

# Exploring the Use of LLMs for Agent Planning Strengths and Weaknesses

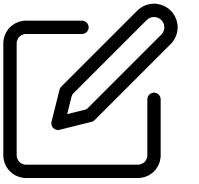
DAVIDE MODOLO  
20/03/2025  
*Supervisor PAOLO GIORGINI*



UNIVERSITÀ  
DI TRENTO

# Context

# Large Language Models' Capabilities



**Text Generation**

Their Scope

1. Context



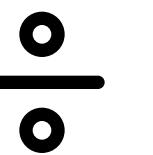
UNIVERSITÀ  
DI TRENTO

# Large Language Models' Capabilities



## Text Generation

Their Scope



## Math Reasoning

Emerging Behavior

### 1. Context

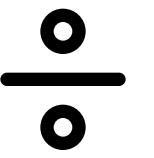


UNIVERSITÀ  
DI TRENTO

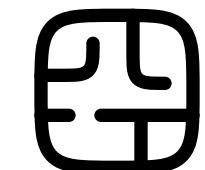
# Large Language Models' Capabilities



**Text Generation**  
Their Scope



**Math Reasoning**  
Emerging Behavior



**Planning Abilities**  
Emerging Behavior

1. Context

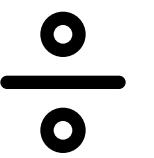


UNIVERSITÀ  
DI TRENTO

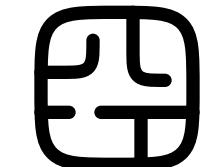
# Large Language Models' Capabilities



**Text Generation**  
Their Scope



**Math Reasoning**  
Emerging Behavior



**Planning Abilities**  
Emerging Behavior

...

**More**

1. Context



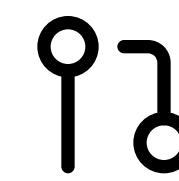
UNIVERSITÀ  
DI TRENTO

# LLM-based Planning



UNIVERSITÀ  
DI TRENTO

# LLM-based Planning

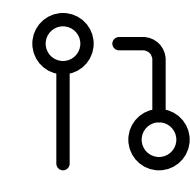


**Chain of  
Thought**  
Reasoning<sup>1</sup>

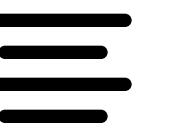
<sup>1</sup> *Chain-of-thought prompting elicits reasoning in large language models* - Wei et al., 2022



# LLM-based Planning



**Chain of  
Thought**  
Reasoning<sup>1</sup>



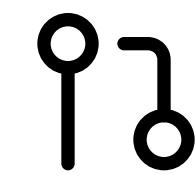
**Few-Shots**  
Prompting<sup>2</sup>

<sup>1</sup> *Chain-of-thought prompting elicits reasoning in large language models* - Wei et al., 2022

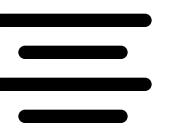
<sup>2</sup> *PDDL planning with pretrained large language models* - Silver et al., 2022



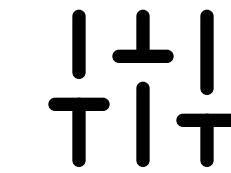
# LLM-based Planning



**Chain of Thought**  
Reasoning<sup>1</sup>



**Few-Shots**  
Prompting<sup>2</sup>



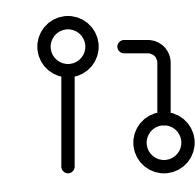
**Fine-Tuning**  
Models<sup>3</sup>

<sup>1</sup> *Chain-of-thought prompting elicits reasoning in large language models* - Wei et al., 2022

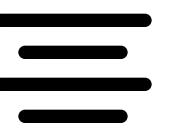
<sup>2</sup> *PDDL planning with pretrained large language models* - Silver et al., 2022

<sup>3</sup> *Unlocking Large Language Model's Planning Capabilities with Maximum Diversity Fine-tuning* - Wenjun Li et al., 2024

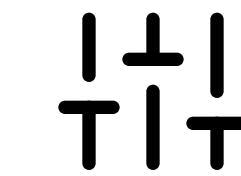
# LLM-based Planning



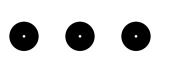
**Chain of Thought**  
Reasoning<sup>1</sup>



**Few-Shots**  
Prompting<sup>2</sup>



**Fine-Tuning**  
Models<sup>3</sup>



**More<sup>4</sup>**

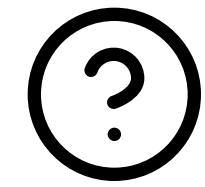
<sup>1</sup> *Chain-of-thought prompting elicits reasoning in large language models* - Wei et al., 2022

<sup>2</sup> *PDDL planning with pretrained large language models* - Silver et al., 2022

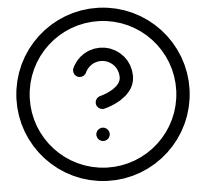
<sup>3</sup> *Unlocking Large Language Model's Planning Capabilities with Maximum Diversity Fine-tuning* - Wenjun Li et al., 2024

<sup>4</sup> *Efficient Sequential Decision Making with Large Language Models* - Dingyang Chen et al., 2024

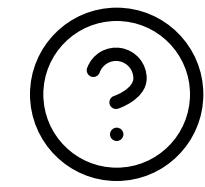
# Research Questions



What happens if we strip everything prior away?

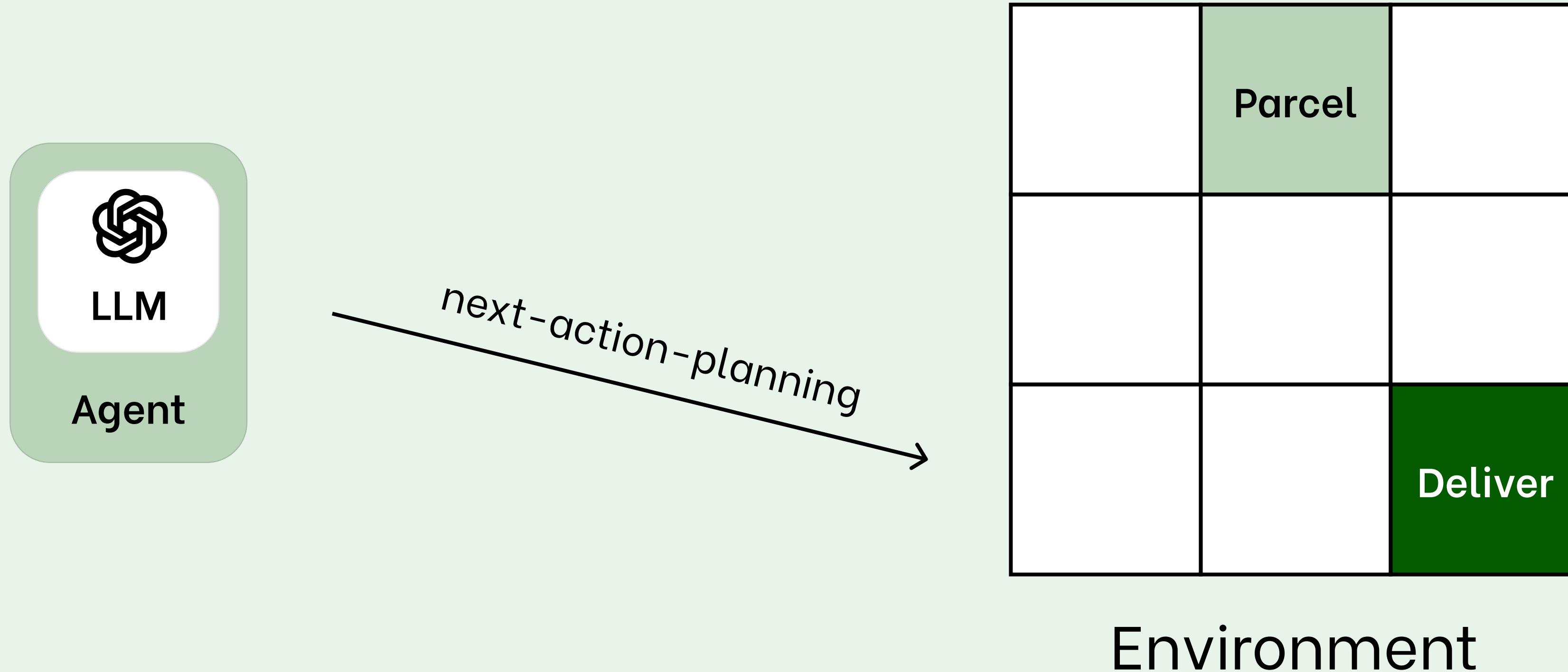


Can an LLM, without additional training or frameworks, effectively plan and navigate in an unknown environment?



How well can LLMs' generative capabilities make sequential decisions in such environments?

# Idea



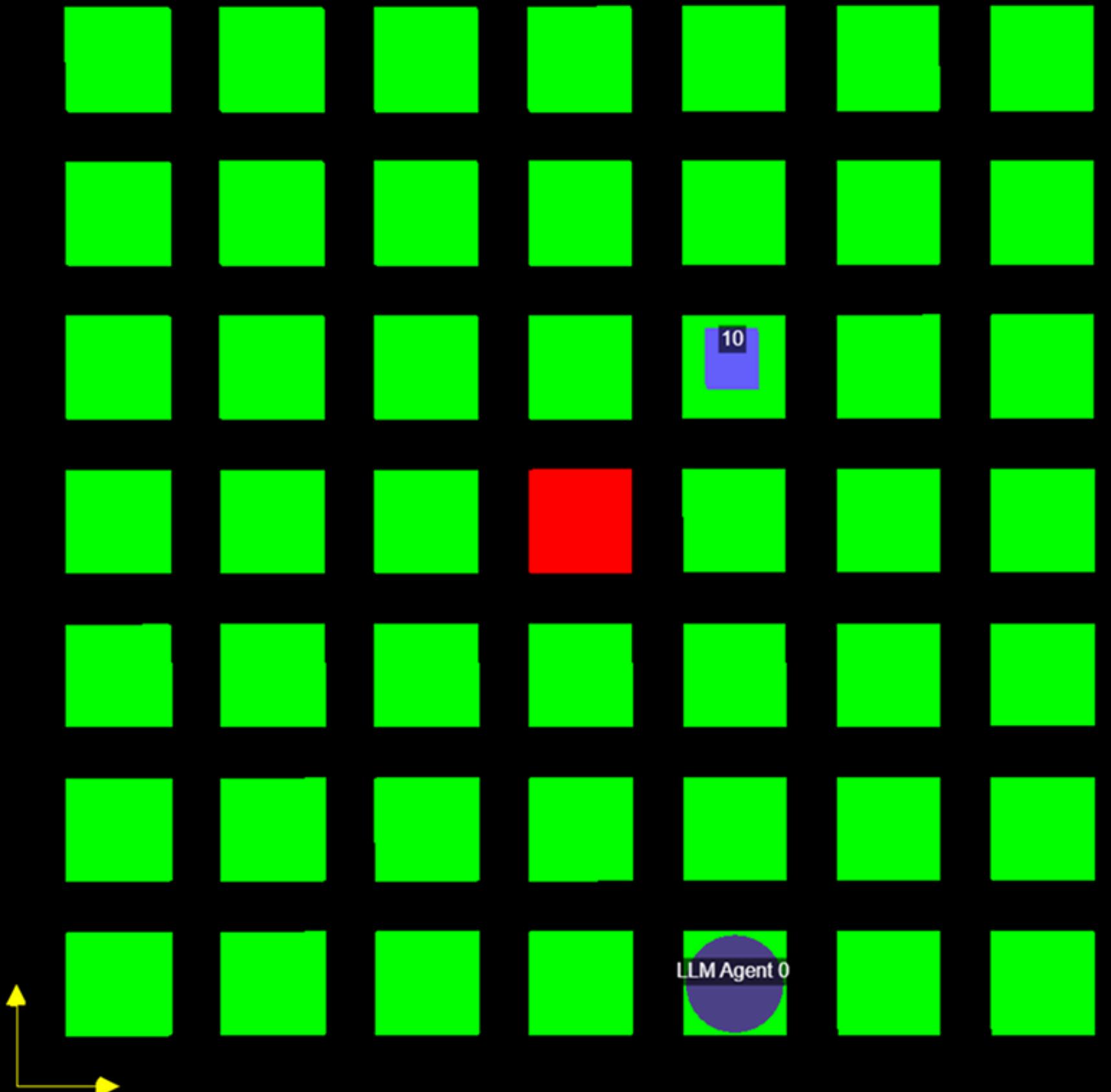
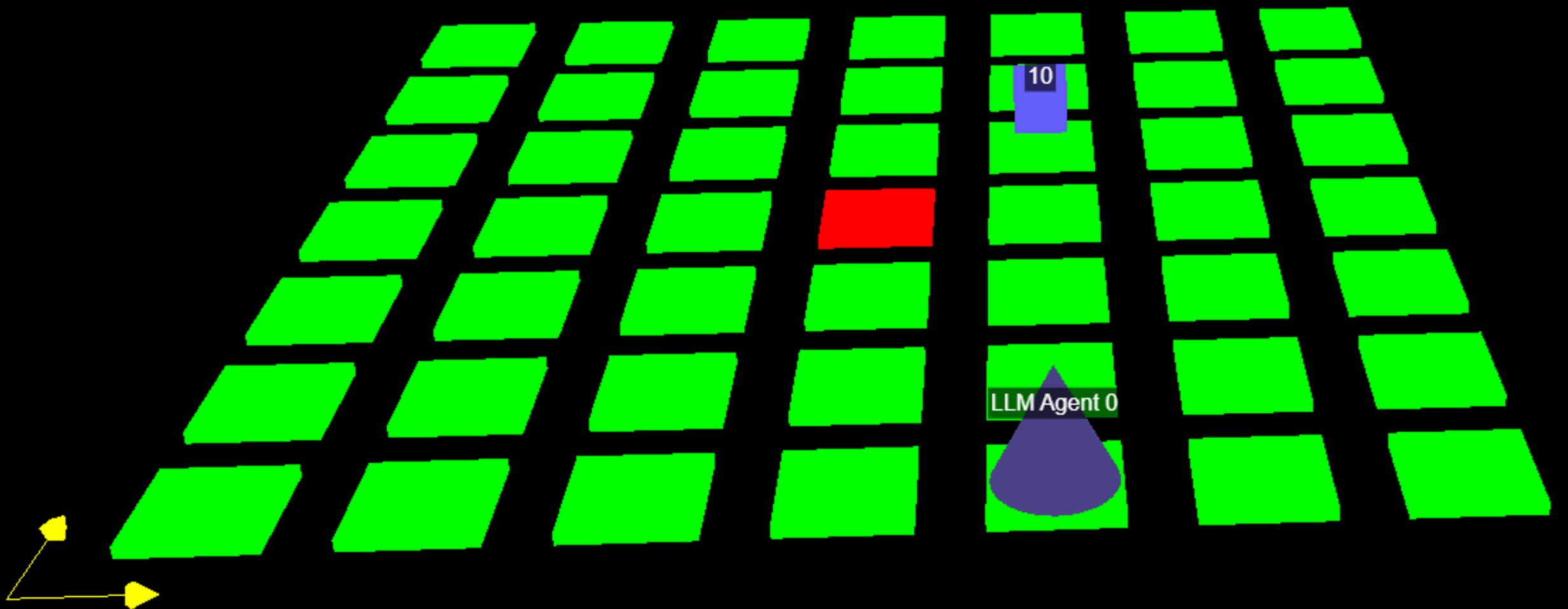
## 1. Context

# Setting

# Deliveroo.js

## Educational Game

Parcels spawn all around the map. The goal is to pickup and deliver them.

