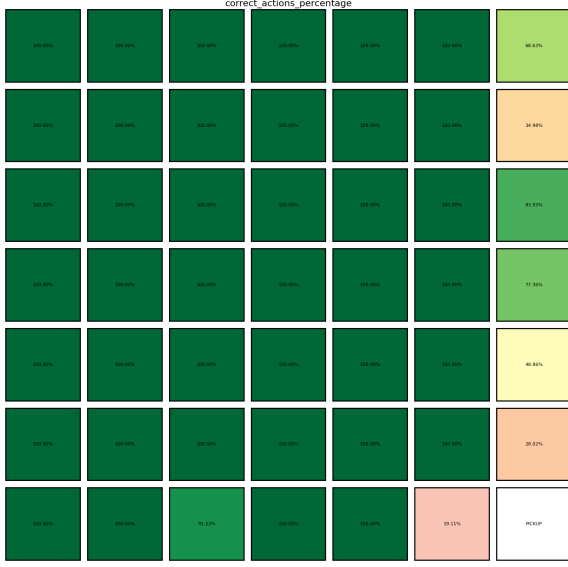


Figure 6.32: Correctness Pickup Bottom Right

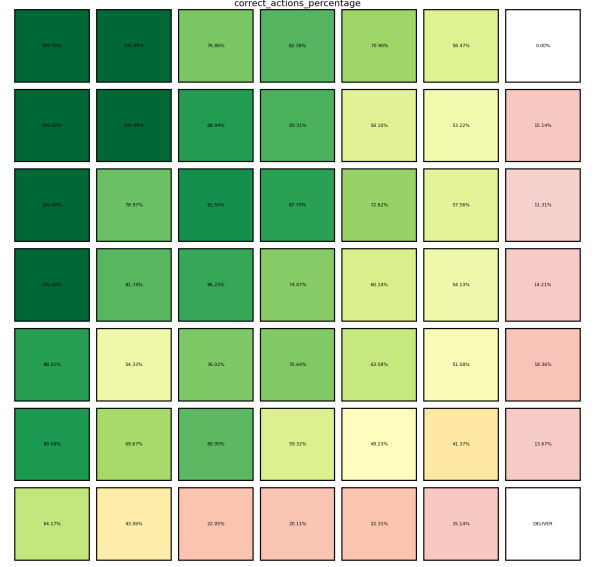


Figure 6.33: Correctness Deliver Bottom Right

Figure 6.34: Correctness heatmaps for stateless agent with pickup and deliver goals in the bottom right corner of the map

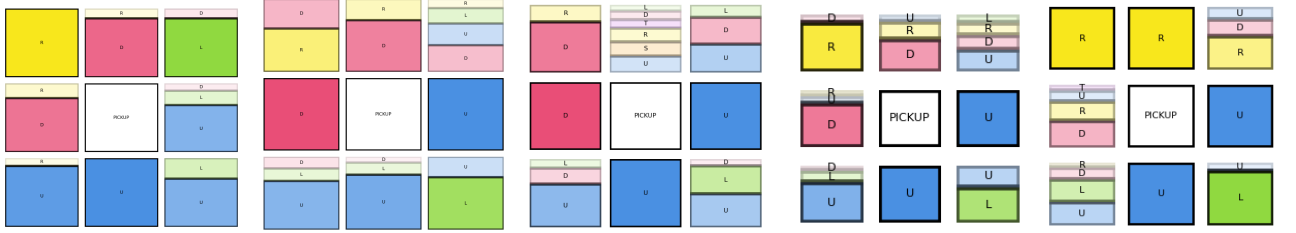


Figure 6.35: 3x3 Figure 6.36: 5 × 5 Figure 6.37: 7 × 7 Figure 6.38: 13 × 13 Figure 6.39: 21x21

Figure 6.40: Heatmaps of the 8 tiles around the goal for different map sizes

It is already visible in Figure 6.18 but it can be seen better in figure 6.40. Even if the values in the cells around the goal are quite similar, but not identical. This discrepancy may arise because the LLM distributes its attention across the entire prompt and smaller maps occupy a smaller portion of the prompt.

6.2.4 Goal position comparison

A final analysis was conducted to determine whether the placement of the delivery goal in different sections of the map affects the overall accuracy of the stateless agent.

Table 6.3 presents performance metrics for the 21x21 grid in three configurations:

- **Top Right (TR):** The delivery goal is positioned in the top-right corner;
- **Center (CN):** The delivery goal is placed in the center of the map;
- **Bottom Right (BR):** The delivery goal is located in the bottom-right corner.

	top1	top2	top3	top1%	top2%	top3%
21x21_TR	378	429	433	0.859	0.975	0.984
21x21_CN	294	394	428	0.668	0.895	0.973
21x21_BR	371	416	417	0.843	0.945	0.948

Table 6.3: Performance metrics for different 21x21 configurations

From the data, we can draw several conclusions:

- The top-right and bottom-right configurations exhibit similar performance, with top1 accuracy at 85.9% and 84.3%, respectively, as stated in Section 5.2.5;
- The center goal configuration shows a notable drop in top1 accuracy to 66.8%, confirming our earlier observation that finding a specific information hidden in the middle of a long text is more challenging;
- The difference between the top2 and top3 accuracy rates across the three configurations indicates that, even when the correct action is not the most probable choice, it is still frequently within the top-ranked options, meaning the model maintains a reasonable degree of decision-making reliability.

These results reinforce the hypothesis that stateless agents struggle most when they must infer goal locations based on spatial relationships rather than relying on explicit goal markers. Additionally, it confirms that the layout of the prompt itself (how the map is structured in textual form) plays a role in how effectively the LLM can interpret spatial cues.

