

# Performance of Large Language Models in Domain-Specific and Underrepresented Languages: A Case Study on the Transportation Domain and Dutch Language

Authors: Thi M.D. Tran<sup>1</sup>, Davy Janssens<sup>1</sup>, Geert Wets<sup>1</sup>, Tom Brijs<sup>1</sup>, Burcu Can<sup>2</sup>, Jan Vuurstaek<sup>1</sup>, Lien Aerts<sup>1</sup>, Wim Ectors<sup>1</sup>

<sup>1</sup> UHasselt, Transportation Research Institute (IMOB), UHasselt, Martelarenlaan 42, 3500 Hasselt, Belgium

<sup>2</sup> Computing Science and Mathematics, University of Stirling, Stirling, UK, FK9 4LA

## INTRODUCTION

### Motivations

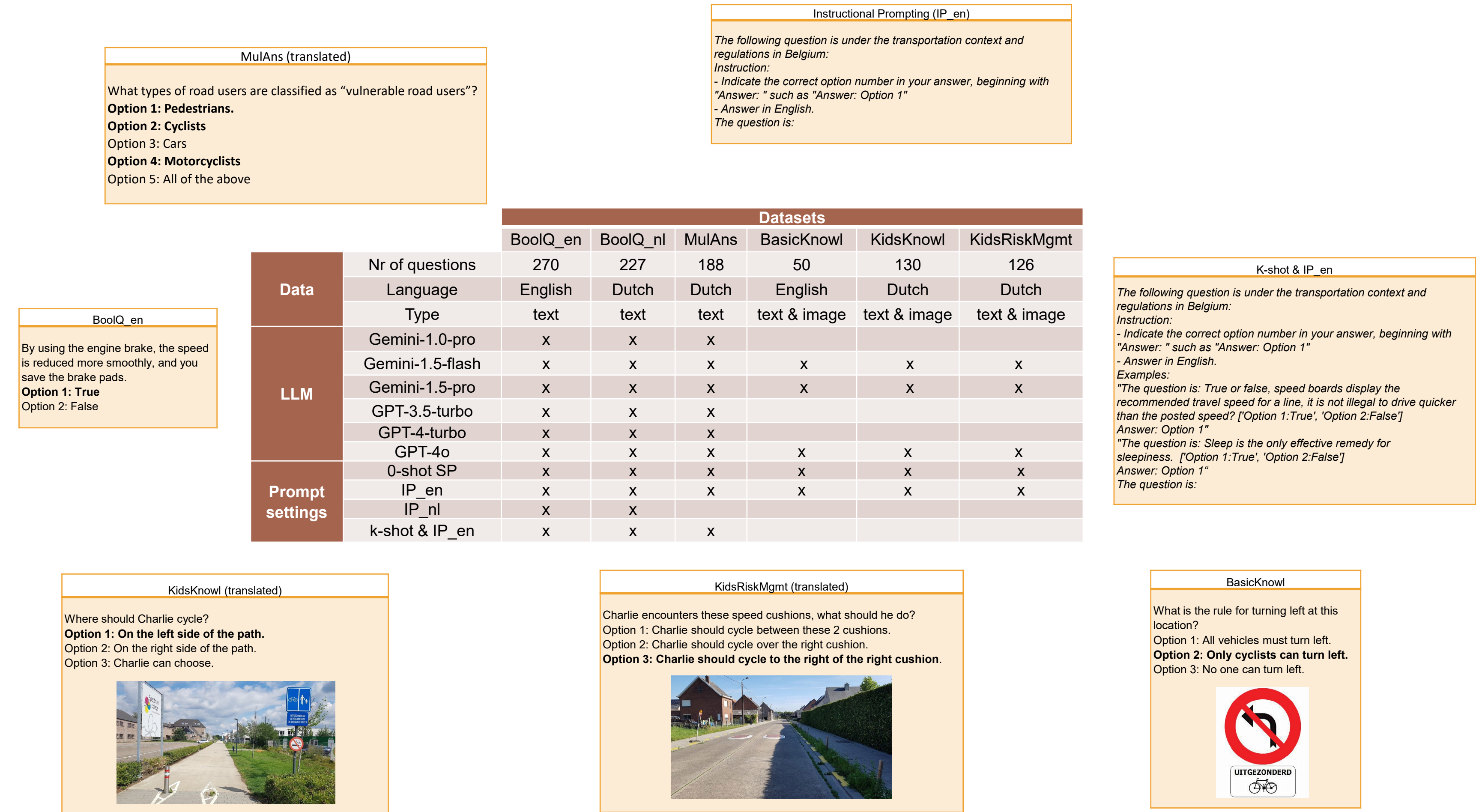
- LLMs excel in general tasks, primarily in English, but their performance in domain-specific reasoning and underrepresented languages (like Dutch) remains underexplored.
- Cross-lingual capabilities in specialized domains have not been widely studied.