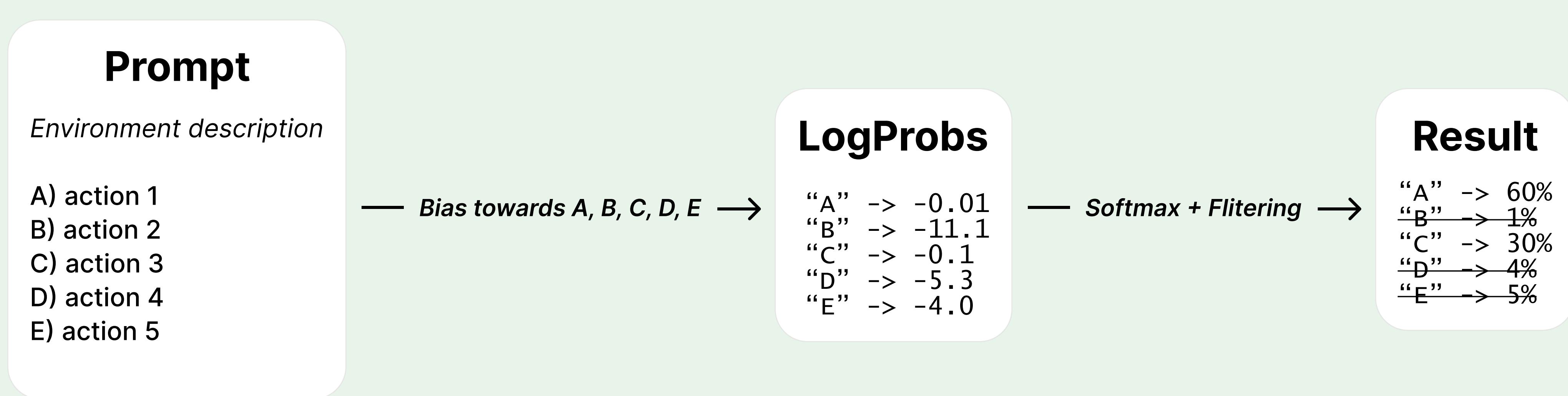


## 2. Setting



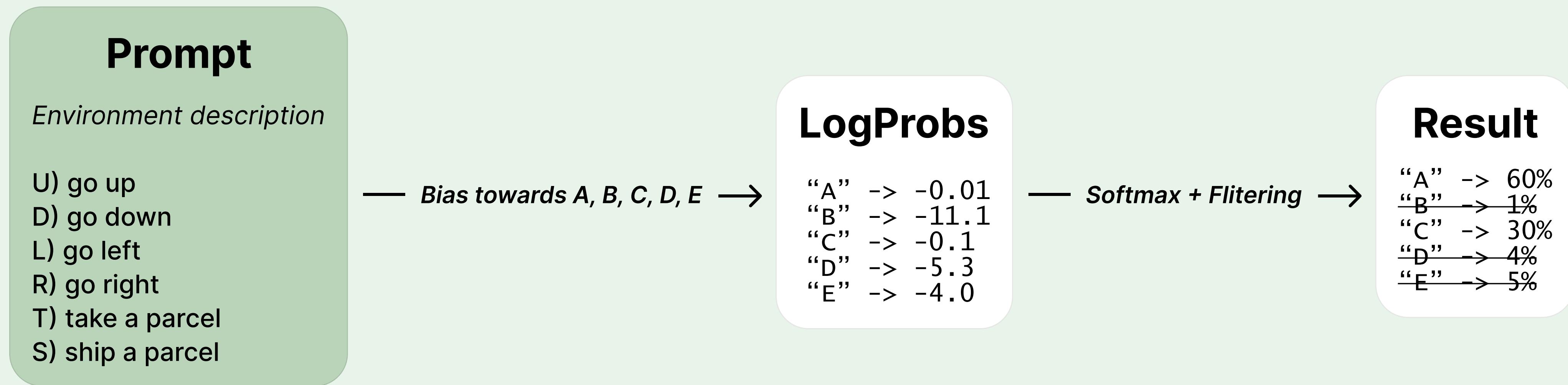
UNIVERSITÀ  
DI TRENTO

# KnowNo Uncertainty Framework<sup>1</sup>



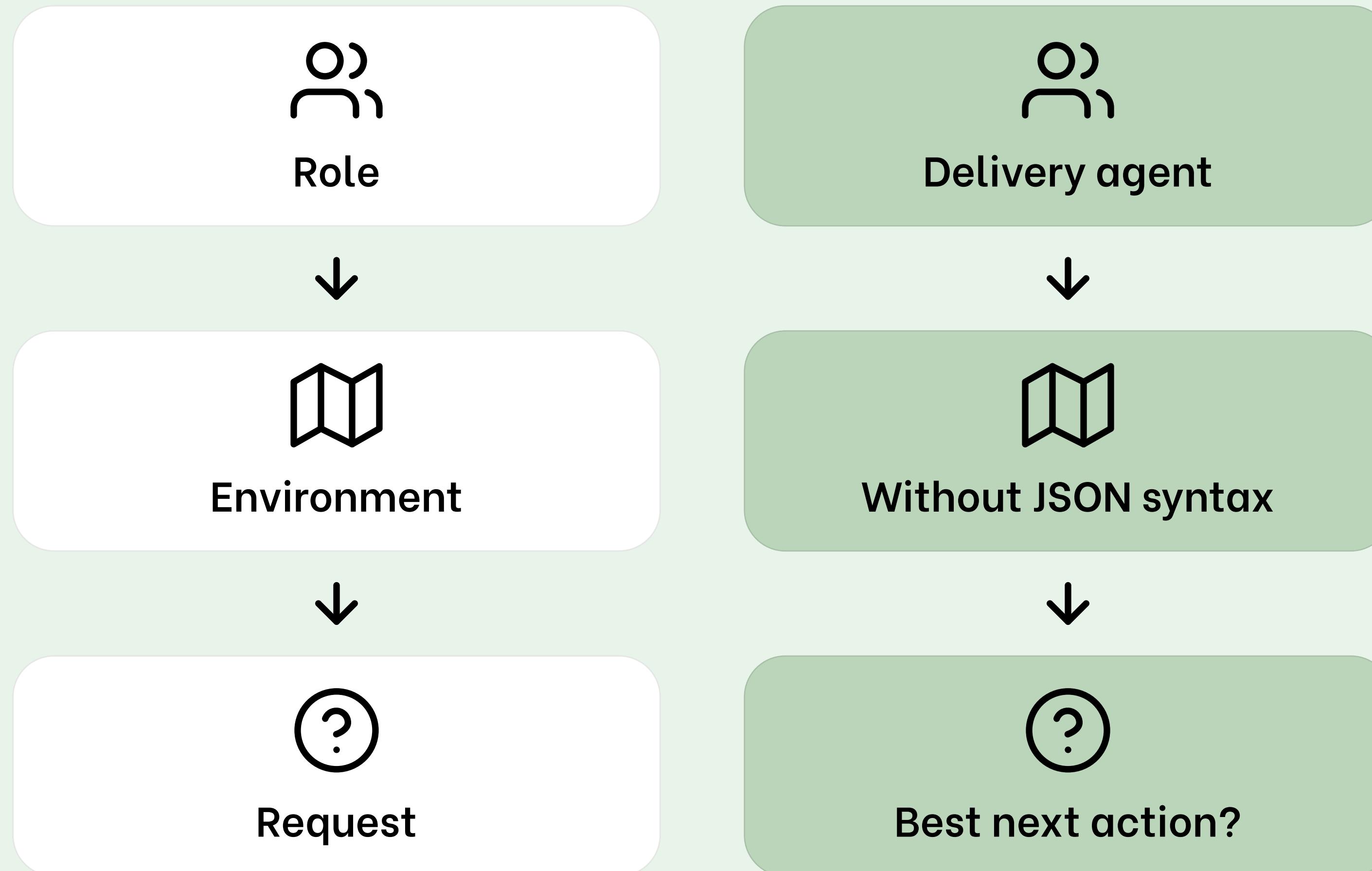
<sup>1</sup>Robots That Ask For Help: Uncertainty Alignment for Large Language Model Planners – Allen Z. Ren et al., 2023

# KnowNo Uncertainty Framework<sup>1</sup>



<sup>1</sup>Robots That Ask For Help: Uncertainty Alignment for Large Language Model Planners – Allen Z. Ren et al., 2023

# Prompting Strategy



# Model

	GPT-4o	GPT-4o-mini
top1%	77%	84%
top2%	95%	91%
top3%	96%	92%

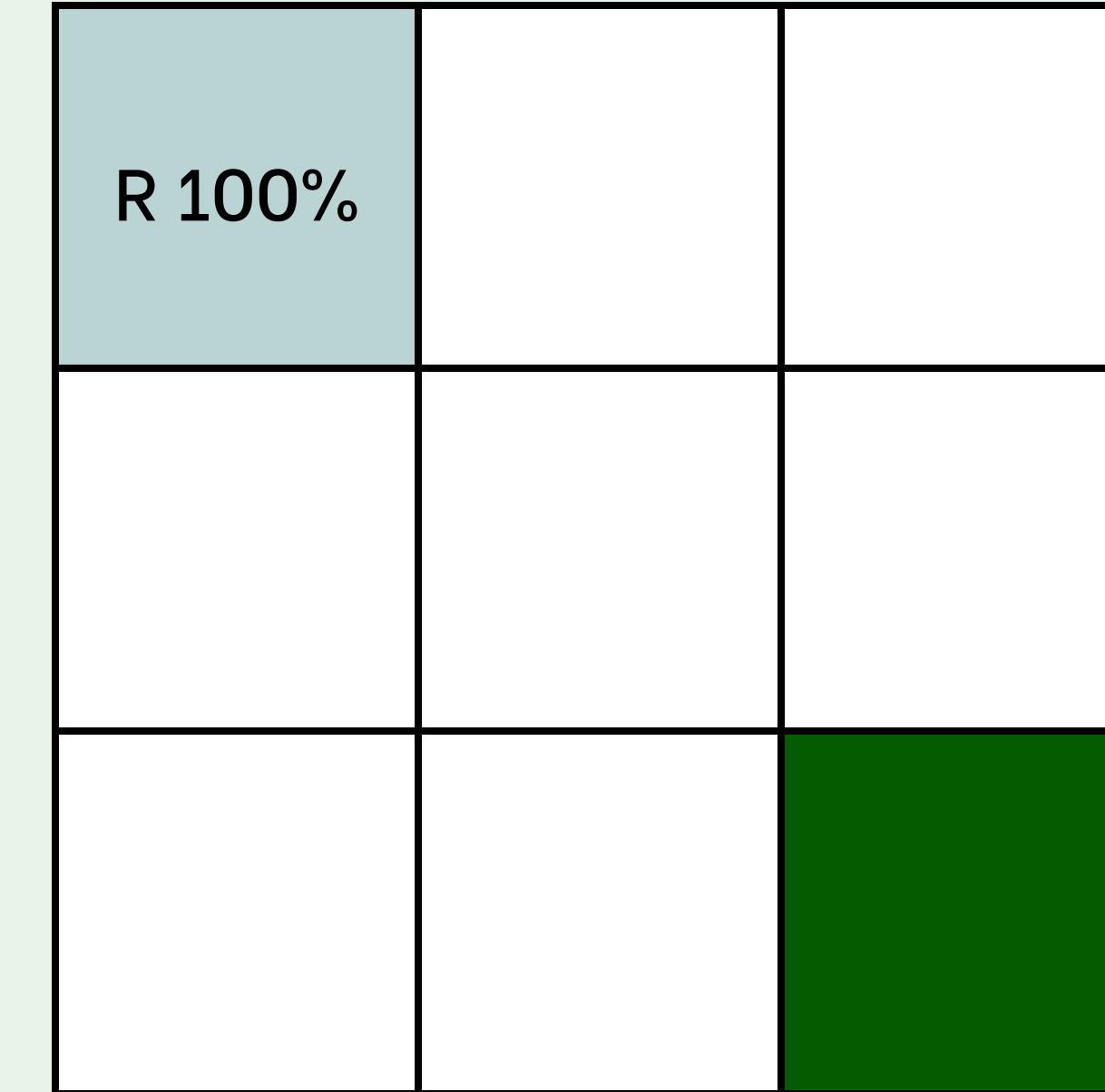
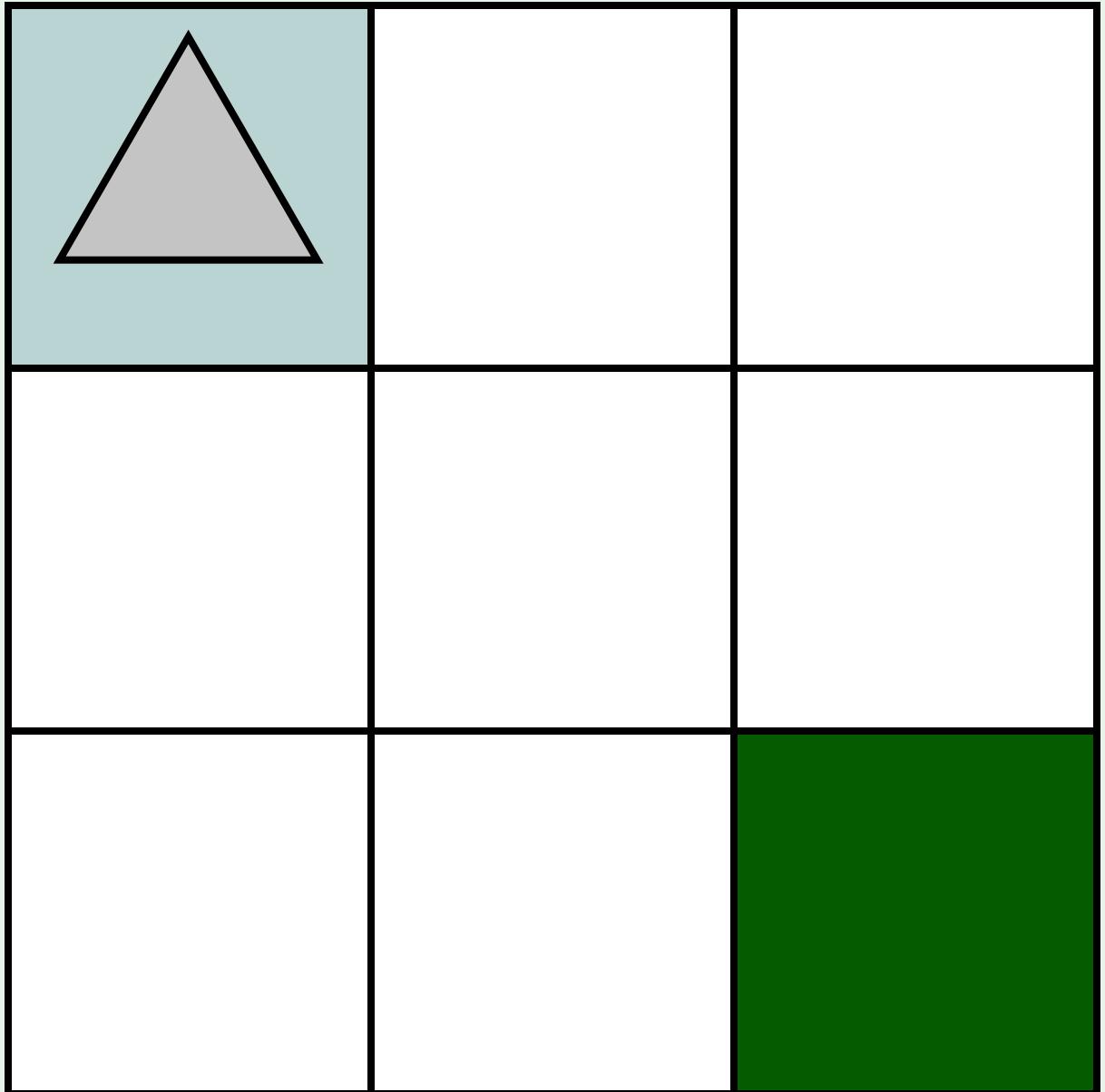
- OpenAI models
- GPT-4o was the best
- GPT-4o-mini selected for price/performance