6.3 Stateful

In the stateful approach, we keep track of the chat history, meaning that, at every step, all the past states and actions are available to the LLM. This allows the agent to “build” an internal representation of the map and track its movements over time, enabling more informed decision-making. The stateful approach integrates historical context, which is especially beneficial for tasks where understanding the progression of events is crucial; in this case it allows the agent to better understand the map orientation.

One key advantage of a stateful agent is its ability to recover from errors. With access to previous interactions, the LLM can identify patterns such as repetitive loops or suboptimal paths. Recognizing these patterns, the agent is able to adjust its strategy, replan its path, and ultimately increase the chance of reaching the goal even after encountering execution errors or obstacles.

This translates to the fact that the agent is able to reach the goal even when the stateless output does not keep any correct action in a specific tile.

However, a stateful design does carry potential challenges. One significant issue is the token limit imposed by the LLM. As the conversation history grows, managing and condensing the context without losing essential details becomes critical. Strategies such as summarizing less relevant parts of the dialogue or periodically resetting portions of the context are often necessary to remain within token constraints while still maintaining effective decision-making. They have been tested in the form of sending the map in the first call, then wait and send it again every 5 or 10 messages, but it either made the agent not performing as desired or the sending of the map was still so frequent that the token limit was reached fast.

6.3.1 Path Visualization

A dedicated visualization script has been developed to illustrate the agent’s path on the map. This script identifies all optimal paths connecting the start and the goal, selects the one sharing the highest number of common cells with the agent’s path, and then calculates both the percentage of overlapping cells and the relative difference in path length. Figure 6.41 displays an example output, highlighting that although the agent eventually reaches the goal, it encounters difficulties navigating in the vicinity of the goal, as discussed in Section 6.2.1.

6.4 “Pathfinding”

Up to this point, we have observed that, in almost every instance, the KnowNo framework consistently discarded actions related to the agent’s goals, such as picking up and delivering parcels. To determine whether these goal-related actions still contributed significantly to the agent’s uncertainty, we designed an experiment where the objective was reduced to simply reaching a specific tile.

To implement this, we removed the *pickup* and *delivery* actions from the prompt while keeping all other conditions unchanged. The goal was set up as just reaching a specific tile, similarly to what

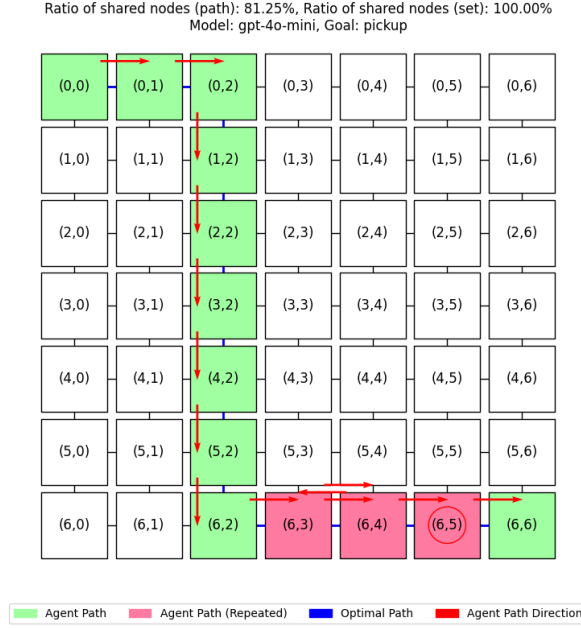


Figure 6.41: Path visualization for stateful agent with pickup goal in the bottom right cell

the “Pickup goal” requires. By doing so, we aimed to isolate and examine the role of these actions in influencing the agent’s uncertainty. The results of this experiment are presented in Figure 6.44. This approach aligns with our systematic methodology, as previously discussed in Section 4.5, where we progressively simplified the agent’s goals to identify potential limitations in decision-making.

A meaningful comparison can be drawn between these results and those obtained from a scenario in which the agent operated on a 5×5 grid with a pickup goal located at the center. Since both setups shared the same prompt structure, this comparison allows us to assess whether the presence of goal-related actions had a significant impact on the agent’s uncertainty. The results from the original 5×5 pickup scenario, previously discussed, can be found in Figure 6.7 and Figure 6.11.