### Objectives

- Enhance understanding of cross-lingual capabilities in specialized domains.
- Explore transfer learning potential for underrepresented languages like Dutch.
- Aid in selecting effective LLM foundation models for domain-specific applications.
- Provide performance benchmarks for LLMs in Dutch for transportation tasks.

## EXPERIMENTAL SETUP

- Question-answer format, extracted from our teaching and training materials at the School of Transportation Sciences and the Transportation Research Institute, UHasselt, Belgium
- 991 questions distributed across six datasets
- Include text only, text & image



## EVALUATION METHOD

$$\text{global accuracy (GA) (\%)} = \frac{\text{total correct responses}}{\text{total questions}} \times 100$$

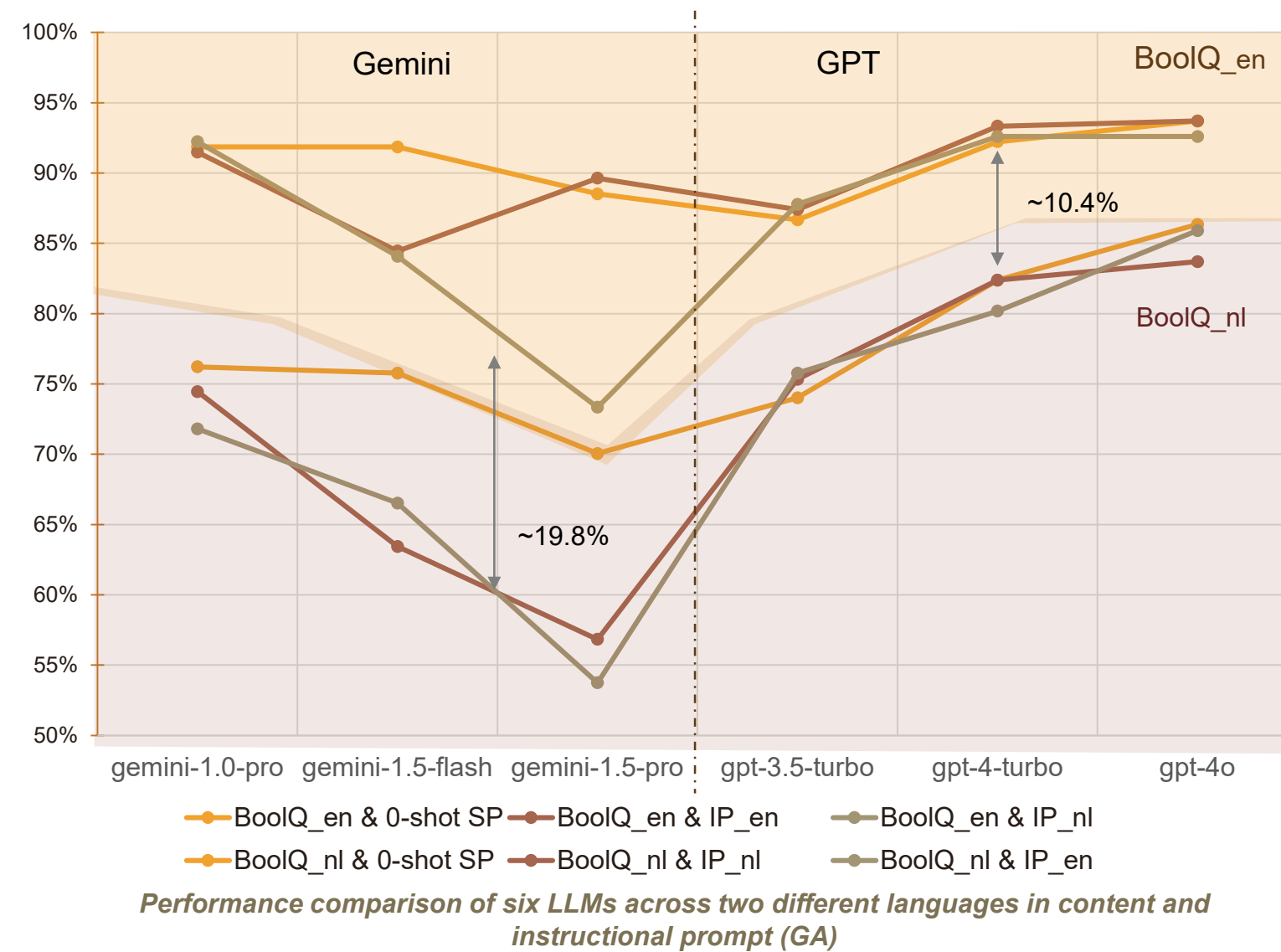
$$\text{local accuracy (LC) (\%)} = \frac{\sum \text{nr of correct indices per response}}{\sum \text{nr of indices per question}} \times 100$$