## 2. Setting

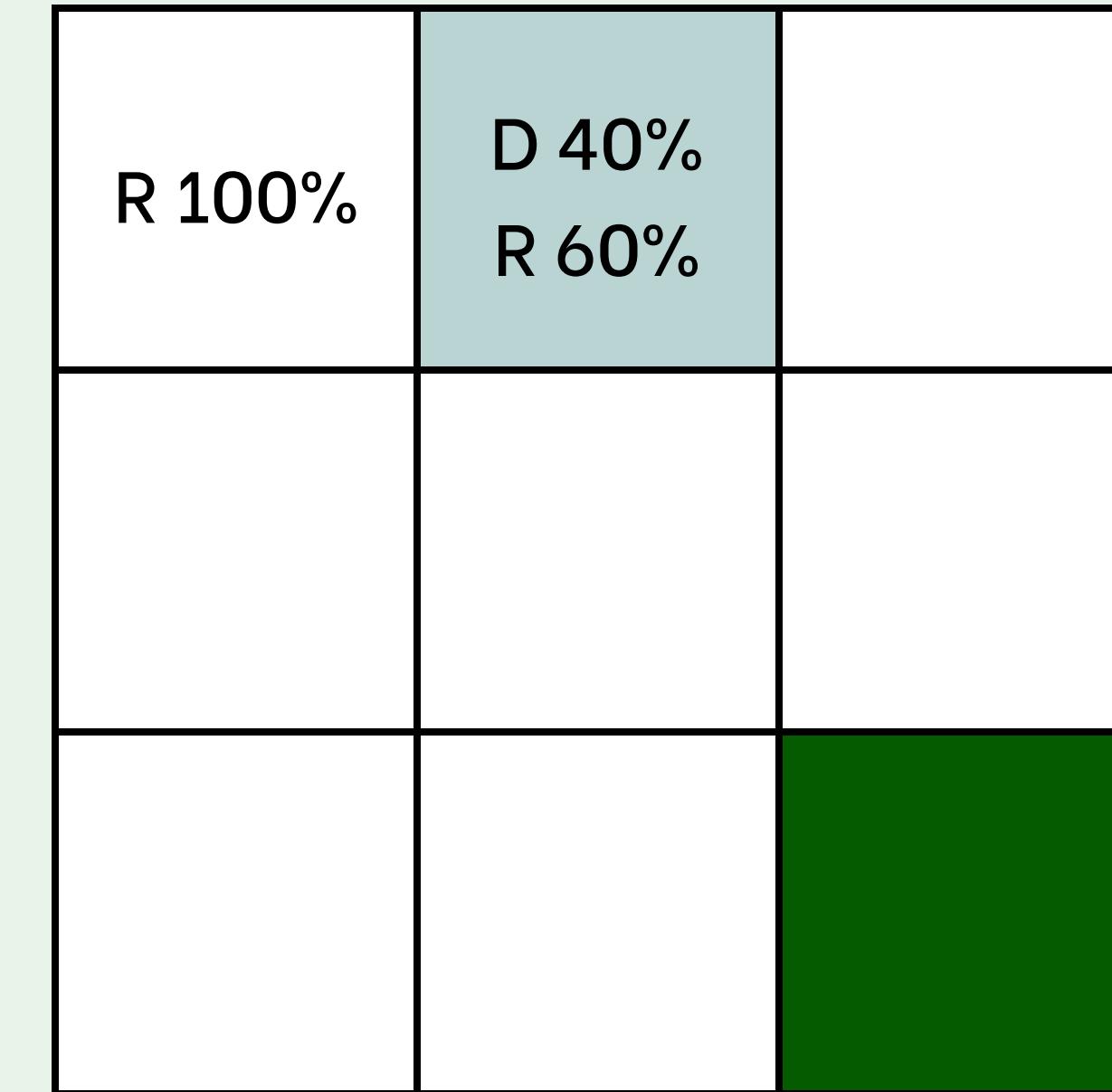
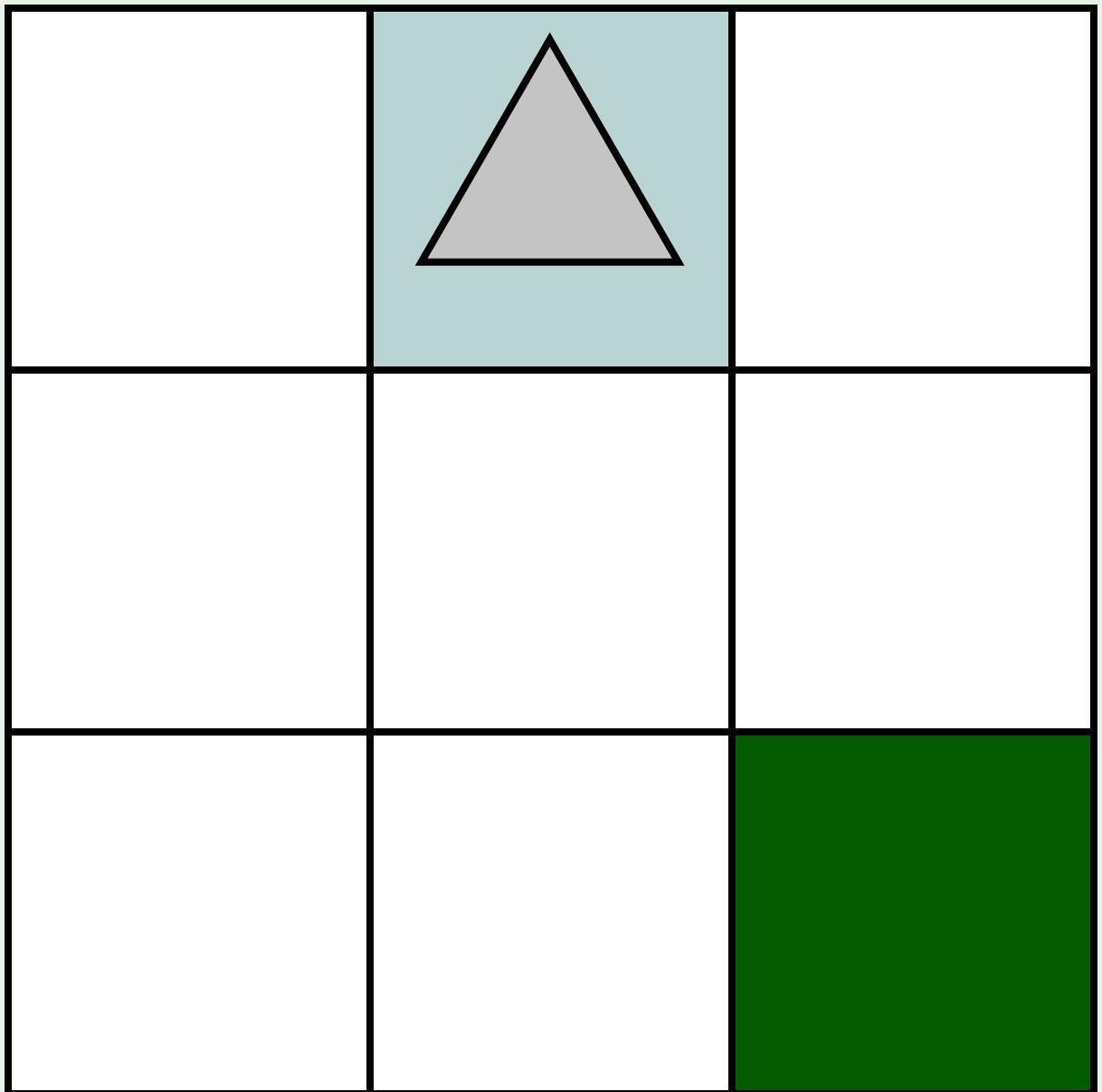


# Data Collection

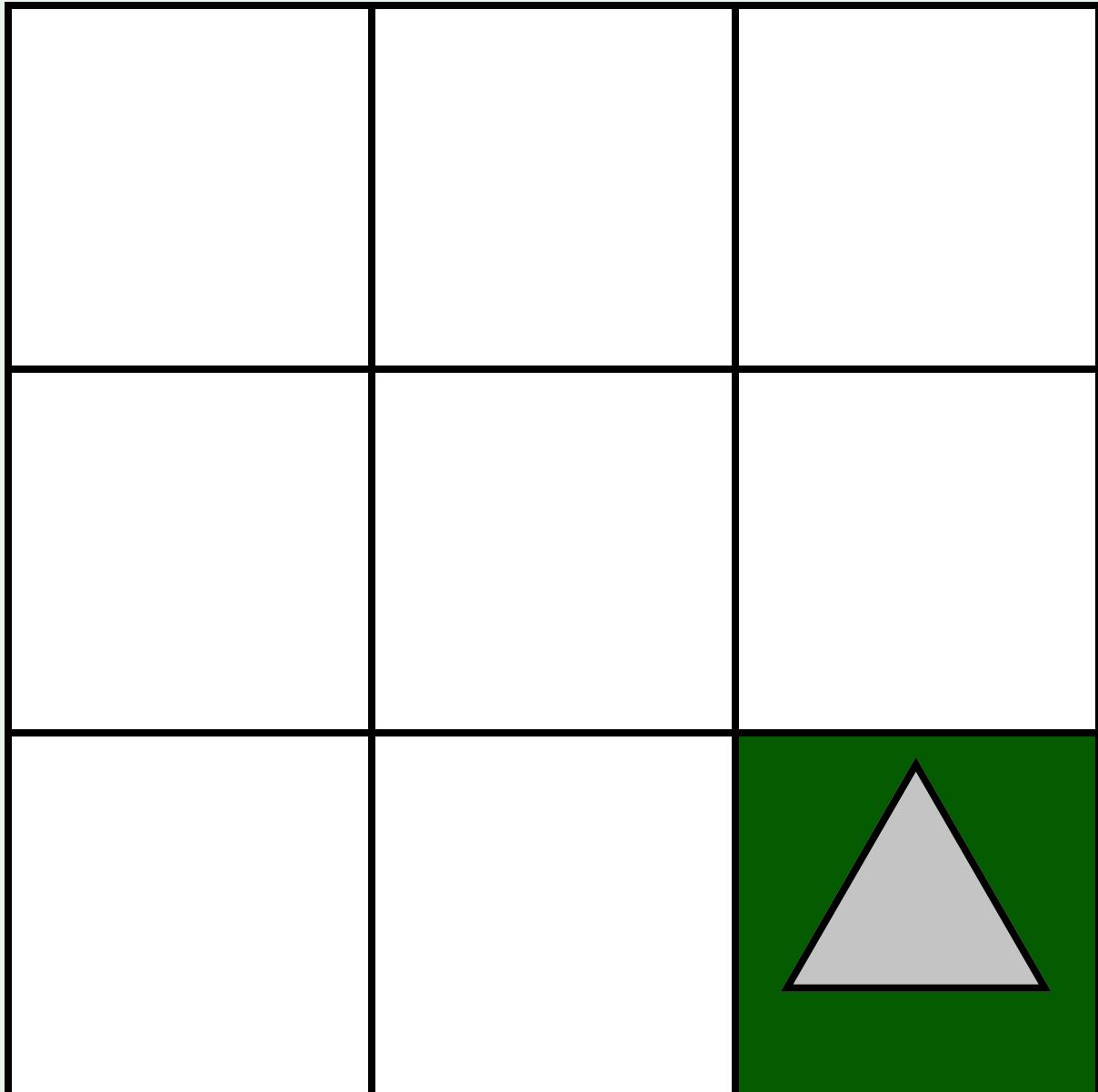
# Heatmaps Creation #1



# Heatmaps Creation #1

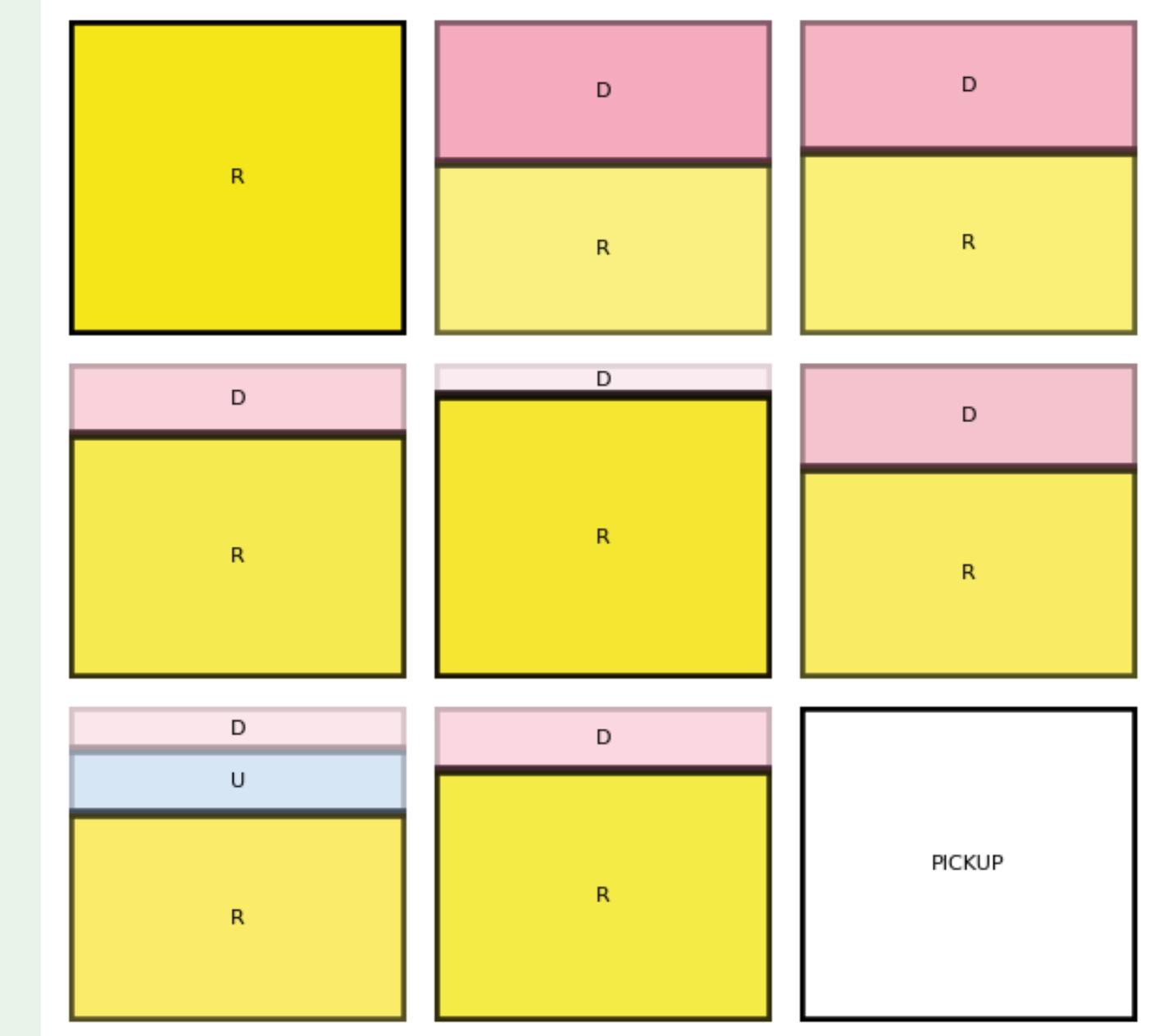
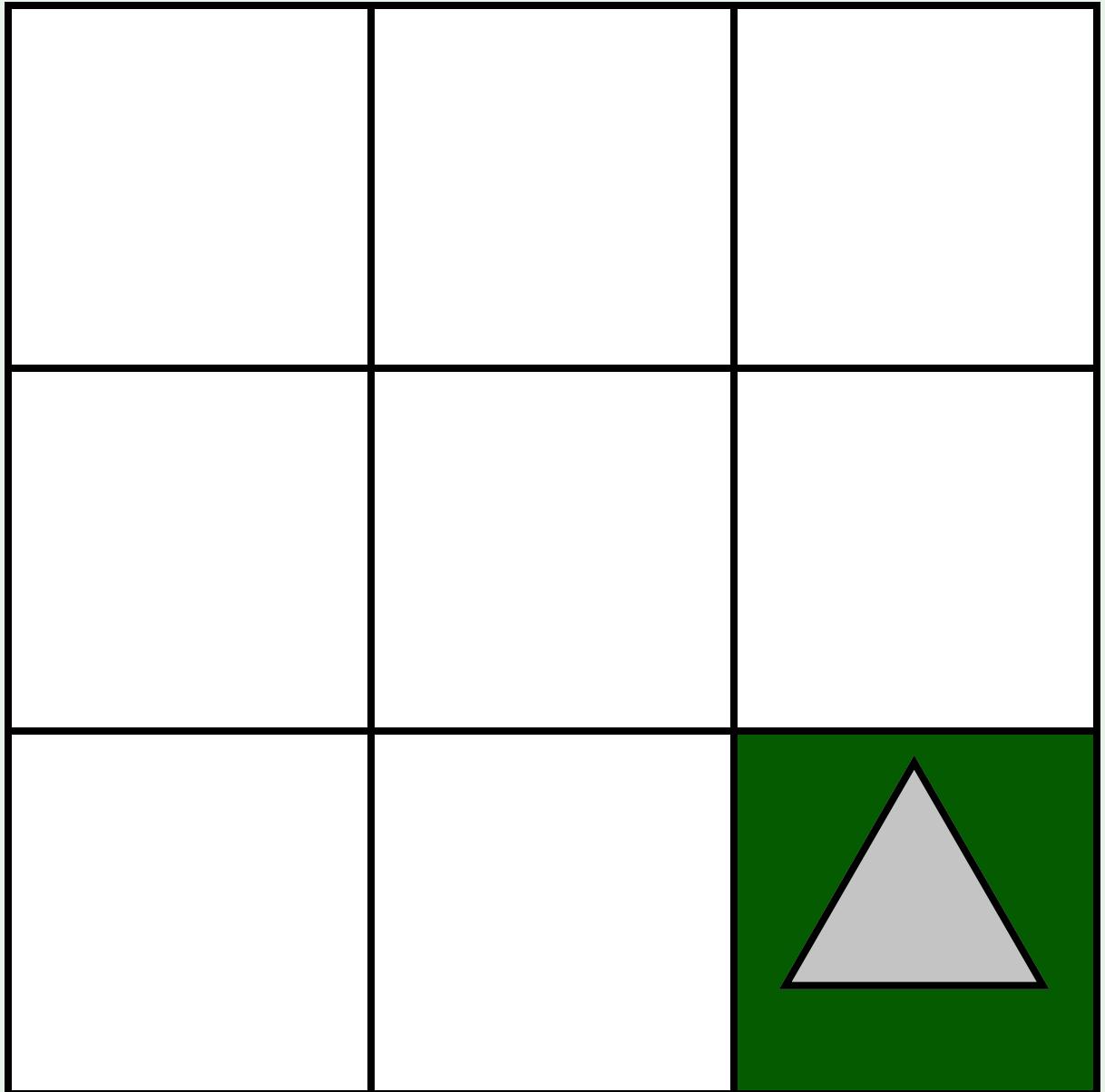