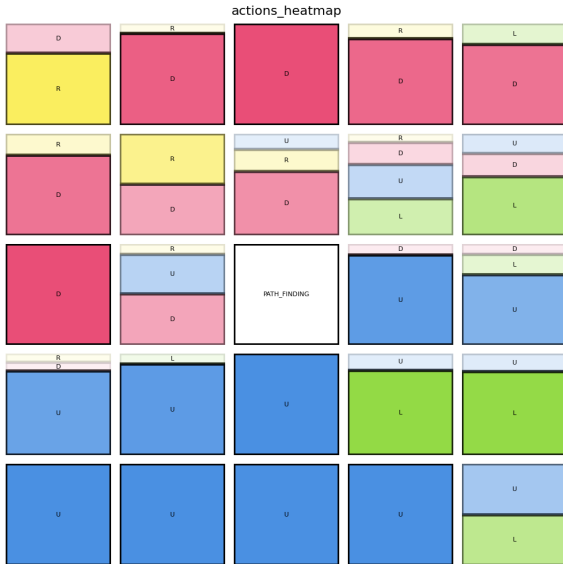


Figure 6.42: Heatmap for pathfinding

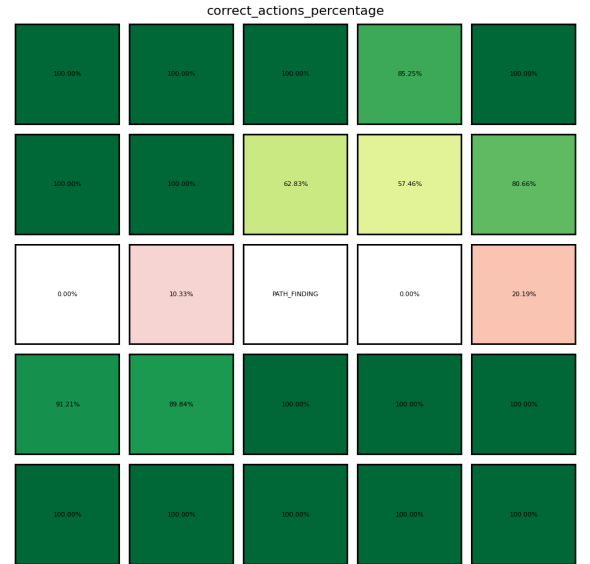


Figure 6.43: Correctness for pathfinding

Figure 6.44: Heatmaps for stateless agent with pathfinding goal

The outcomes of both experiments are remarkably similar, revealing identical areas of high and low uncertainty. This suggests that goal-related actions are not the primary source of uncertainty in the agent’s decision-making process. Furthermore, it reinforces the idea that the structure of the prompt, carefully designed based on the literature reviewed as explained in Section 5.2, is not a key

factor contributing to uncertainty. Instead, the observed uncertainty may be more closely linked to the inherent limitations of the LLM itself.

6.5 Stateless & Stateful - Performance Summary

Experimental results in this study indicate that the stateful configuration outperforms the stateless one in terms of accuracy and goal attainment, particularly in larger maps. The ability to trace its previous actions and incorporate historical context leads to a more adaptive and flexible approach, which is evident in tasks requiring dynamic path corrections and strategic foresight. During the first stages of the experiments, where we tried the limits of the token input size, the stateless agent was not able at all to achieve any goal in maps bigger than 21x21, while the stateful agent was on the correct path but blocked by the token limit. This is a clear demonstration of the stateful agent’s superior performance in handling complex spatial navigation tasks.

Additionally, the stateful method aligns well with the inherent strengths of large language models as few-shot learners [4]. The continuous integration of historical data and real-time inputs allows the agent to refine its internal representation of the map, which significantly improves its overall navigational performance. This dynamic adjustment provides resilience against uncertainties common in spatial navigation tasks.

Practically speaking, the stateful agent demonstrates a clear advantage in navigating the “problematic” areas of the maps, particularly in the top-right quadrant, where the stateless version exhibits significant uncertainty. This improvement arises from the stateful agent’s ability to retain memory of its past actions, enabling it to recognize and correct mistakes. If it takes a step in the wrong direction, it can backtrack to a previous cell and attempt an alternative route, ultimately increasing its chances of reaching the goal efficiently.

However, as the map size increases, both versions encounter growing difficulties, with the top-right quadrant consistently presenting the most challenges. This issue stems from the nature of decision-making in that region: many cells offer multiple actions with nearly equal probabilities, making it harder for the agent to determine the optimal path. While the relative frequency of these problematic cells remains constant as the map scales up, the absolute number of such cells increases substantially. For example, in a small 3x3 map, a single problematic cell is manageable, as the agent can quickly backtrack and reach the goal. In contrast, in a much larger 21x21 map, there could be as many as 49 problematic cells. Even though this proportion is the same in percentage terms, the larger map exacerbates the challenge. With more problematic cells spread across a wider area, the agent is more likely to make a series of unfortunate steps, potentially getting trapped in a suboptimal loop and failing to recover efficiently. This issue is particularly severe for the stateless agent, but even the stateful agent, despite its ability to backtrack, can struggle to escape if the problematic region is too large.

A summary of their performance are presented in the table 6.4, where the agent was tasked to pickup a parcel in the bottom left cell and started in the top right cell (since it is the most problematic area). The recorded action count includes the pickup action.

Map Size	Stateless	Stateful	Manhattan Distance
13 × 13	41 actions	39 actions	25 actions
7 × 7	27 actions	21 actions	13 actions
5 × 5	14 actions	11 actions	9 actions
3x3	11 actions	9 actions	5 actions

Table 6.4: Number of actions in different map sizes using GPT-4o-mini

In summary, although the stateful approach can introduce computational overhead due to increased context size, its advantages in enhancing navigational accuracy, recovery from errors, and adaptive decision-making render it a vital component in modern agent design. Still, the limitation of the token limit is a significant drawback that may be overcome by more powerful models, that may introduce better results out of the box as well as a longer context window.

Still, we think this is a great result since we are relying only on the generative capabilities of these models. This demonstrates the potential of LLMs to perform complex spatial reasoning tasks without any specialized training or fine-tuning, solely based on their inherent language understanding and generation abilities.

6.6 Insights from the Closest Cell to the Goal Approach

As detailed in Section 4.5, this approach aimed to simplify the agent’s decision-making process by breaking it down into two sequential steps: (1) identifying the best adjacent cell to move toward and (2) selecting the correct action to reach that cell. The motivation behind this decomposition was to reduce the complexity of the global path-finding task and test whether the agent could make more reliable local decisions.

Although this method was not central to our primary experiments, we include it here for completeness, as analyzing its shortcomings provided useful insights into the agent’s limitations.

Despite its structured nature, this method introduced new challenges rather than resolving the agent’s decision-making issues. In particular, two key limitations emerged:

- **Ambiguity in selecting the “best” cell:** In many cases, more than one neighboring cell reduced the distance to the goal, making it unclear which one the agent should prioritize and more difficult to us to record structured data to analyze. As illustrated in Figure 6.45, the decision was not always straightforward, and slight variations in prompt phrasing could lead the LLM to prefer different paths, even when multiple valid choices existed;
- **Compounding errors in high-uncertainty areas:** This two-step approach inadvertently introduced an additional layer of failure. If the agent misidentified the best neighboring cell, it would inherently lead to a suboptimal move, even if the second step—choosing the action—was executed perfectly. This doubled the impact of errors, particularly in ambiguous or low-information regions of the map, where uncertainty was already high.