**Example Question:**  
What is allowed when under the influence of alcohol as a cyclist?  
**Option 1: Leave the bicycle behind and walk home.**  
Option 2: Proceed to cycle home.  
**Option 3: Push the bicycle home.**  
Response from an LLM: "The correct answer is Option 3"  
➤ Ground truth: [1,0,1]; LLM's answer: [0,0,1]  
➤ The global accuracy is 0%  
➤ The local accuracy 60%

## FINDINGS

### Impact of Language on LLM's Performance in the Transportation Domain

- LLM generally performs better with English than Dutch content.
- Language performance gap varies
- GPT models handle different content languages better than Gemini models
- Average accuracy variance: GPT models (~10.4%) vs. Gemini models (~19.8%).
- Language of the instruction plays a small impact

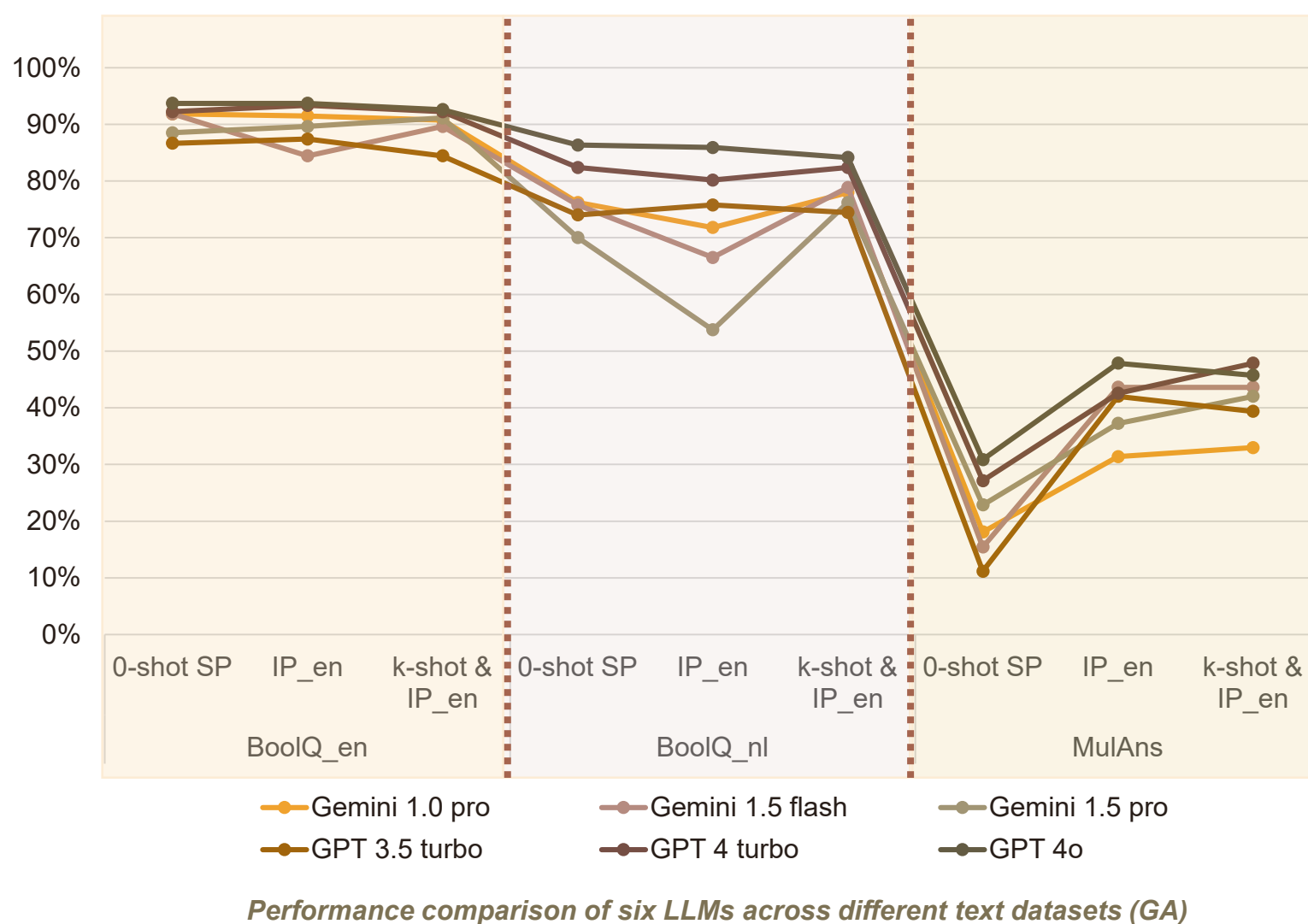


### Impact of Few-Shot Prompting

- Mixed results:
- Gemini models:** Performance improves when examples follow instructions
  - GPT models:** accuracy for GPT-4 turbo

### Impact of Instructional Prompting

- Enhances accuracy for multiple-answer questions across all models.
- Useful in handling complex tasks



### Performance Differences Between Question Complexity

- Multiple-answer questions:**
  - Require complex understanding, logical reasoning, and problem-solving.
  - Improved with instructional prompting.
- Boolean questions:**
  - Simpler due to its binary nature
  - Instructional prompting has minimal impact