# Heatmaps Creation #1



R 100%	D 40% R 60%	D 41% R 59%
D 5% R 95%	D 5% R 95%	D 33% R 67%
D 16% U 18% R 66%	D 20% R 80%	

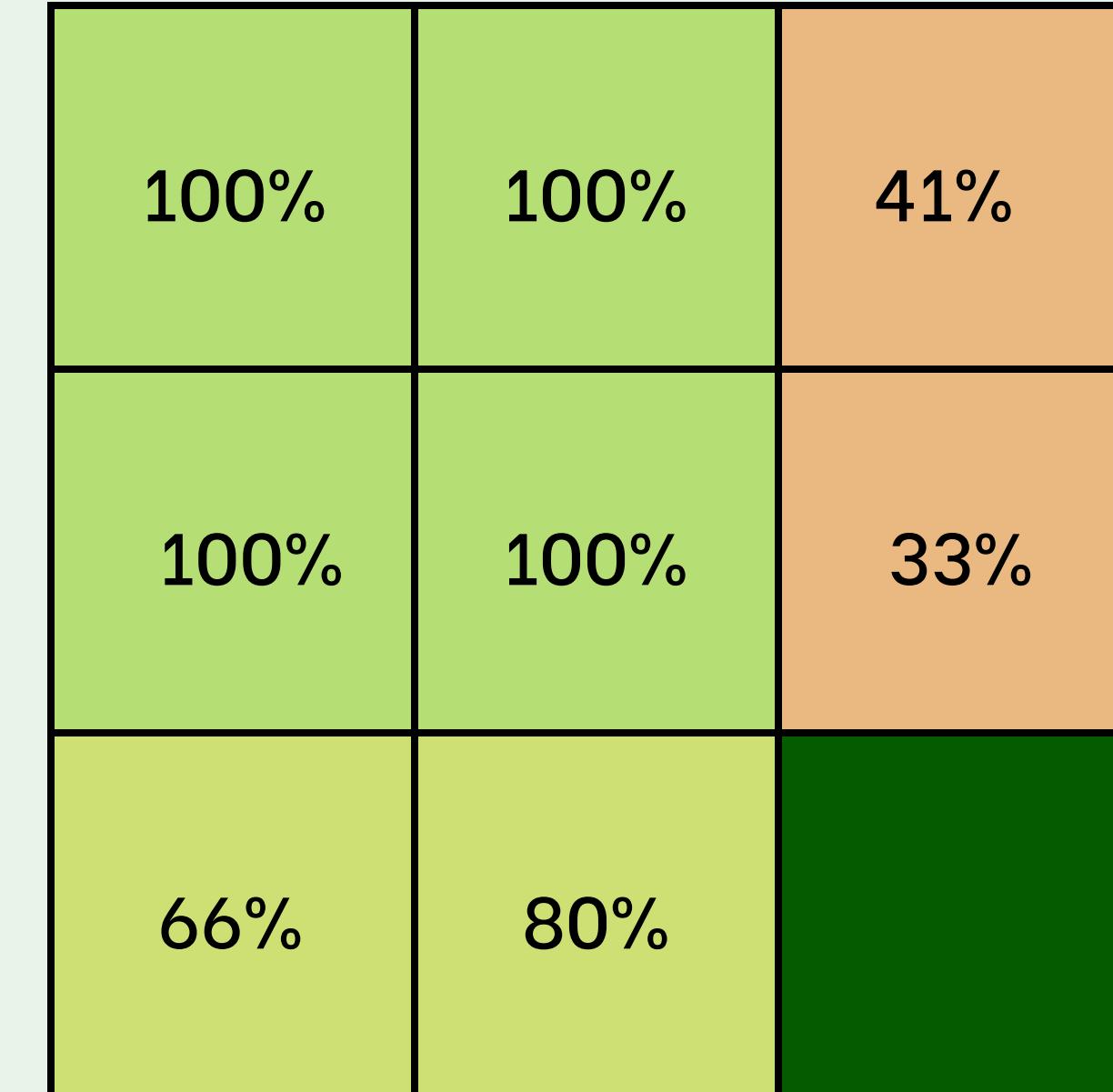
# Heatmaps Creation #1



Right Left Down Up

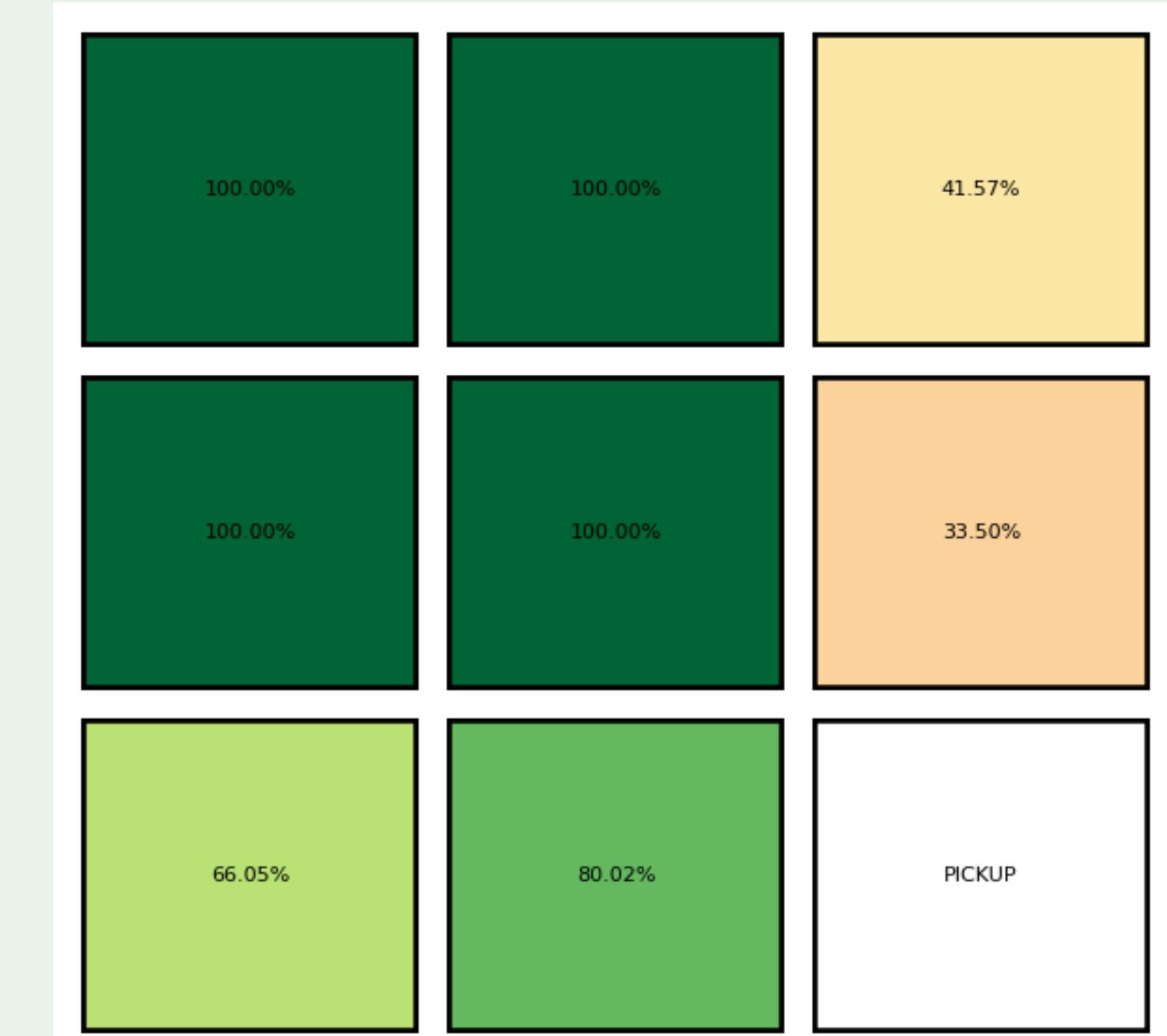
# Heatmaps Creation #2

R 100%	D 40% R 60%	D 41% R 59%
D 5% R 95%	D 5% R 95%	D 33% R 67%
D 16% U 18% R 66%	D 20% R 80%	



# Heatmaps Creation #2

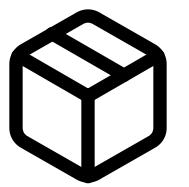
R 100%	D 40% R 60%	D 41% R 59%
D 5% R 95%	D 5% R 95%	D 33% R 67%
D 16% U 18% R 66%	D 20% R 80%	



100%  
0%

# Testing Strategy

## Goals



### Pickup

Goal tile explicit



### Deliver

Goal tile identified in  
map description

Average top1%  
87%

Average top1%  
81%

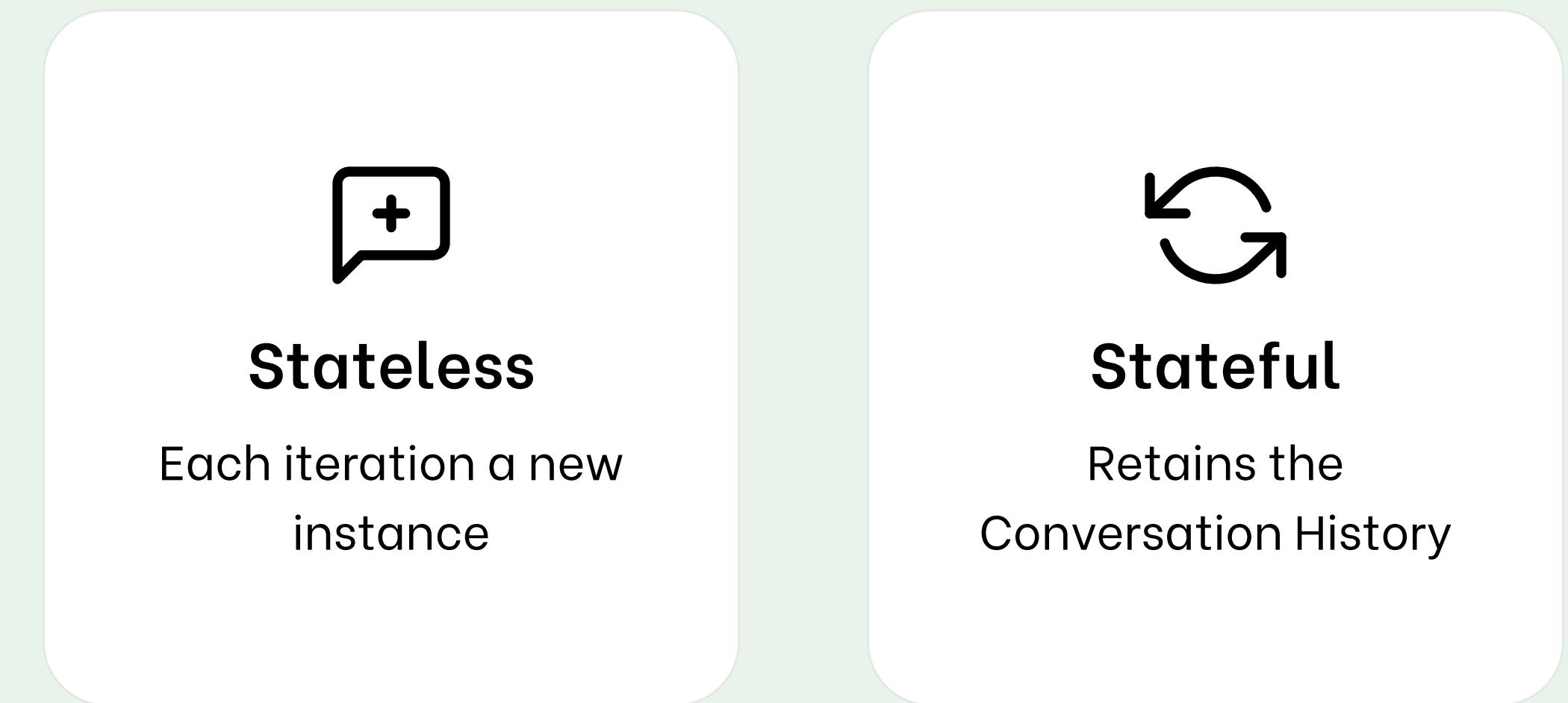
## 3. Data Collection



UNIVERSITÀ  
DI TRENTO

# Testing Strategy

## Agents



## 3. Data Collection



UNIVERSITÀ  
DI TRENTO

# Our Findings

# Map Orientation

“Since we did not provide any info about the orientation,  
how does the LLM perceive it?”

## 4. Findings

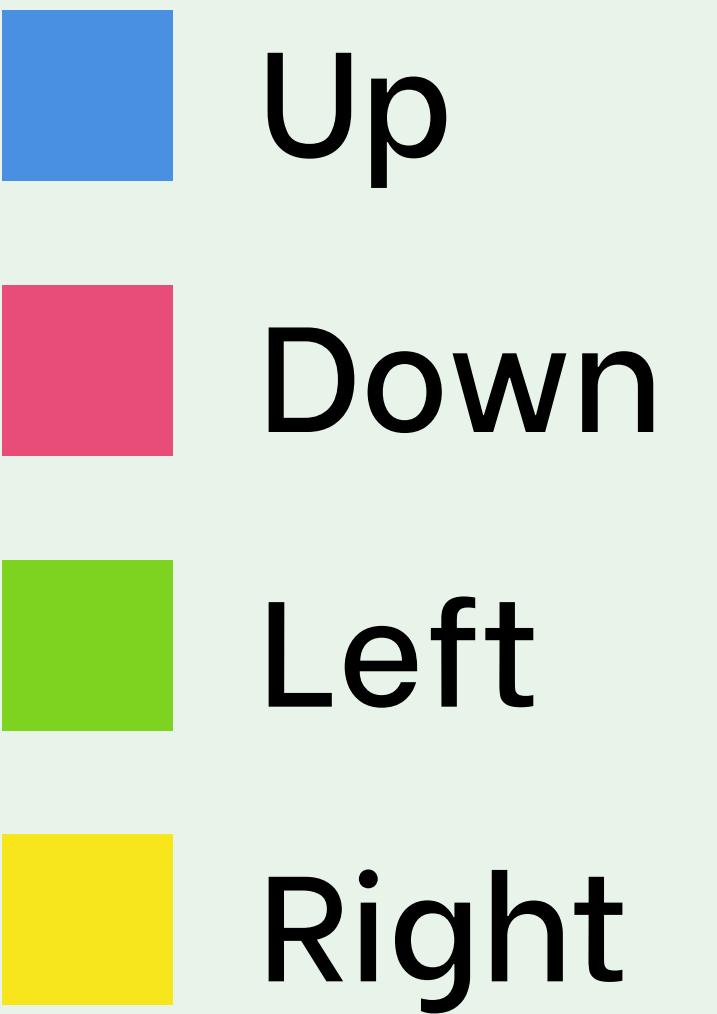
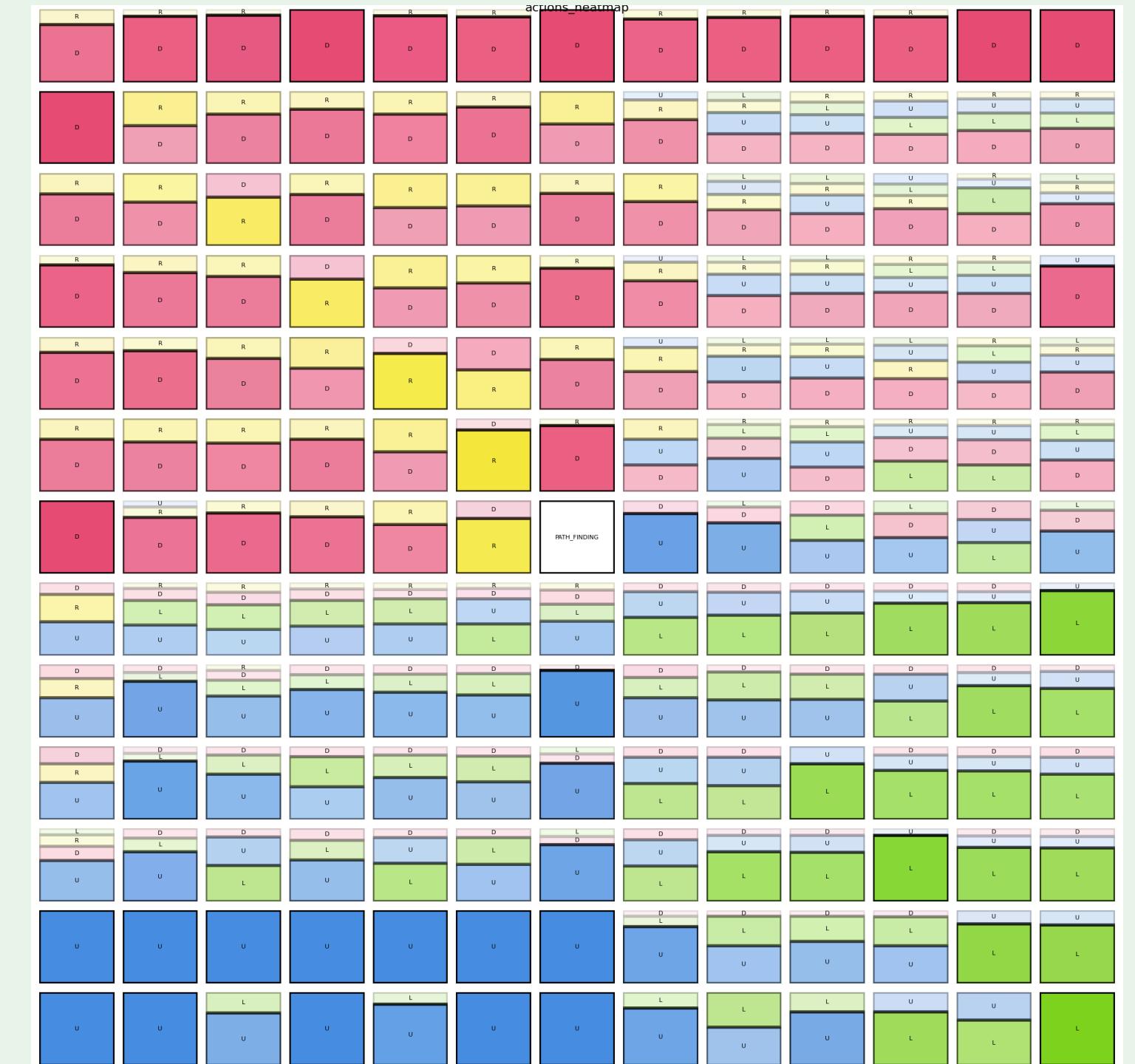


# Map Orientation

(0,0)