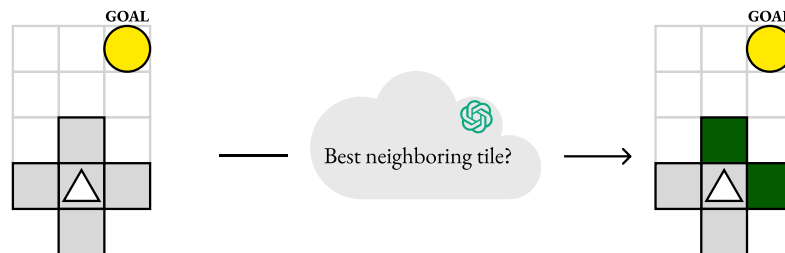


Figure 6.45: Two-steps decision making problem

6.7 Model Comparison

As mentioned in previous sections, our experiments primarily utilized GPT-4o-mini, a smaller variant of the GPT-4o model. This choice was mainly influenced by cost considerations and its performance relative to the full GPT-4o model. We also evaluated GPT-3.5-turbo, an older model, and found that while its overall capabilities were not necessarily worse, it was less likely to select the correct action as the one with the highest probability compared to GPT-4o-mini.

Results obtained with GPT-4o-mini closely align with those of GPT-4o, ensuring that all discussions in this document remain applicable to both models. From Table 6.5, we can see that GPT-3.5-turbo selects the correct action as the highest probability option less frequently than the other models, though its top-2 and top-3 probabilities remain comparable.

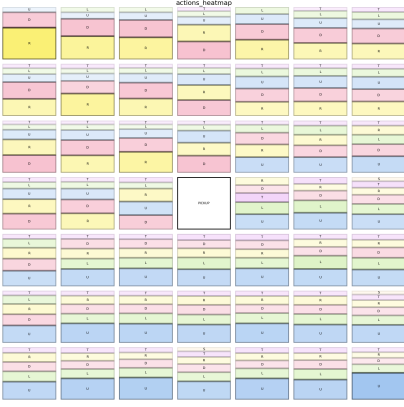


Figure 6.46: GPT-3.5-turbo



Figure 6.47: GPT-4o



Figure 6.48: GPT-4o-mini

Figure 6.49: Heatmaps for different GPT models

	top1	top2	top3	top1%	top2%	top3%
GPT-3.5-turbo	34	46	48	0.708	0.958	1.000
GPT-4o	37	44	44	0.771	0.917	0.917
GPT-4o-mini	37	40	41	0.771	0.833	0.854

Table 6.5: Comparison between different GPT models in a 7×7 map

However, Figures 6.49 and 6.53 clearly show that GPT-3.5-turbo exhibits greater overall uncertainty. Still, this demonstrates that our experimental setup is adaptable to different models and effectively structured to assess the capabilities of LLMs in this task.

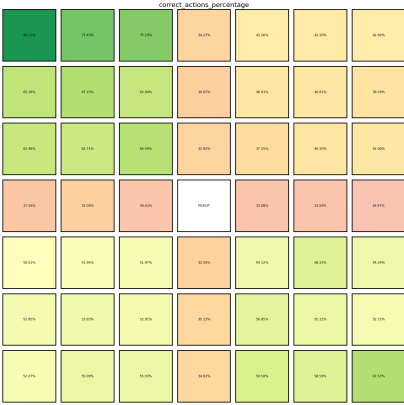


Figure 6.50: GPT-3.5-turbo

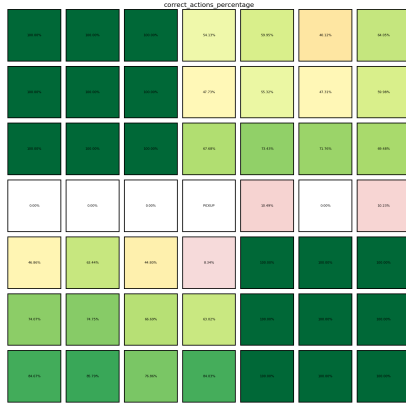


Figure 6.51: GPT-4o

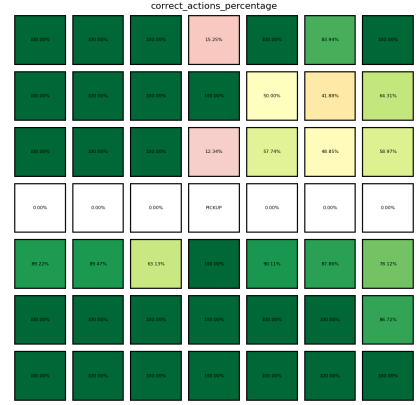


Figure 6.52: GPT-4o-mini

Figure 6.53: Correctness Heatmaps for different GPT models

7 Conclusions

In this work, we have explored the challenges and opportunities associated with navigating grid-based environments using Large Language Models to solve logistics tasks. Through a series of experiments, we examined both stateless and stateful agent configurations, highlighting the strengths and limitations of each approach, the former relying only on a snapshot of the current environment state, while the latter maintains a record of past interactions.

Our analysis revealed that stateless agent, although capable of reaching the goal, struggles with uncertainty when inferring spatial relationships solely from the prompt as the map size grows. This approach has been the foundation of this thesis, since it allowed us to map the uncertainty over the entire environment, highlighting limitations and biases linked to the environment description.

In contrast, the stateful agent, which maintains a record of past actions and environments, shows an improved ability to handle complex navigation tasks by leveraging historical context. This feature allows it to better understand and adapt to changes in the environment, recover from suboptimal paths and make more informed decisions, particularly in dense or larger grid scenarios.