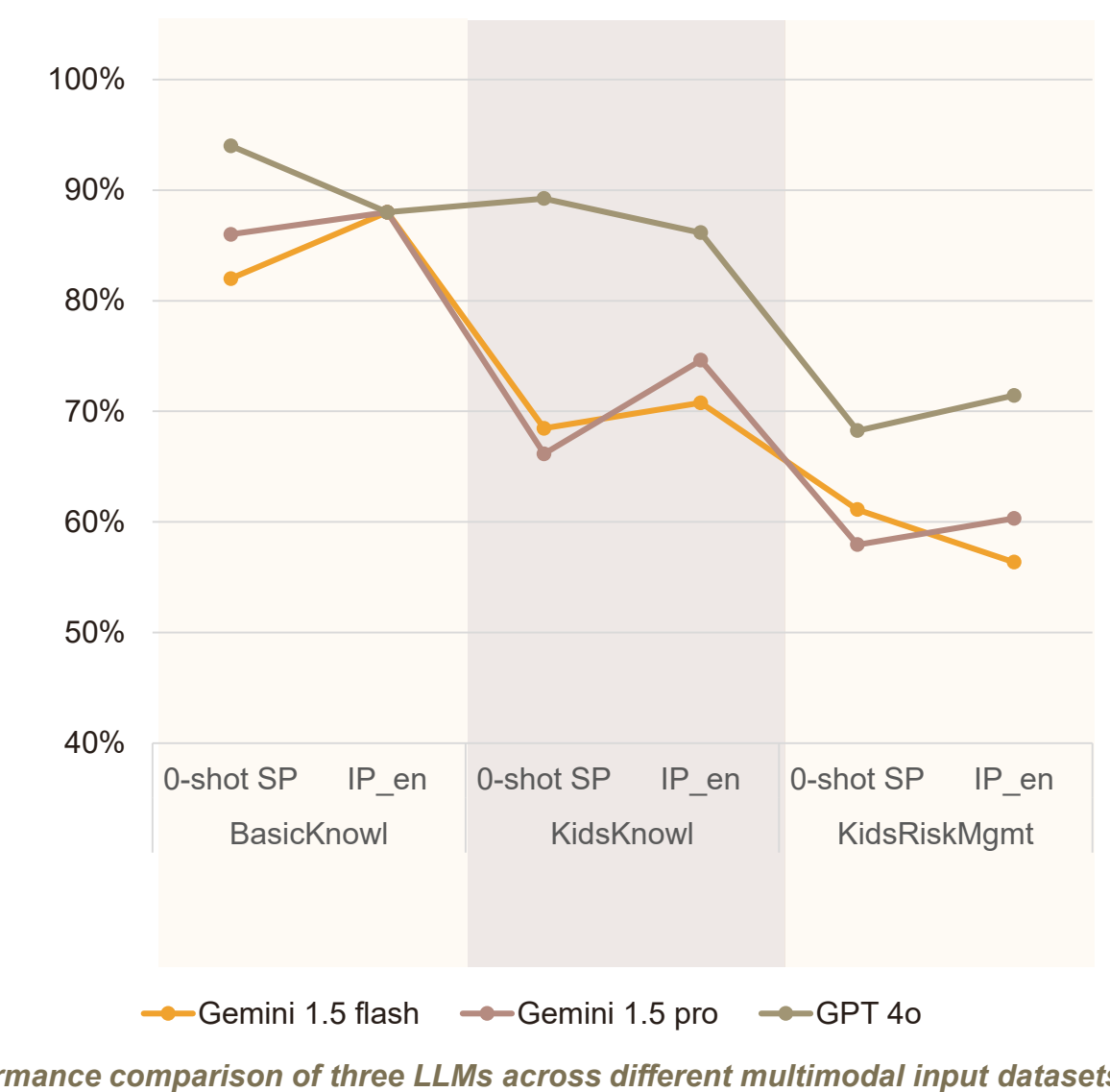
Models	0-shot SP		IP_en		k-shot & IP_en	
	GA	LA	GA	LA	GA	LA
Gemini-1.0-pro	18.09%	64.55%	31.38%	73.21%	32.98%	72.75%
Gemini-1.5-flash	15.43%	65.94%	43.62%	77.48%	43.62%	78.41%
Gemini-1.5-pro	22.87%	64.67%	37.23%	72.86%	42.02%	77.37%
GPT-3.5-turbo	11.17%	63.86%	42.02%	77.83%	39.36%	76.33%
GPT-4-turbo	27.13%	72.17%	42.55%	80.14%	47.87%	80.37%
GPT-4o	30.85%	74.13%	47.87%	82.22%	45.74%	83.03%

Global accuracy (GA) and local accuracy (LA) of MulAns dataset

- Local accuracy aligns with global accuracy
- Performance above random guessing (>50% LA)
- The gap between GP & LA's → potential for improvement with IP

### Performance on Text and Image-Based Transportation Tasks

- Better than random.
- GPT-4o** consistently outperformed **Gemini models**
- Better performance with common knowledge (BasicKnowl).
- Decreased performance with specialized knowledge (KidsKnowl).
- Lowest performance on complex, domain-specific tasks (KidsRiskMgmt)



## CONCLUSION

### Performance on Transportation Tasks:

- outperformed random guessing in both text-only and text-image scenarios.

### Language Sensitivity:

- Performed better with English content than Dutch.
- Less sensitivity to language differences of GPT's than Gemini's.

### Model Comparison:

- GPT-4o consistently.
- Gemini models, (Gemini 1.5 Pro): fluctuating performance and higher sensitivity to language.

### Implications:

- Provides a deeper understanding of LLM performance in transportation tasks, especially in Dutch.
- Offers valuable insights for selecting a suitable LLM for tasks involving specialized domains and underrepresented languages.

## LIMITATIONS

- Evaluated six Gemini and GPT models; excluded open-source multimodal LLMs
- Due to closed source, limited insights on the LLM architecture, flexibility
- Limited number of datasets and transportation scenarios.

## FUTURE WORK

- Expand the scope of evaluation to include a broader range of datasets and scenarios in transportation.
- Incorporate open-source multimodal LLMs
- Improve performances by fine-tuning in the Dutch language, fine-tuning LLMs tailored for specialized transportation tasks and contexts, optimizing prompting

## ACKNOWLEDGMENTS

- The Flemish government for funding this project
- Support from colleagues for question acquisition and data annotation.