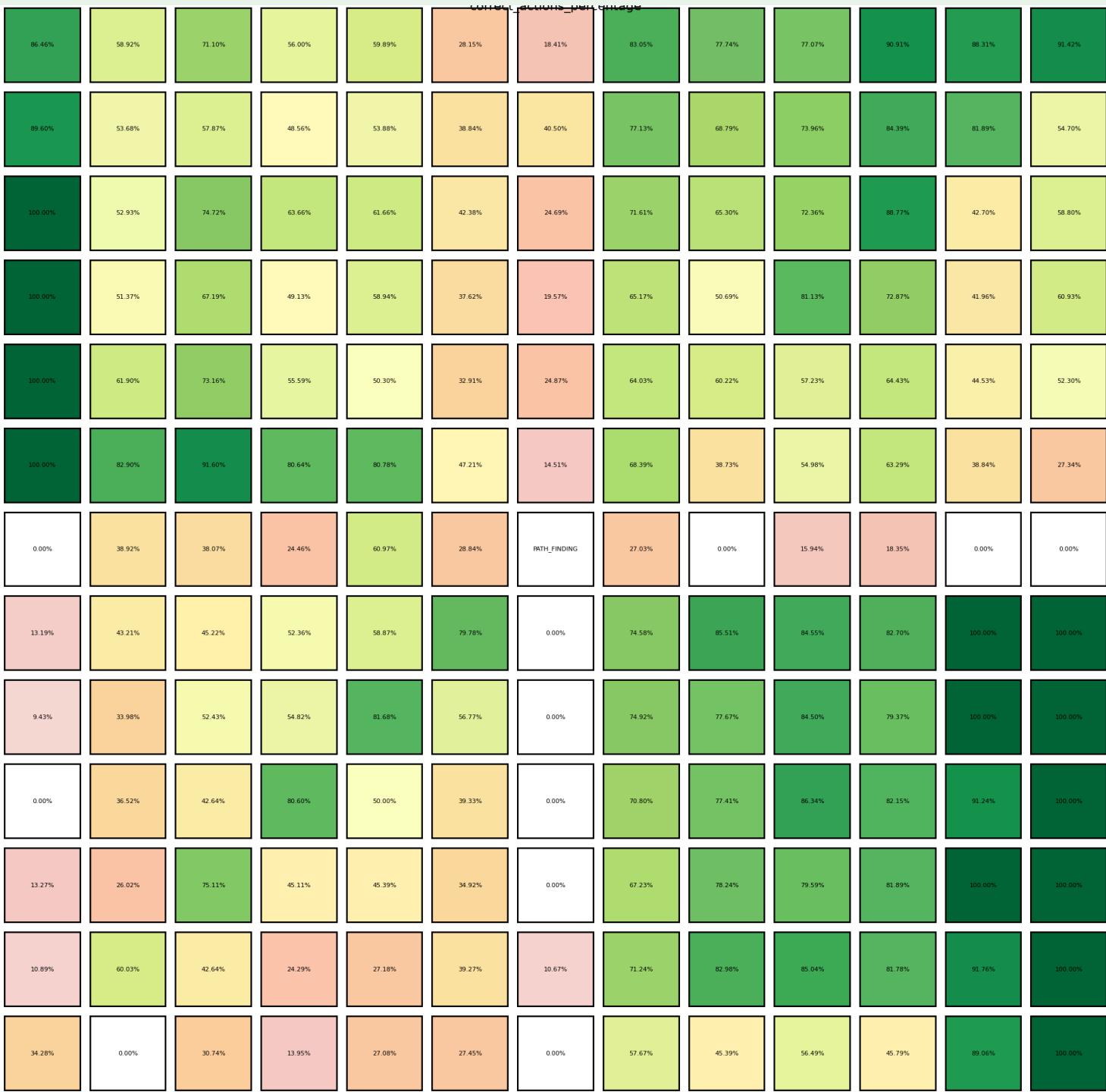
(0,0)



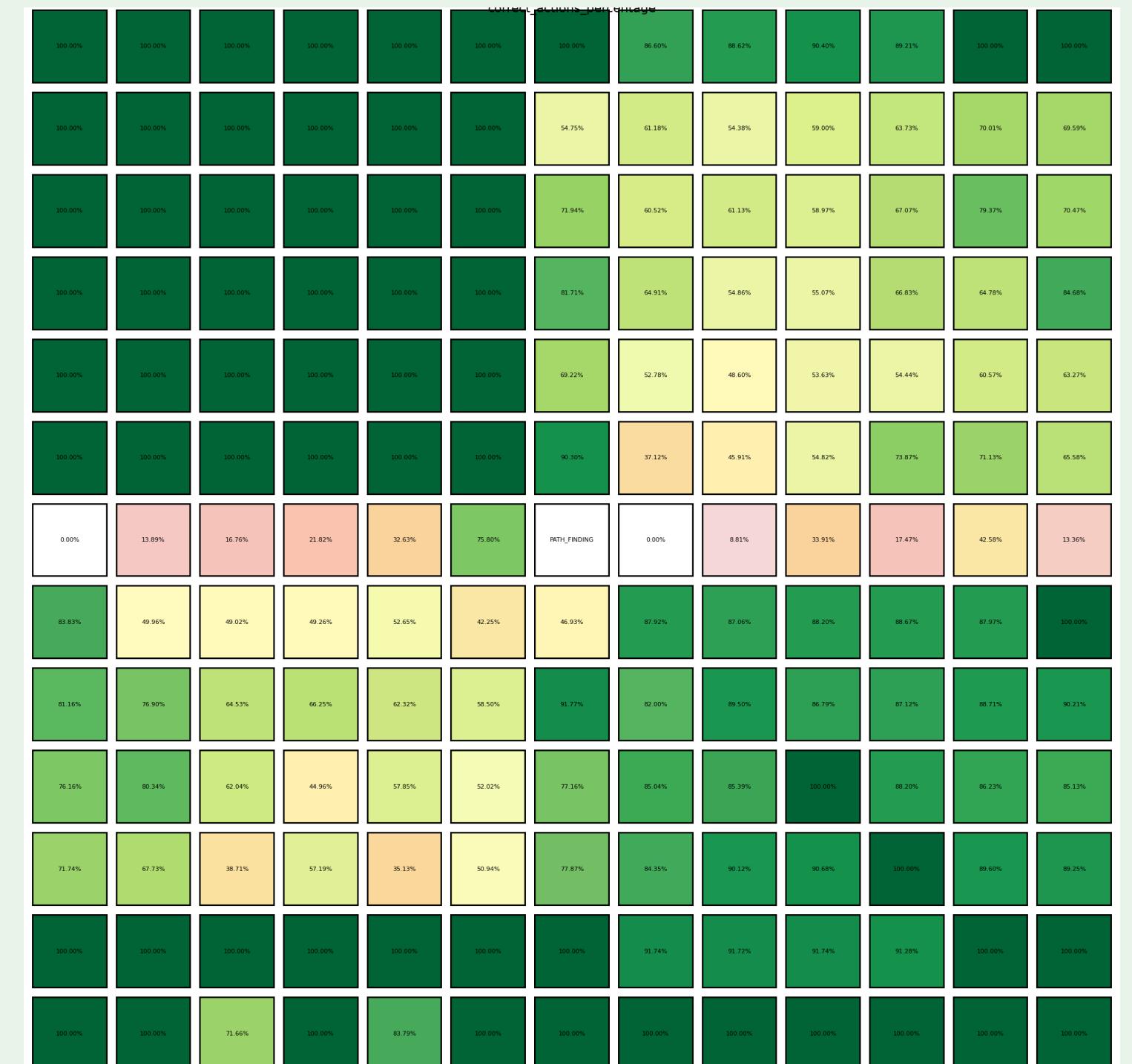
## 4. Findings

# Map Orientation

(0,0)



(0,0)



100%



0%

## 4. Findings

# Map Orientation

	Bottom-Left Origin	Top-Left Origin
top1%	62%	92%
top2%	92%	97%
top3%	93%	99%

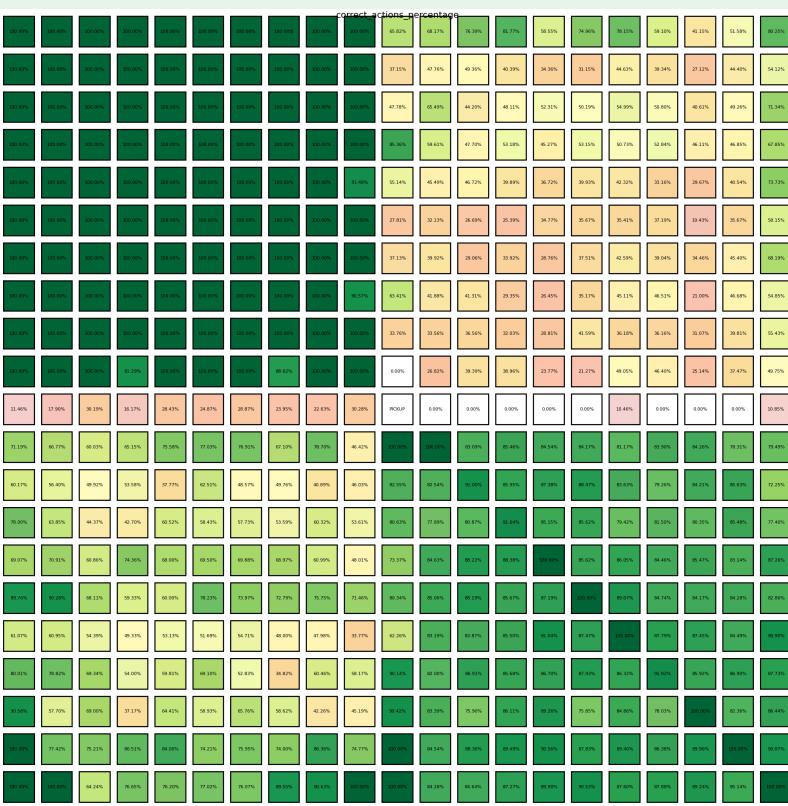
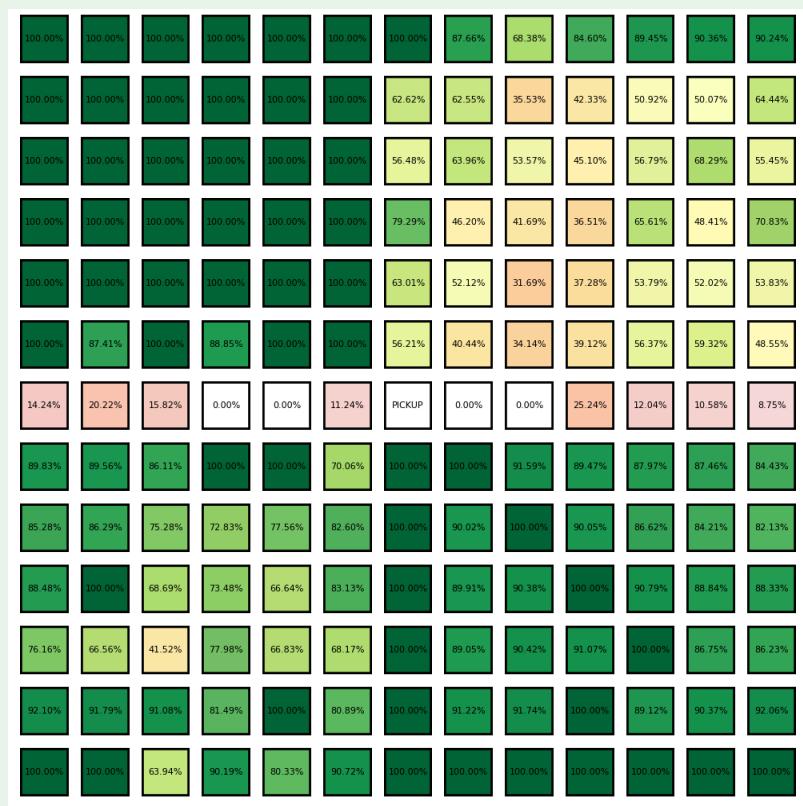
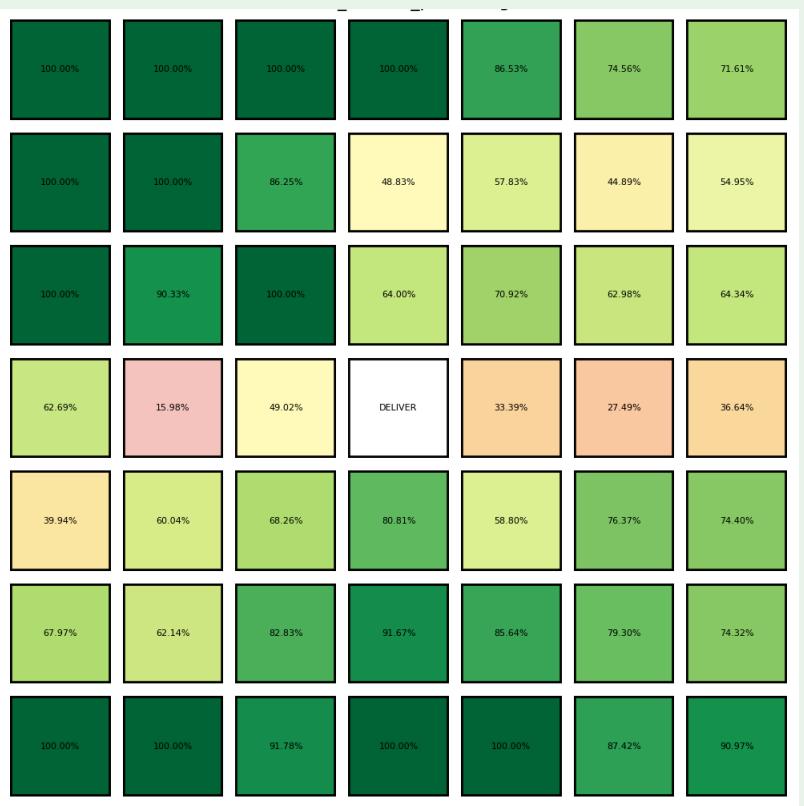
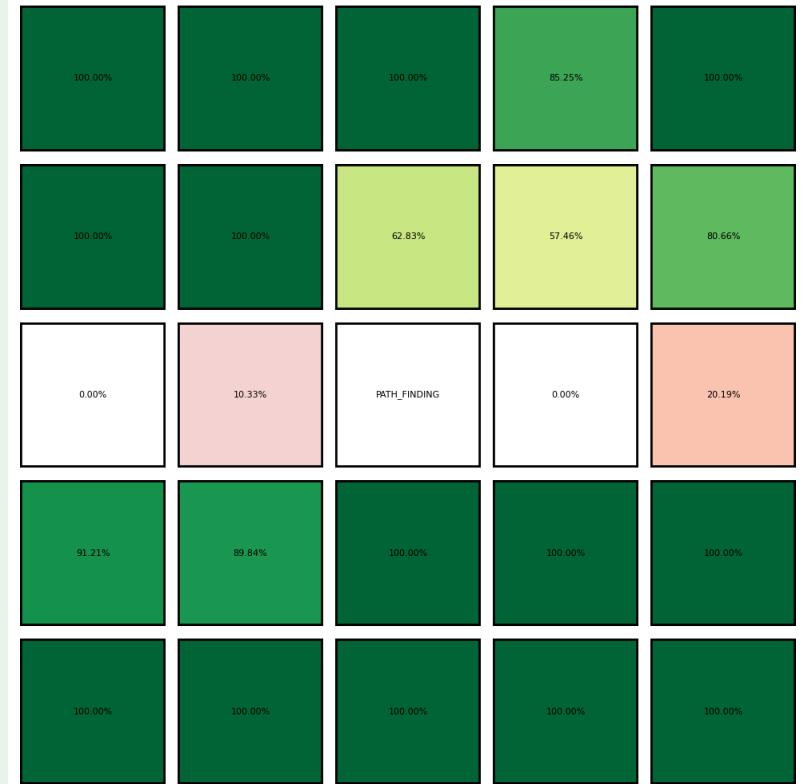
## 4. Findings