Our findings underscore the importance of leveraging historical context when dealing with dynamic and complex environments. While stateless approaches offer simplicity, they may fall short in scenarios where resolving subtle spatial cues is critical. The stateful framework, despite imposing additional constraints due to token limitations, demonstrates significant advantages in accuracy and adaptability. Future research should focus on optimizing the balance between context size and performance, possibly through advanced summarization techniques or selective memory retention.

Some limitations to this thesis include the possibility that the prompt, although designed in accordance with the literature, may not be optimal. Additionally, the behavior of OpenAI’s API `logit bias` parameter is subject to changes, which could impact consistency. Most of the tests were conducted with a single model, and the results were then generalized to all models, which may not fully capture variations. Furthermore, the environment was not truly dynamic; while it was treated as such by not relying on a specific parser, the agent was not actually tested in a fast-changing environment. Context size limitations also proved to be crucial in the agent’s performance, and some biases were observed in the agent’s understanding of the environment.

Overall, we were able to demonstrate the potential of LLMs’ generative capabilities in grid-based navigation tasks, highlighting the importance of context and memory in achieving optimal performance, while also identifying common weaknesses among existing models.

There are also several promising directions for future research, such as:

- Exploring alternative ways to leverage uncertainty at each step, such as using it as a trigger for an expert system to take over or incorporating it as a variable in a more complex function that adjusts the output dynamically (e.g., if the delta between the two highest probabilities actions is below a certain epsilon, it could be treated as *no uncertainty*);
- Investigating the use of multimodal models to achieve a richer representation of the environment, as analyzing a 2D map as a mere sequence may significantly reduce performance;
- Testing reasoning-based models, which, while requiring more time per step, during which the environment may change, could benefit from future faster models capable of reasoning more efficiently;
- Integrate Retrieval-Augmented Generation techniques to enhance the agent’s decision-making process by adding precomputed favorable paths that can be used to guide the agent towards the goal in specific scenarios;

- Conducting a broader comparison across all the latest available models, including open-source ones, to establish a dedicated benchmark for evaluating their capabilities;
- Continuously testing newer models, given that performance improvements have been observed across successive versions.

Bibliography

- [1] Aparna Dhinakaran. The Needle In a Haystack Test: Evaluating the Performance of LLM RAG Systems. <https://arize.com/blog-course/the-needle-in-a-haystack-test-evaluating-the-performance-of-llm-rag-systems/>. Accessed 11/03/2025.
- [2] Evan Becker and Stefano Soatto. Cycles of thought: Measuring llm confidence through stable explanations, 2024.
- [3] Bernd Bohnet, Azade Nova, Aaron T Parisi, Kevin Swersky, Katayoon Goshvadi, Hanjun Dai, Dale Schuurmans, Noah Fiedel, and Hanie Sedghi. Exploring and benchmarking the planning capabilities of large language models, 2024.
- [4] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. Language models are few-shot learners, 2020.
- [5] Davide Modolo. GitHub Project Repository. https://github.com/davidemodolo/master_the_sis_project/. 12/03/2025.
- [6] DeepSeek-AI, Daya Guo, Dejian Yang, Haowei Zhang, Junxiao Song, Ruoyu Zhang, Runxin Xu, Qihao Zhu, Shirong Ma, Peiyi Wang, Xiao Bi, Xiaokang Zhang, Xingkai Yu, Yu Wu, Z. F. Wu, Zhibin Gou, Zhihong Shao, Zhuoshu Li, Ziyi Gao, Aixin Liu, Bing Xue, Bingxuan Wang, Bochao Wu, Bei Feng, Chengda Lu, Chenggang Zhao, Chengqi Deng, Chenyu Zhang, Chong Ruan, Damai Dai, Deli Chen, Dongjie Ji, Erhang Li, Fangyun Lin, Fucong Dai, Fuli Luo, Guangbo Hao, Guanting Chen, Guowei Li, H. Zhang, Han Bao, Hanwei Xu, Haocheng Wang, Honghui Ding, Huajian Xin, Huazuo Gao, Hui Qu, Hui Li, Jianzhong Guo, Jiashi Li, Jiawei Wang, Jingchang Chen, Jingyang Yuan, Junjie Qiu, Junlong Li, J. L. Cai, Jiaqi Ni, Jian Liang, Jin Chen, Kai Dong, Kai Hu, Kaige Gao, Kang Guan, Kexin Huang, Kuai Yu, Lean Wang, Lecong Zhang, Liang Zhao, Litong Wang, Liyue Zhang, Lei Xu, Leyi Xia, Mingchuan Zhang, Minghua Zhang, Minghui Tang, Meng Li, Miaojuan Wang, Mingming Li, Ning Tian, Panpan Huang, Peng Zhang, Qiancheng Wang, Qinyu Chen, Qiushi Du, Ruiqi Ge, Ruisong Zhang, Ruizhe Pan, Runji Wang, R. J. Chen, R. L. Jin, Ruyi Chen, Shanghao Lu, Shangyan Zhou, Shanhuang Chen, Shengfeng Ye, Shiyu Wang, Shuiping Yu, Shunfeng Zhou, Shuting Pan, S. S. Li, Shuang Zhou, Shaoqing Wu, Shengfeng Ye, Tao Yun, Tian Pei, Tianyu Sun, T. Wang, Wangding Zeng, Wanjia Zhao, Wen Liu, Wenfeng Liang, Wenjun Gao, Wenqin Yu, Wentao Zhang, W. L. Xiao, Wei An, Xiaodong Liu, Xiaohan Wang, Xiaokang Chen, Xiaotao Nie, Xin Cheng, Xin Liu, Xin Xie, Xingchao Liu, Xinyu Yang, Xinyuan Li, Xuecheng Su, Xuheng Lin, X. Q. Li, Xiangyue Jin, Xiaojin Shen, Xiaosha Chen, Xiaowen Sun, Xiaoxiang Wang, Xinnan Song, Xinyi Zhou, Xianzu Wang, Xinxia Shan, Y. K. Li, Y. Q. Wang, Y. X. Wei, Yang Zhang, Yanhong Xu, Yao Li, Yao Zhao, Yaofeng Sun, Yaohui Wang, Yi Yu, Yichao Zhang, Yifan Shi, Yiliang Xiong, Ying He, Yishi Piao, Yisong Wang, Yixuan Tan, Yiyang Ma, Yiyuan Liu, Yongqiang Guo, Yuan Ou, Yuduan Wang, Yue Gong, Yuheng Zou, Yujia He, Yunfan Xiong, Yuxiang Luo, Yuxiang You, Yuxuan Liu, Yuyang Zhou, Y. X. Zhu, Yanhong Xu, Yanping Huang, Yaohui Li, Yi Zheng, Yuchen Zhu, Yunxian Ma, Ying Tang, Yukun Zha, Yuting Yan, Z. Z. Ren, Zehui Ren, Zhangli Sha, Zhe Fu, Zhean Xu,