# Our Findings

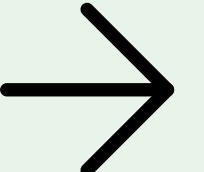
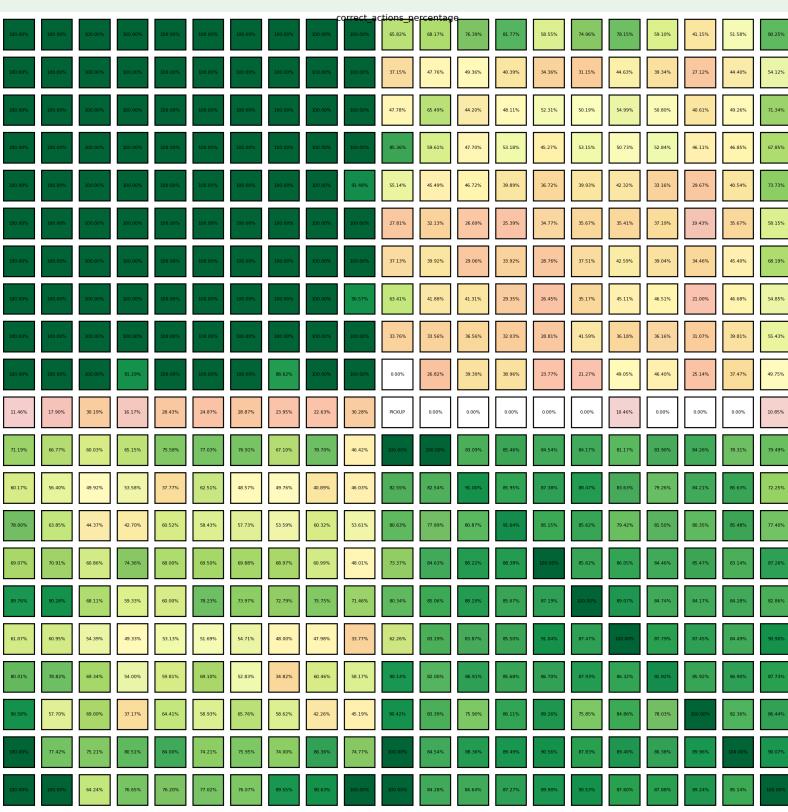
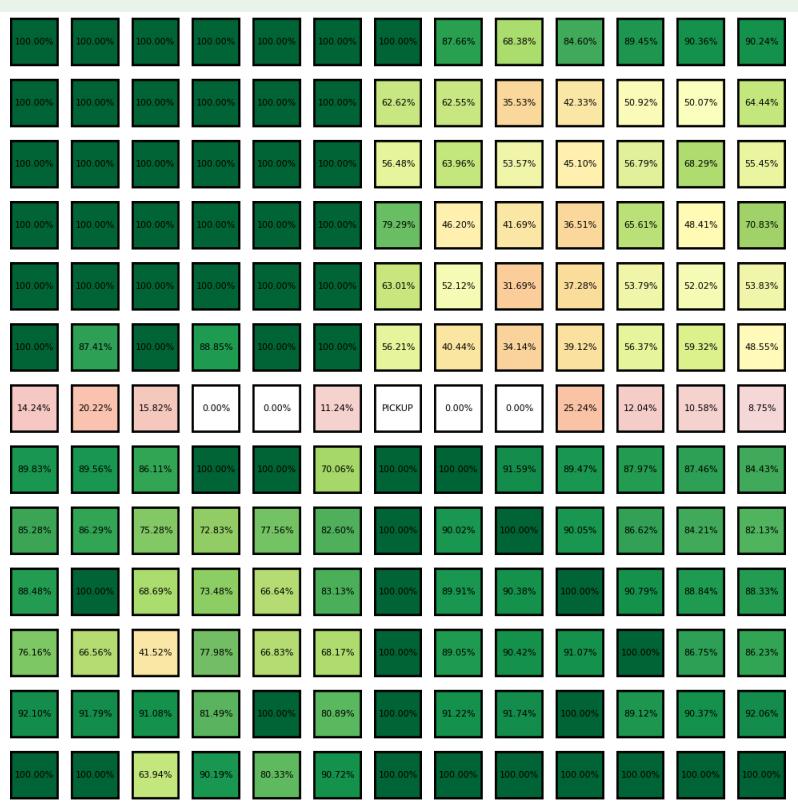
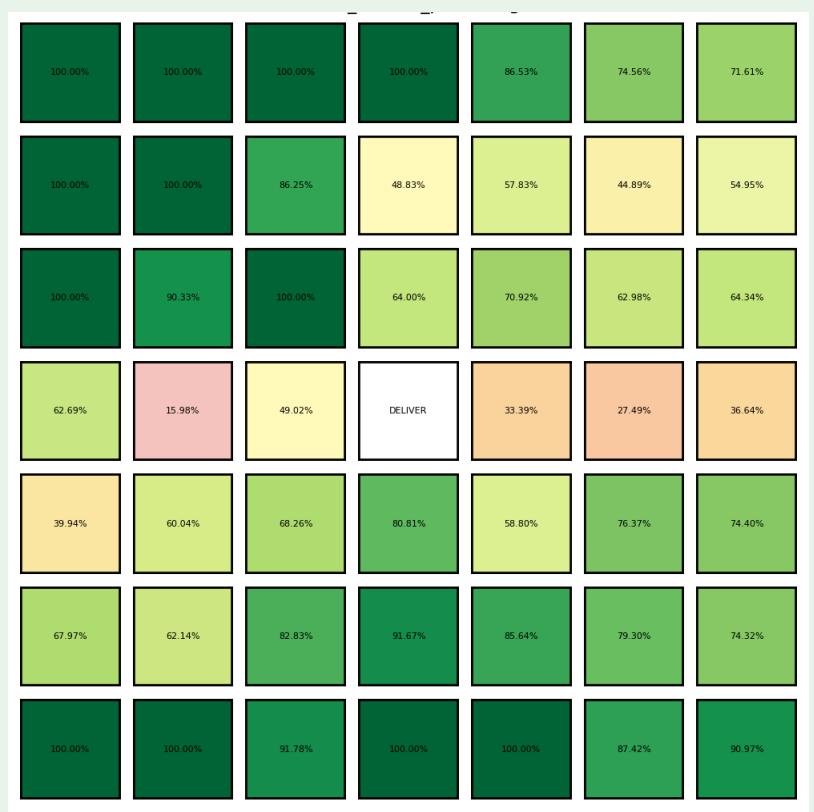
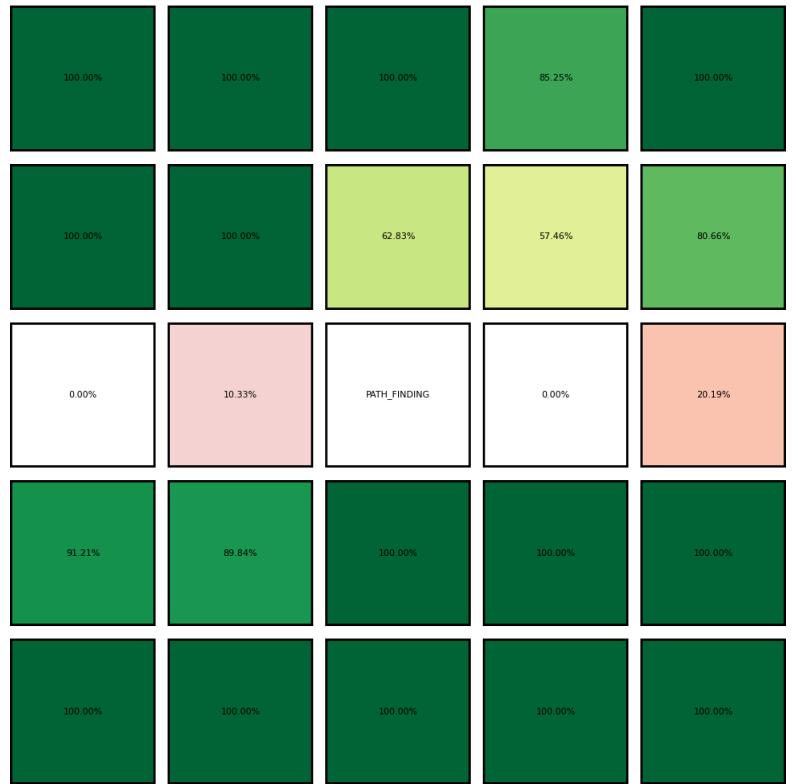
# Common Uncertainty Patterns

# Common Uncertainty Patterns #1

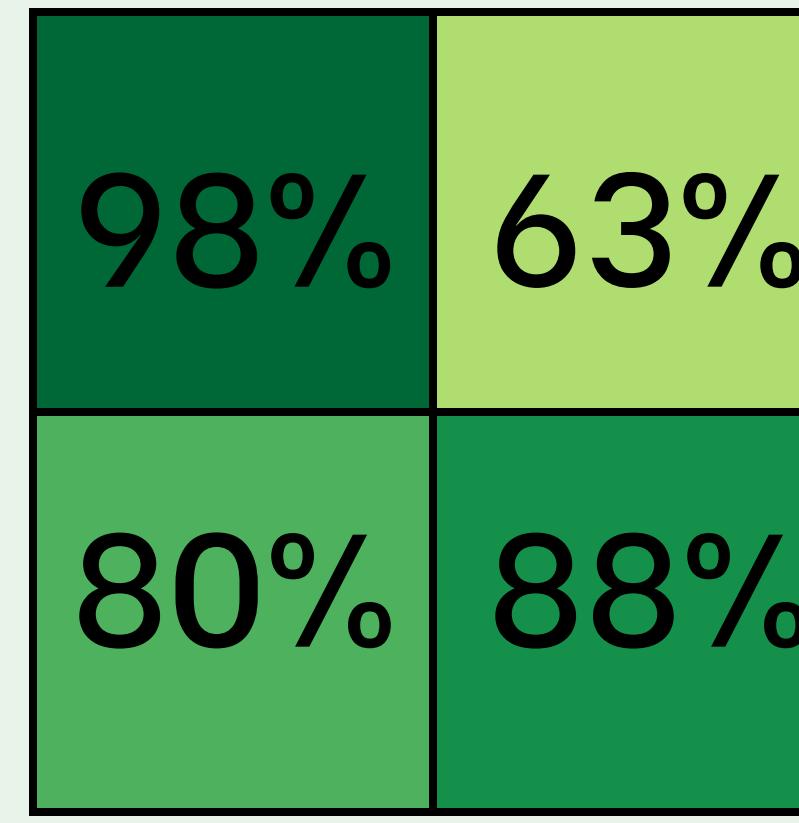


## 4. Findings

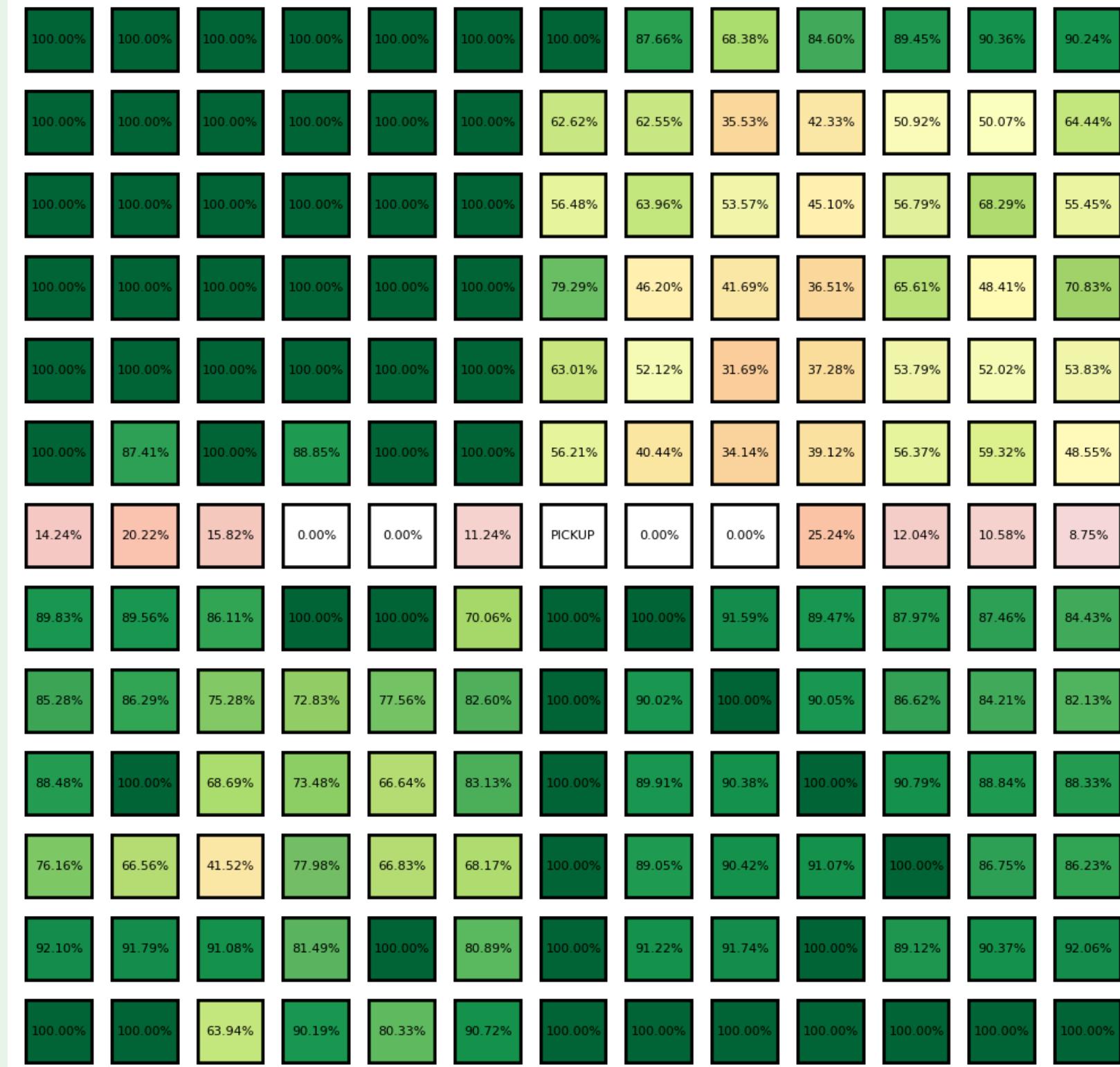
# Common Uncertainty Patterns #1



Average

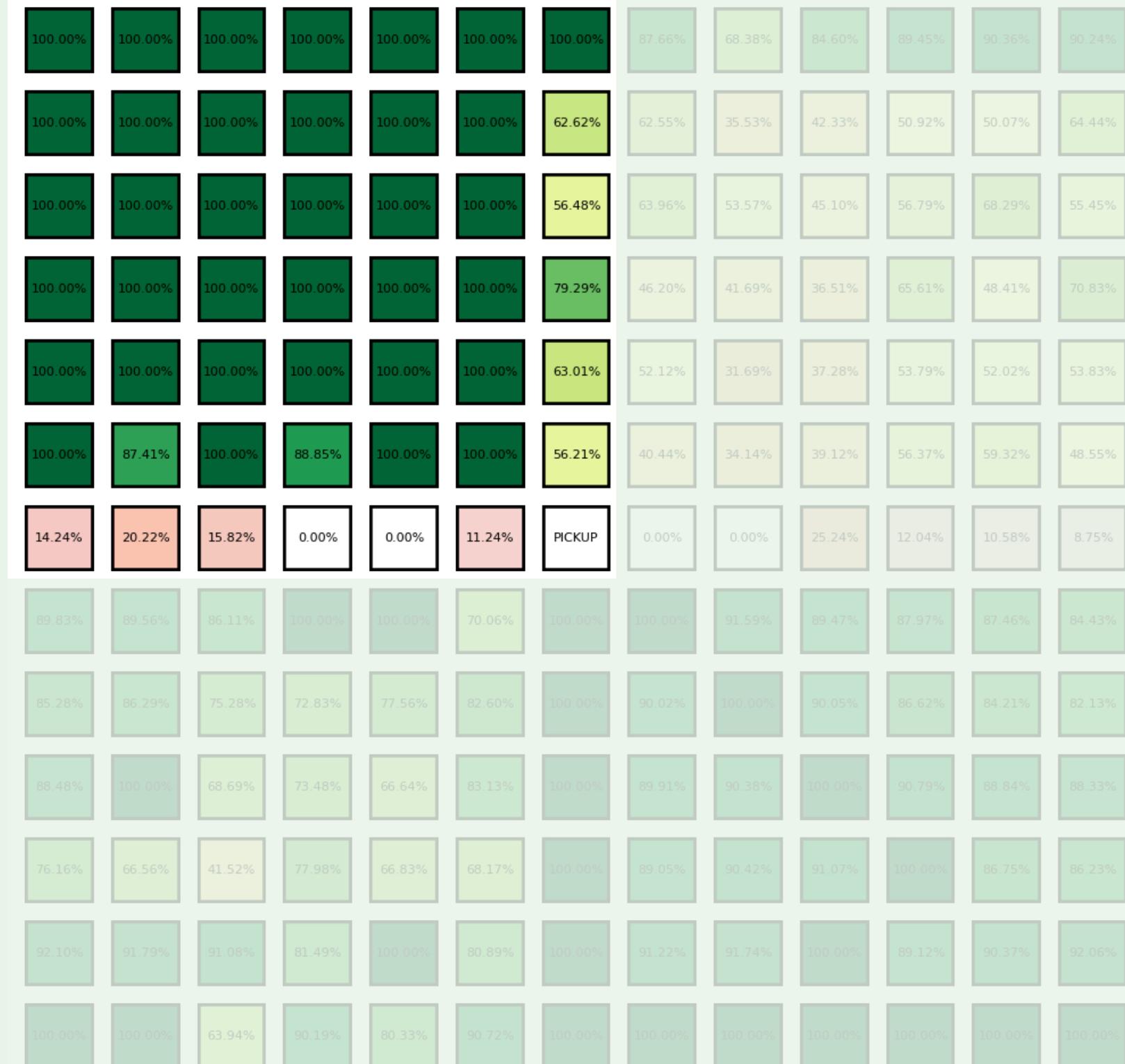


# Common Uncertainty Patterns #1.1



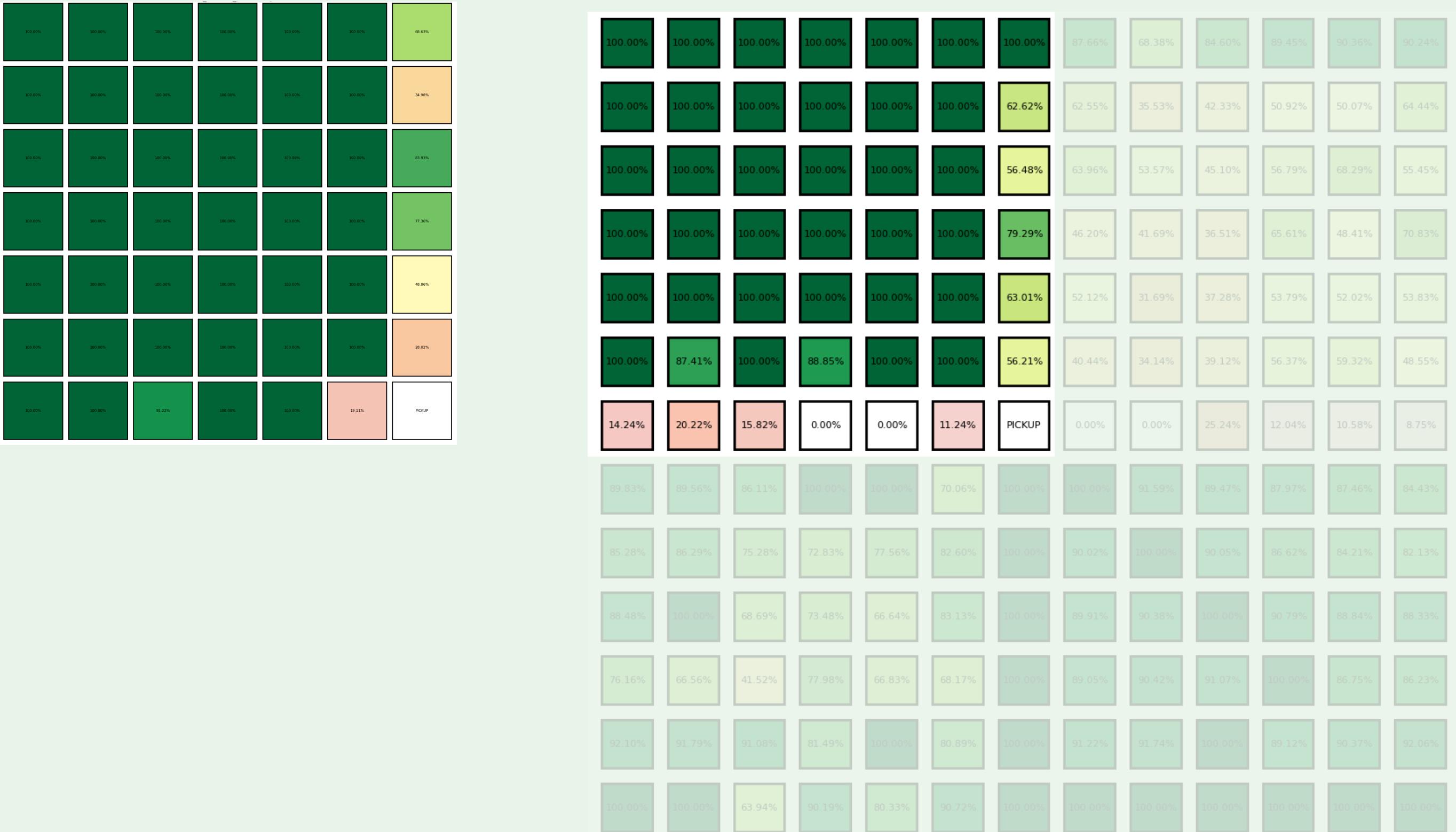
## 4. Findings

# Common Uncertainty Patterns #1.1



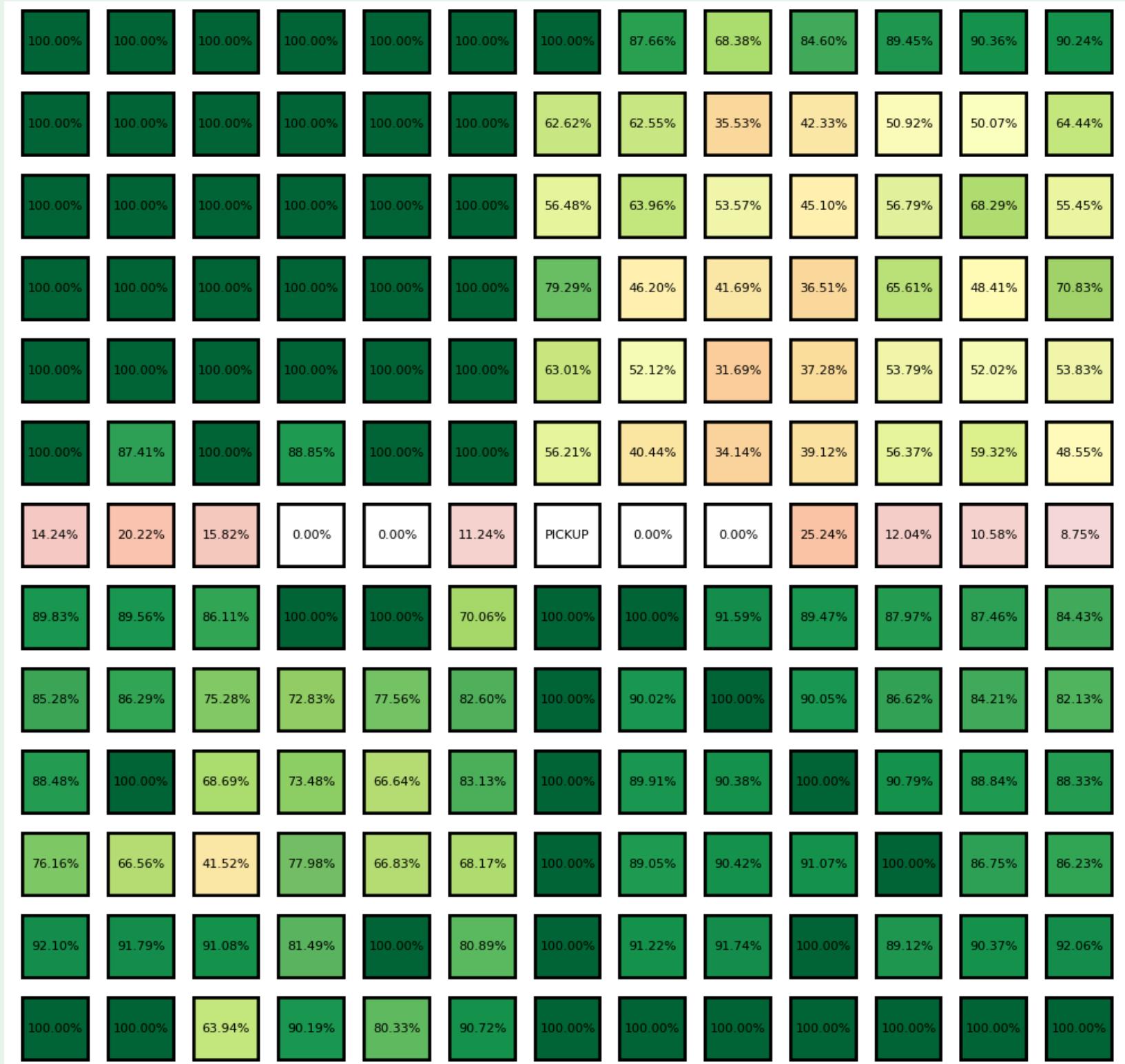
## 4. Findings

# Common Uncertainty Patterns #1.1



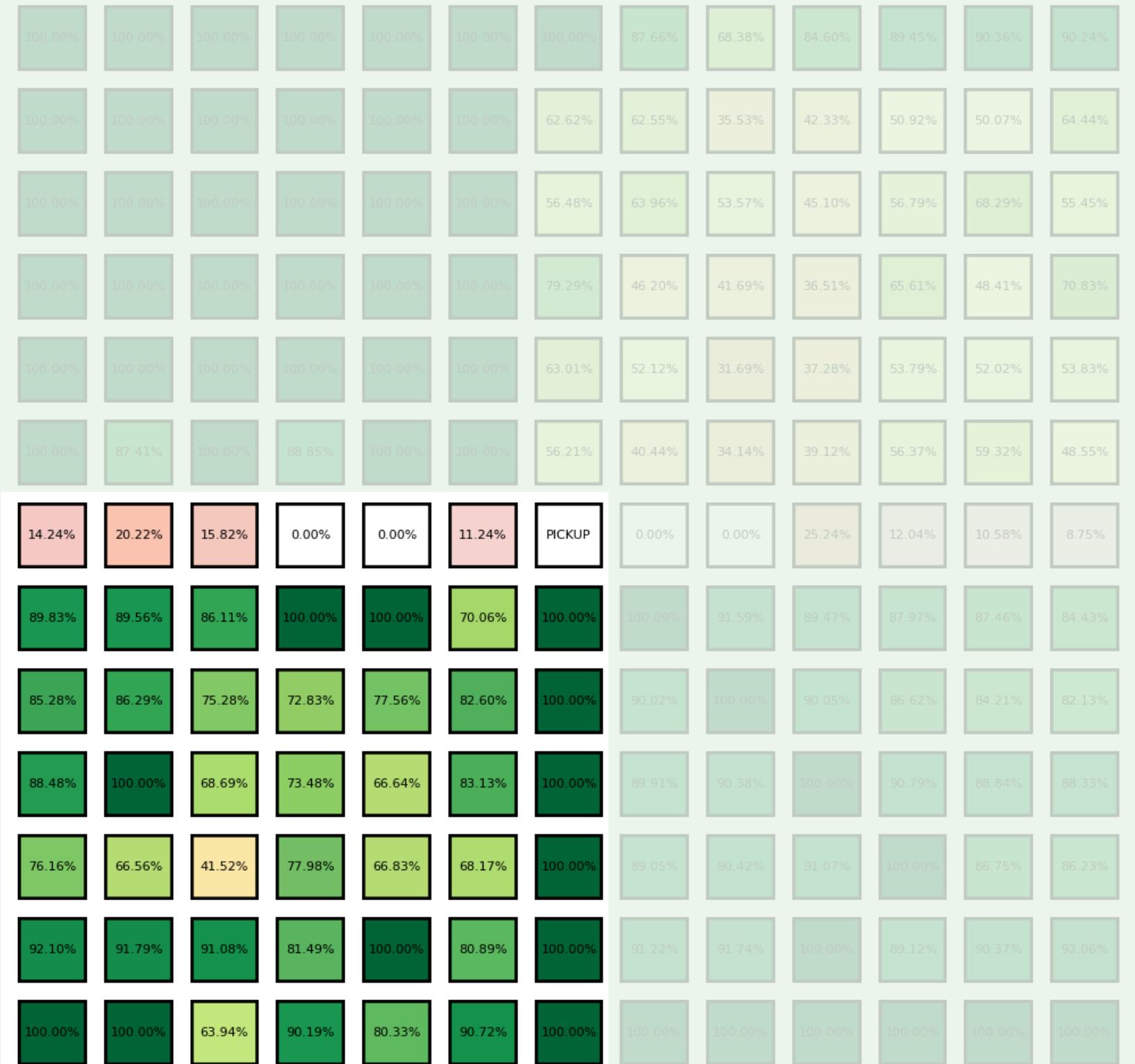
## 4. Findings

# Common Uncertainty Patterns #1.2



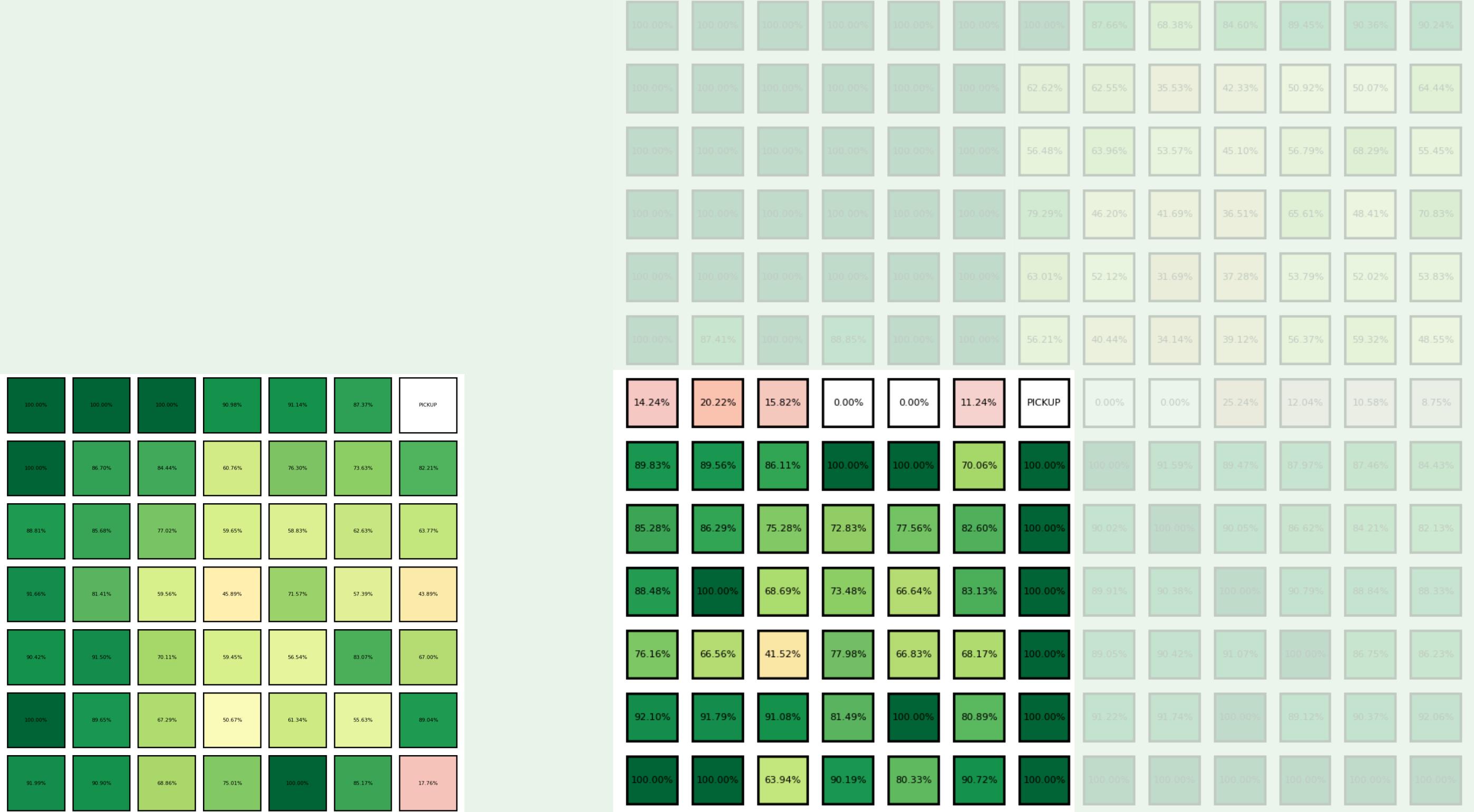
## 4. Findings

# Common Uncertainty Patterns #1.2



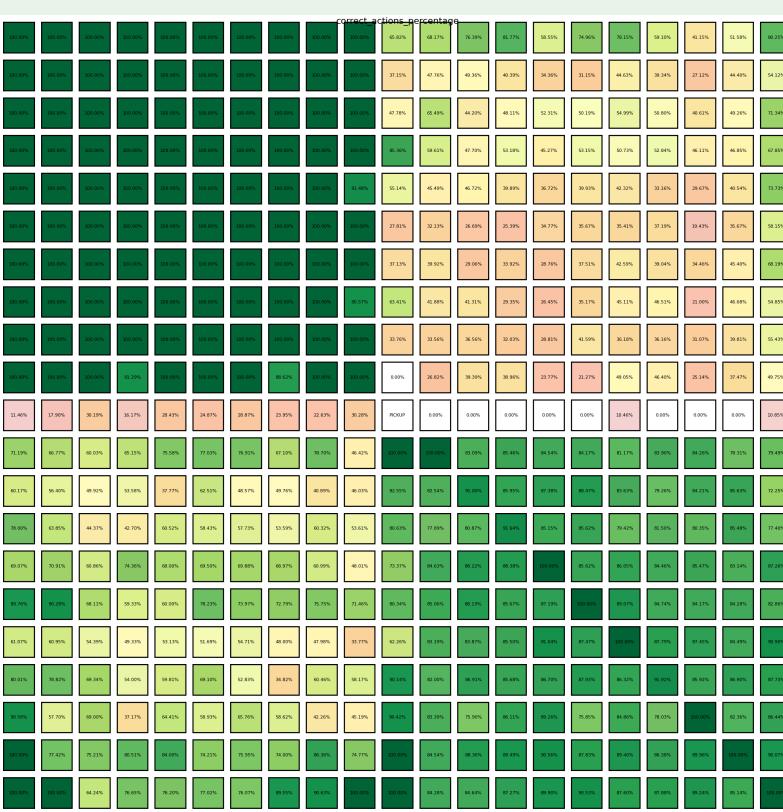
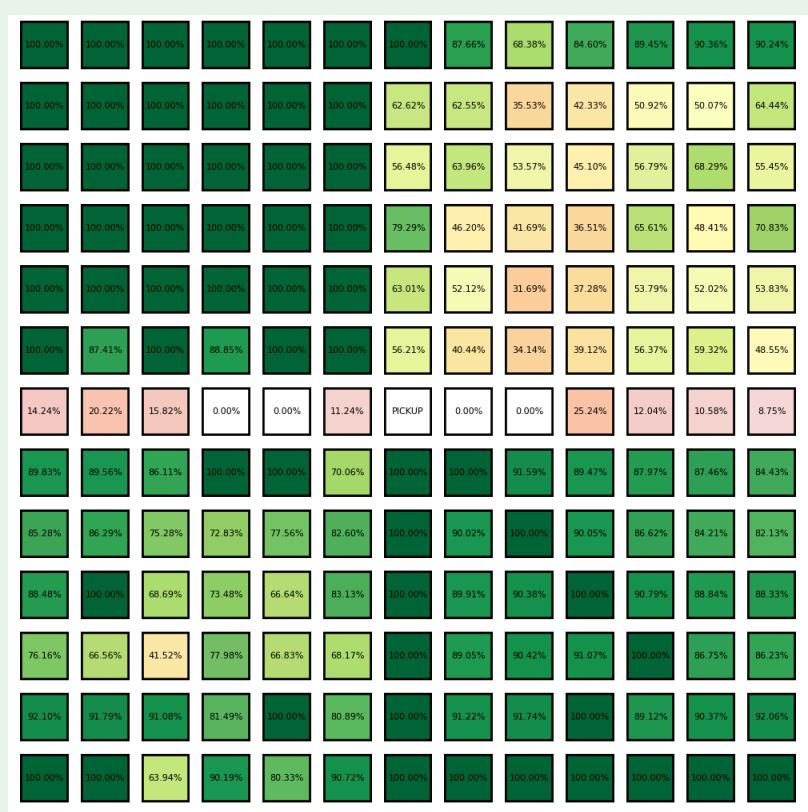
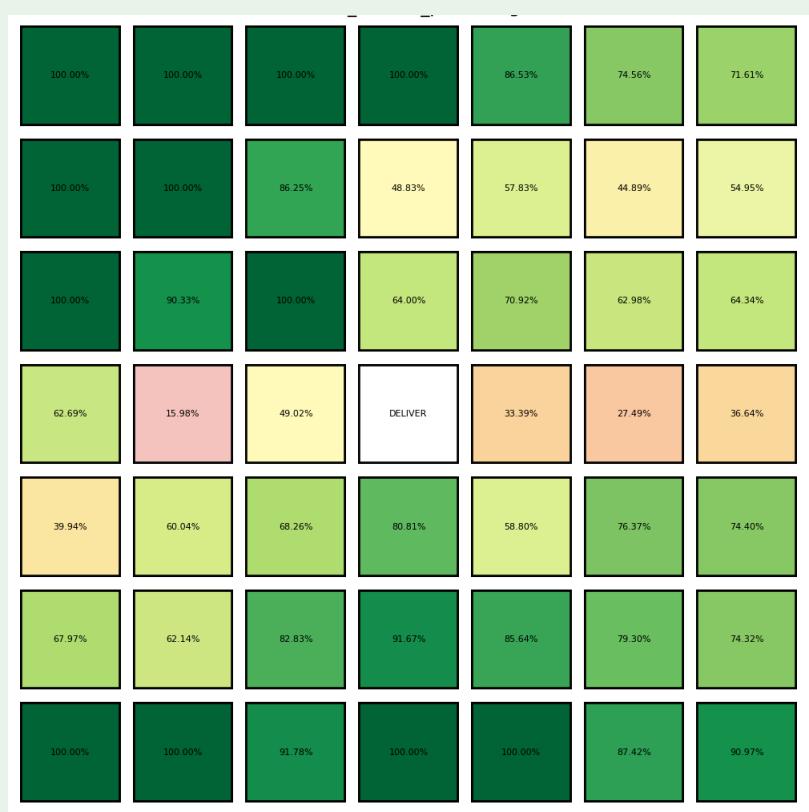
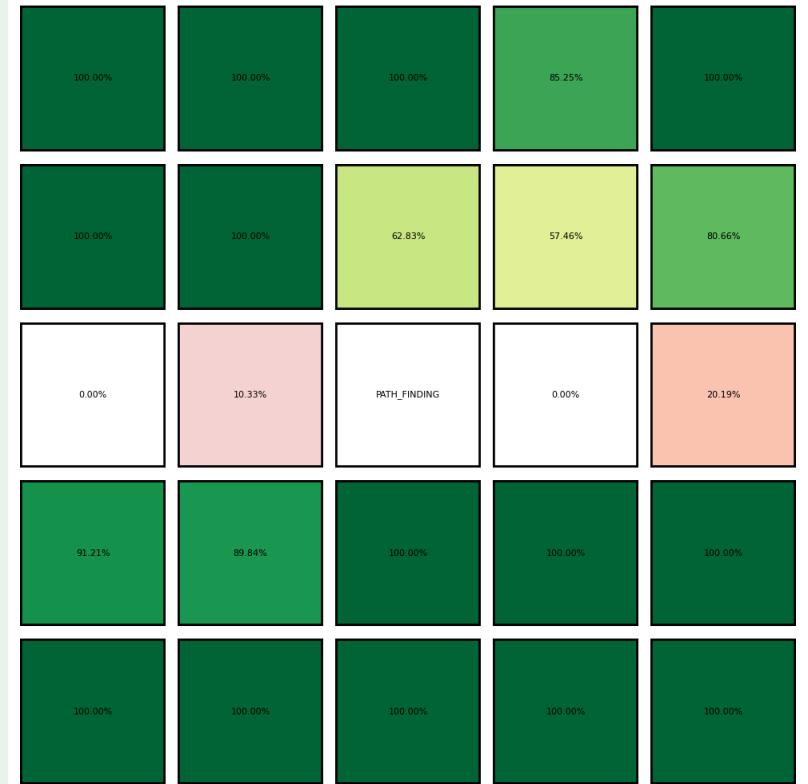
## 4. Findings

# Common Uncertainty Patterns #1.2

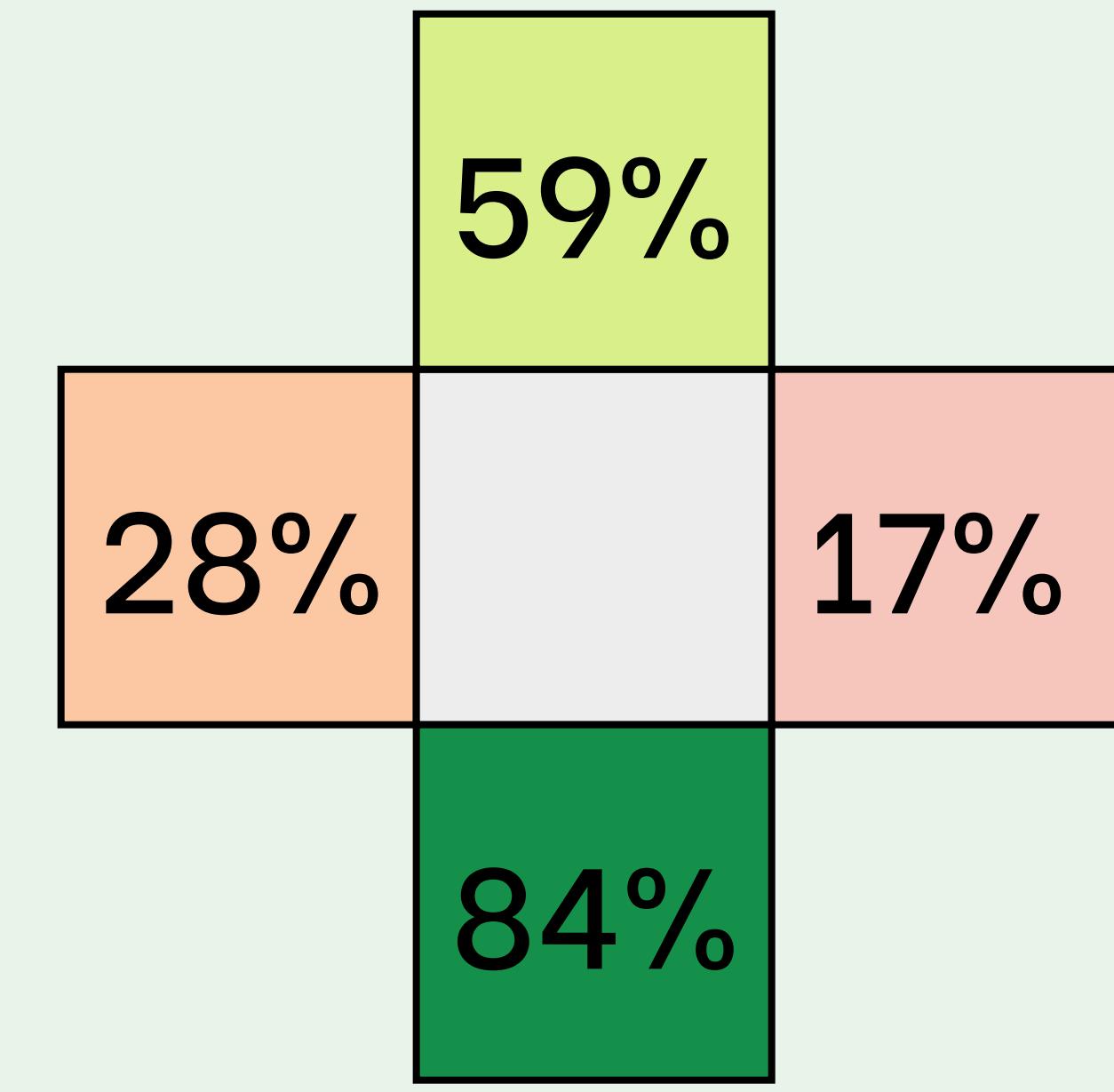


## 4. Findings

# Common Uncertainty Patterns #2



Average

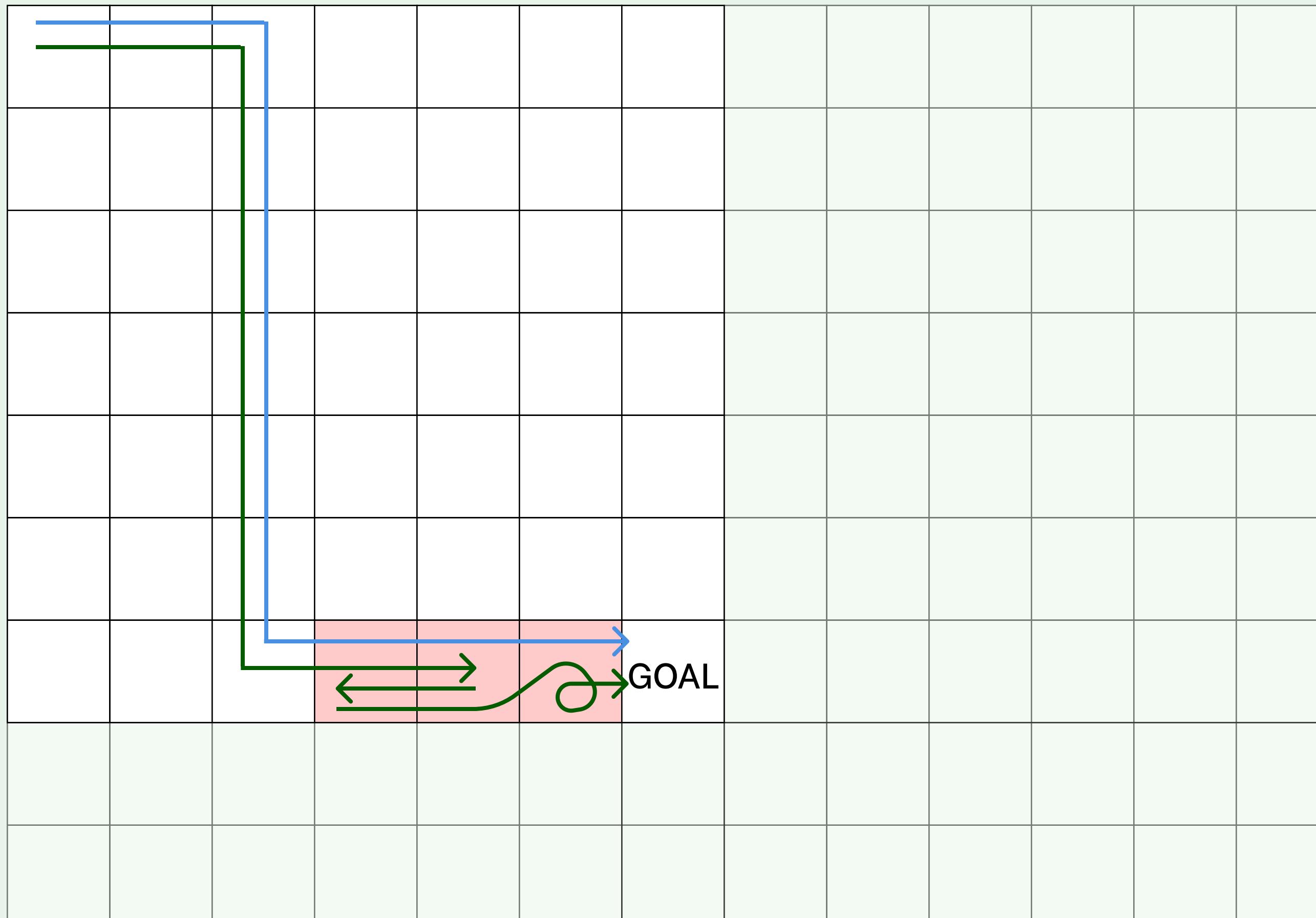


## 4. Findings



UNIVERSITÀ  