- Zhenda Xie, Zhengyan Zhang, Zhewen Hao, Zhicheng Ma, Zhigang Yan, Zhiyu Wu, Zihui Gu, Zijia Zhu, Zijun Liu, Zilin Li, Ziwei Xie, Ziyang Song, Zizheng Pan, Zhen Huang, Zhipeng Xu, Zhongyu Zhang, and Zhen Zhang. Deepseek-r1: Incentivizing reasoning capability in llms via reinforcement learning, 2025.
- [7] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
 - [8] Fabian Gloeckle, Badr Youbi Idrissi, Baptiste Rozière, David Lopez-Paz, and Gabriel Synnaeve. Better and faster large language models via multi-token prediction, 2024.
 - [9] Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014.
 - [10] Jonathan Ho, Ajay Jain, and Pieter Abbeel. Denoising diffusion probabilistic models, 2020.
 - [11] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural Computation*, 9:1735–1780, 11 1997.
 - [12] Wenlong Huang, Pieter Abbeel, Deepak Pathak, and Igor Mordatch. Language models as zero-shot planners: Extracting actionable knowledge for embodied agents, 2022.
 - [13] Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. Survey of hallucination in natural language generation. *ACM Computing Surveys*, 55(12):1–38, March 2023.
 - [14] Diederik P Kingma and Max Welling. Auto-encoding variational bayes, 2022.
 - [15] Aobo Kong, Shiwan Zhao, Hao Chen, Qicheng Li, Yong Qin, Ruiqi Sun, Xin Zhou, Enzhi Wang, and Xiaohang Dong. Better zero-shot reasoning with role-play prompting, 2024.
 - [16] Alex Krizhevsky, Ilya Sutskever, and Geoffrey Hinton. Imagenet classification with deep convolutional neural networks. *Neural Information Processing Systems*, 25, 01 2012.
 - [17] John Lafferty, Andrew McCallum, and Fernando Pereira. Conditional random fields: Probabilistic models for segmenting and labeling sequence data. pages 282–289, 01 2001.
 - [18] Patrick Lewis, Ethan Perez, Aleksandra Piktus, Fabio Petroni, Vladimir Karpukhin, Naman Goyal, Heinrich Küttler, Mike Lewis, Wen tau Yih, Tim Rocktäschel, Sebastian Riedel, and Douwe Kiela. Retrieval-augmented generation for knowledge-intensive nlp tasks, 2021.
 - [19] Xuanqi Liu and Zhuotao Liu. Llms can understand encrypted prompt: Towards privacy-computing friendly transformers, 2023.
 - [20] Ashman Mehra, Snehanshu Saha, Vaskar Raychoudhury, and Archana Mathur. Deliverai: Reinforcement learning based distributed path-sharing network for food deliveries, 2024.
 - [21] Ahmed Njifenjou, Virgile Sucas, Bassam Jabaian, and Fabrice Lefèvre. Role-play zero-shot prompting with large language models for open-domain human-machine conversation, 2024.
 - [22] OpenAI, Josh Achiam, Steven Adler, Sandhini Agarwal, Lama Ahmad, Ilge Akkaya, Florencia Leoni Aleman, Diogo Almeida, Janko Altschmidt, Sam Altman, Shyamal Anadkat, Red Avila, Igor Babuschkin, Suchir Balaji, Valerie Balcom, Paul Baltescu, Haiming Bao, Mohammad Bavarian, Jeff Belgum, Irwan Bello, Jake Berdine, Gabriel Bernadett-Shapiro, Christopher Berner, Lenny Bogdonoff, Oleg Boiko, Madelaine Boyd, Anna-Luisa Brakman, Greg Brockman, Tim Brooks, Miles Brundage, Kevin Button, Trevor Cai, Rosie Campbell, Andrew Cann, Brittany Carey, Chelsea Carlson, Rory Carmichael, Brooke Chan, Che Chang, Fotis Chantzis, Derek Chen, Sully Chen, Ruby Chen, Jason Chen, Mark Chen, Ben Chess, Chester Cho, Casey Chu, Hyung Won Chung, Dave Cummings, Jeremiah Currier, Yunxing Dai, Cory Decareaux, Thomas

Degry, Noah Deutsch, Damien Deville, Arka Dhar, David Dohan, Steve Dowling, Sheila Dunning, Adrien Ecoffet, Atty Eleti, Tyna Eloundou, David Farhi, Liam Fedus, Niko Felix, Simón Posada Fishman, Juston Forte, Isabella Fulford, Leo Gao, Elie Georges, Christian Gibson, Vik Goel, Tarun Gogineni, Gabriel Goh, Rapha Gontijo-Lopes, Jonathan Gordon, Morgan Grafstein, Scott Gray, Ryan Greene, Joshua Gross, Shixiang Shane Gu, Yufei Guo, Chris Hallacy, Jesse Han, Jeff Harris, Yuchen He, Mike Heaton, Johannes Heidecke, Chris Hesse, Alan Hickey, Wade Hickey, Peter Hoeschele, Brandon Houghton, Kenny Hsu, Shengli Hu, Xin Hu, Joost Huizinga, Shantanu Jain, Shawn Jain, Joanne Jang, Angela Jiang, Roger Jiang, Haozhun Jin, Denny Jin, Shino Jomoto, Billie Jonn, Heewoo Jun, Tomer Kaftan, Lukasz Kaiser, Ali Kamali, Ingmar Kanitscheider, Nitish Shirish Keskar, Tabarak Khan, Logan Kilpatrick, Jong Wook Kim, Christina Kim, Yongjik Kim, Jan Hendrik Kirchner, Jamie Kiros, Matt Knight, Daniel Kokotajlo, Lukasz Kondraciuk, Andrew Kondrich, Aris Konstantinidis, Kyle Kosic, Gretchen Krueger, Vishal Kuo, Michael Lampe, Ikai Lan, Teddy Lee, Jan Leike, Jade Leung, Daniel Levy, Chak Ming Li, Rachel Lim, Molly Lin, Stephanie Lin, Mateusz Litwin, Theresa Lopez, Ryan Lowe, Patricia Lue, Anna Makanju, Kim Malfacini, Sam Manning, Todor Markov, Yaniv Markovski, Bianca Martin, Katie Mayer, Andrew Mayne, Bob McGrew, Scott Mayer McKinney, Christine McLeavey, Paul McMillan, Jake McNeil, David Medina, Aalok Mehta, Jacob Menick, Luke Metz, Andrey Mishchenko, Pamela Mishkin, Vinnie Monaco, Evan Morikawa, Daniel Mossing, Tong Mu, Mira Murati, Oleg Murk, David Mély, Ashvin Nair, Reiichiro Nakano, Rajeev Nayak, Arvind Neelakantan, Richard Ngo, Hyeonwoo Noh, Long Ouyang, Cullen O’Keefe, Jakub Pachocki, Alex Paino, Joe Palermo, Ashley Pantuliano, Giambattista Parascandolo, Joel Parish, Emy Parparita, Alex Passos, Mikhail Pavlov, Andrew Peng, Adam Perelman, Filipe de Avila Belbute Peres, Michael Petrov, Henrique Ponde de Oliveira Pinto, Michael, Pokorny, Michelle Pokrass, Vitchyr H. Pong, Tolly Powell, Alethea Power, Boris Power, Elizabeth Proehl, Raul Puri, Alec Radford, Jack Rae, Aditya Ramesh, Cameron Raymond, Francis Real, Kendra Rimbach, Carl Ross, Bob Rotsted, Henri Roussez, Nick Ryder, Mario Saltarelli, Ted Sanders, Shibani Santurkar, Girish Sastry, Heather Schmidt, David Schnurr, John Schulman, Daniel Selsam, Kyla Sheppard, Toki Sherbakov, Jessica Shieh, Sarah Shoker, Pranav Shyam, Szymon Sidor, Eric Sigler, Maddie Simens, Jordan Sitkin, Katarina Slama, Ian Sohl, Benjamin Sokolowsky, Yang Song, Natalie Staudacher, Felipe Petroski Such, Natalie Summers, Ilya Sutskever, Jie Tang, Nikolas Tezak, Madeleine B. Thompson, Phil Tillet, Amin Tootoonchian, Elizabeth Tseng, Preston Tuggle, Nick Turley, Jerry Tworek, Juan Felipe Cerón Uribe, Andrea Vallone, Arun Vijayvergiya, Chelsea Voss, Carroll Wainwright, Justin Jay Wang, Alvin Wang, Ben Wang, Jonathan Ward, Jason Wei, CJ Weinmann, Akila Welihinda, Peter Welinder, Jiayi Weng, Lilian Weng, Matt Wiethoff, Dave Willner, Clemens Winter, Samuel Wolrich, Hannah Wong, Lauren Workman, Sherwin Wu, Jeff Wu, Michael Wu, Kai Xiao, Tao Xu, Sarah Yoo, Kevin Yu, Qiming Yuan, Wojciech Zaremba, Rowan Zellers, Chong Zhang, Marvin Zhang, Shengjia Zhao, Tianhao Zheng, Juntang Zhuang, William Zhuk, and Barret Zoph. Gpt-4 technical report, 2024.

- [23] Alec Radford, Jong Wook Kim, Tao Xu, Greg Brockman, Christine McLeavey, and Ilya Sutskever. Robust speech recognition via large-scale weak supervision, 2022.
- [24] Alec Radford and Karthik Narasimhan. Improving language understanding by generative pre-training. 2018.
- [25] Alec Radford, Jeff Wu, Rewon Child, David Luan, Dario Amodei, and Ilya Sutskever. Language models are unsupervised multitask learners. 2019.
- [26] Aditya Ramesh, Mikhail Pavlov, Gabriel Goh, Scott Gray, Chelsea Voss, Alec Radford, Mark Chen, and Ilya Sutskever. Zero-shot text-to-image generation, 2021.
- [27] Anand S. Rao and Michael P. Georgeff. BDI agents: From theory to practice. In Victor R. Lesser and Les Gasser, editors, *1st International Conference on Multi Agent Systems (ICMAS 1995)*, pages 312–319, San Francisco, CA, USA, 12-14 June 1995. The MIT Press.