DI TRENTO

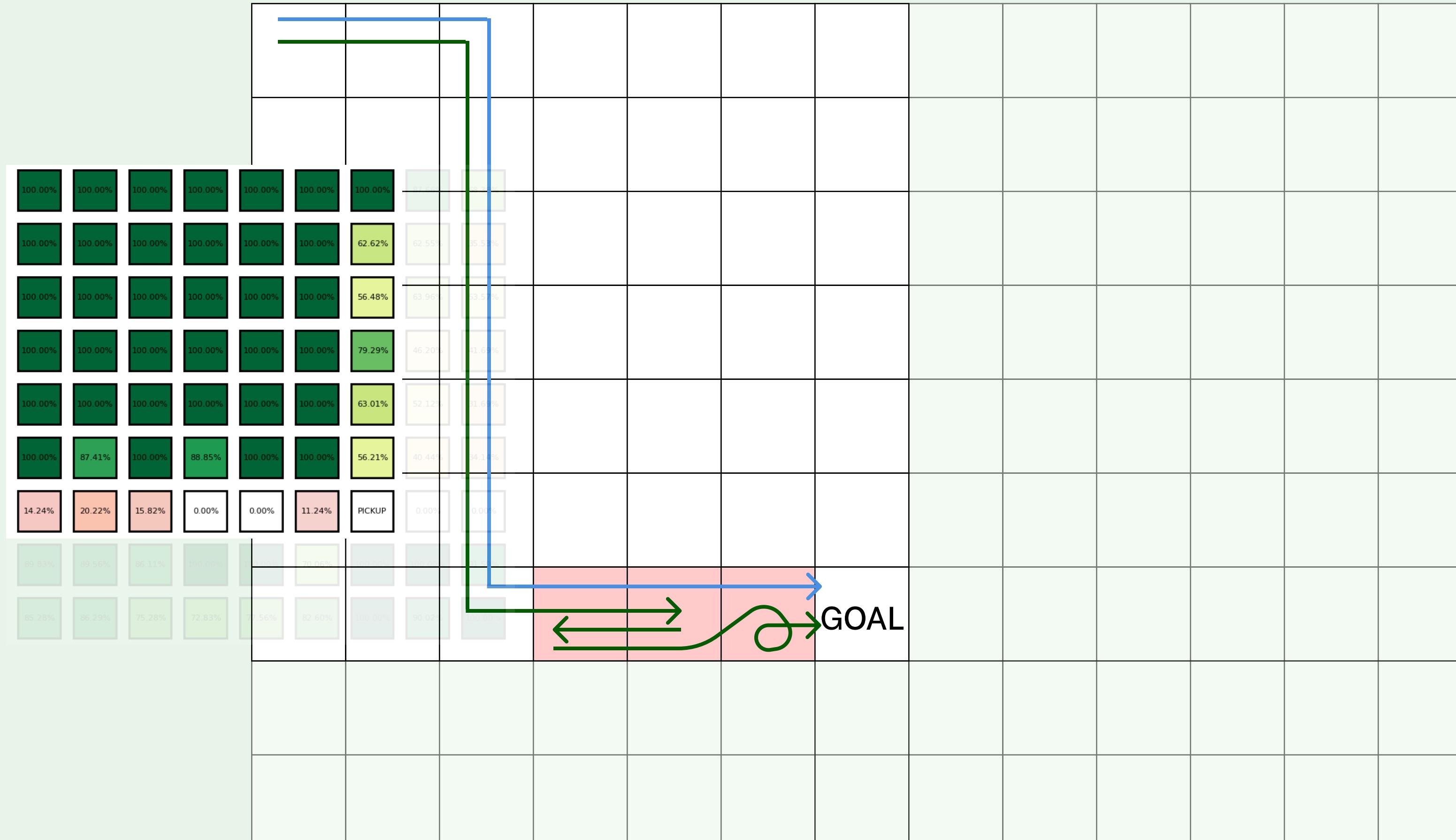
# Stateful Example



█ Optimal path  
█ Our path  
█ Repeated cells

Shared Nodes: 100%  
Ratio: 81%

# Stateful Example



Optimal path  
Our path  
Repeated cells  
Shared Nodes: 100%  
Ratio: 81%

# Stateful

Map Size	Stateless	Stateful	Difference
13x13	41	39	-5%
7x7	27	21	-29%
5x5	14	11	-27%
3x3	11	9	-22%

## 4. Findings



# Conclusions

## Strengths

Effective

~ 20% Less action required  
with history

Retrieving goals in big maps

Better models, better results

Similar % error as the size  
increase

## Weaknesses

Limited Explainability

Prone to error near the goal

Consistent problematic zones

Context size is a limitation

Similar % error as the size  
increase, but real number is a  
problem

## 4. Findings



# Thank You

## Exploring the Use of LLMs for Agent Planning: Strengths and Weaknesses

# Thank You

## Exploring the Use of LLMs for Agent Planning: Strengths and Weaknesses

# HM example where T and S were not always discarded



Example