- [28] Allen Z. Ren, Anushri Dixit, Alexandra Bodrova, Sumeet Singh, Stephen Tu, Noah Brown, Peng Xu, Leila Takayama, Fei Xia, Jake Varley, Zhenjia Xu, Dorsa Sadigh, Andy Zeng, and Anirudha Majumdar. Robots that ask for help: Uncertainty alignment for large language model planners, 2023.
- [29] Tom Silver, Soham Dan, Kavitha Srinivas, Joshua B. Tenenbaum, Leslie Pack Kaelbling, and Michael Katz. Generalized planning in pddl domains with pretrained large language models, 2023.
- [30] Tom Silver, Varun Hariprasad, Reece S Shuttleworth, Nishanth Kumar, Tomás Lozano-Pérez, and Leslie Pack Kaelbling. PDDL planning with pretrained large language models. In *NeurIPS 2022 Foundation Models for Decision Making Workshop*, 2022.
- [31] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023.
- [32] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, Ed H. Chi, Tatsunori Hashimoto, Oriol Vinyals, Percy Liang, Jeff Dean, and William Fedus. Emergent abilities of large language models, 2022.
- [33] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023.
- [34] Michael Wooldridge. *An Introduction to Multiagent Systems*. John Wiley & Sons, 1 edition, 2002. Chapter 2.
- [35] Miao Xiong, Zhiyuan Hu, Xinyang Lu, Yifei Li, Jie Fu, Junxian He, and Bryan Hooi. Can llms express their uncertainty? an empirical evaluation of confidence elicitation in llms, 2024.
- [36] Liping Yang and Ruishi Liang. An ai planning approach to emergency material scheduling using numerical pddl. In *Proceedings of the 2022 International Conference on Artificial Intelligence, Internet and Digital Economy (ICAID 2022)*, pages 47–54. Atlantis Press, 2022.

Appendix A Acronyms

- **AI**: Artificial Intelligence
- **LLM**: Large Language Model
- **RL**: Reinforcement Learning
- **PDDL**: Planning Domain Definition Language
- **DL**: Deep Learning
- **NN**: Neural Network
- **DNN**: Deep Neural Network
- **CNN**: Convolutional Neural Network
- **RNN**: Recurrent Neural Network
- **CP**: Conformal Prediction
- **ML**: Machine Learning
- **RLHF**: Reinforcement Learning from Human Feedback
- **RAG**: Retrieval-Augmented Generation
- **CoT**: Chain-of-Thought

Appendix B Prompts

Prompt

```
1 You are a delivery agent in a web-based game I am going to give you the raw
  information I receive from the server and the possible actions. You have to
  take (pickup) the parcel and ship (deliver) it in a delivery tile.
2 Don't loop using the same moves.
3 If the information does not change, it means you are choosing the wrong actions
  .
4 Raw 'onMap' response: {"width":2,"height":5,"tiles":[{"x":0,"y":0,"delivery":
  false}, {"x":0,"y":1,"delivery":true}, {"x":1,"y":0,"delivery":false}, {"x":1,"
  y":1,"delivery":false}, {"x":2,"y":0,"delivery":false}, {"x":2,"y":1,"delivery
  ":false}, {"x":3,"y":0,"delivery":false}, {"x":3,"y":1,"delivery":false}, {"x
  ":4,"y":0,"delivery":false}, {"x":4,"y":1,"delivery":false}]}
5
6 Raw 'onYou' response: {"id":"75d4e78ed8e","name":"raw_llm_agent","x":3,"y":1,"
  score":0}
7
8 Raw 'onParcelsSensing' response: [{"id":"p0","x":3,"y":0,"carriedBy":null,"
  reward":10}]
9
10 ACTIONS you can do:
11 U): move up
12 D): move down
13 L): move left
14 R): move right
15 T): take the parcel that is in your tile
16 S): ship a parcel (you must be in a delivery=true tile)
17 Don't explain the reasoning and don't add any comment, just provide the action.
18 What is your next action?
```

Listing B.1: Example of a prompt used in the second agent implementation, see Section 4.3

Prompt

```
1 You are a delivery agent in a web-based game and I want to test your ability.
   You are in a grid world (represented with a matrix) with some obstacles and
   some parcels to deliver.
2 Parcels are generated at random on random free spots.
3 The value of the parcels lowers as the time passes, so you should deliver them
   as soon as possible.
4 Your view of the world is limited to a certain distance, so you can only see
   the parcels and the delivery points that are close to you.
5 MAP:
6 1 1 1 1 1
7 1 1 P 1 1
8 1 1 1 1 1
9 2 A 1 1 1
10 1 1 1 1 1
11
12 LEGEND:
13 - A: you (the Agent) are in this position;
14 - 1: you can move in this position;
15 - 2: you can deliver a parcel in this position (and also move there);
16 - P: a parcel is in this position;
17 - X: you are in the same position of a parcel;
18 - Q: you are in the delivery/shipping point;
19
20 ACTIONS you can do:
21 - U: move up
22 - D: move down
23 - L: move left
24 - R: move right
25 - T: take a parcel
26 - S: ship a parcel
27
28 You have 1 parcels to deliver.
29 Important rules:
30 - If you have 0 parcels, you must look for the closest parcel to pick up.
31 - If you are going to deliver >0 parcels and on the way you find 1 parcel, you
   should go and pick it up before shipping.
32 - If you have at least 1 parcel, your goal should be to deliver it/them to the
   closest delivery point. The more parcels you have, the more important it is
   to deliver them as soon as possible.
33 - If you can't see any delivery point, just move around to explore the map
   until one enters your field of view, then go and deliver the parcels.
34 - If there is no parcel in the map, just move around to explore the map until
   one parcel spawns, then go and get it.
35
36 You want to maximize your score by delivering the most possible number of
   parcels. You can pickup multiple parcels and deliver them in the same
   delivery point.
37 Don't explain the reasoning and don't add any comment, just provide the action.
38 Try to not go back and forth, it's a waste of time, so use the conversation
   history to your advantage.
39 Example: if you want to go down, just answer 'D'.
40 The closest delivery point is right and up from you.
41 What is your next action?
```

Listing B.2: Example of a prompt used in the first agent implementation, see Section 4.2