# Toward a Multi Agent Approach for LLM-Based Dynamic Vehicle Control and Communication in Accidental Condition

Rafi Md Ashifujjman
*Graduate School of
Integrated Science and Technology
Shizuoka University*
Hamamatsu, Japan
ashifujjman@gmail.com

Naoki Fukuta
*College of Informatics,Academic Institute
Shizuoka University*
Hamamatsu, Japan
fukuta@inf.shizuoka.ac.jp

*Abstract—*

## I. Introduction

This paper explores the potential of integrating Large Language Models (LLMs) [1] into a multi-agent framework to improve decision making in autonomous vehicles (AVs), particularly in unknown and unsafe domains. Our main concern would be **to explore the use of large language models (LLM) for autonomous driving systems(ADS) as the main decision-making agent within a multi-agent framework to evaluate its reasoning ability to handle long-tail, False positive, and False negative scenarios.**

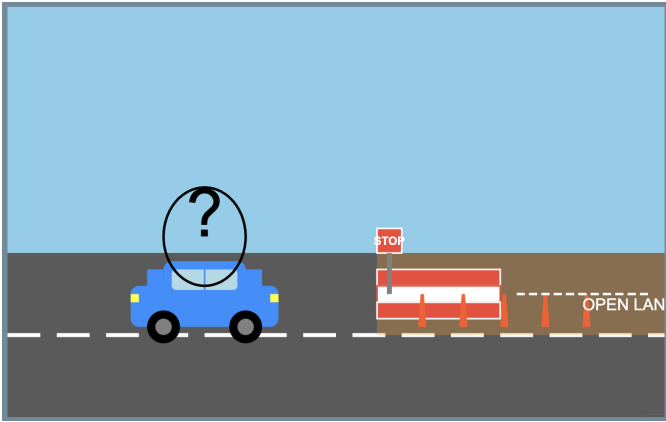## II. Motivation and Background



Fig. 1: What current AV might do in this situation?

Imagine a situation where a road is partially blocked due to construction, with traffic cones and a "STOP" sign indicating obstruction. A typical autonomous vehicle (AV) might see the "STOP" sign and interpret it as a complete road closure, coming to a complete stop because that is what its pre-programmed data tell it to do. However, a human driver in the same situation might use common sense, realize that the road is only partially closed, and safely continue through the open section.
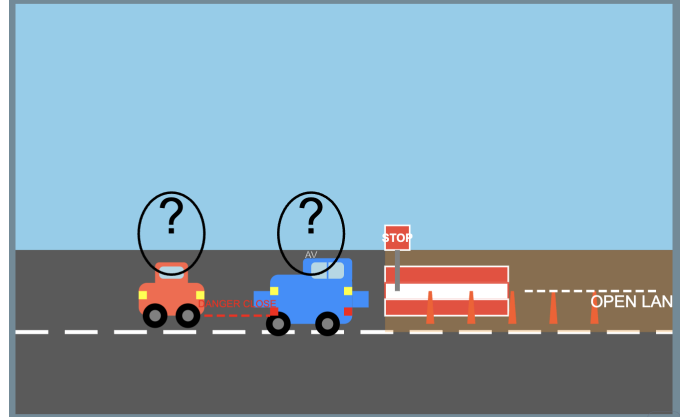


Fig. 2: Lack of communication between vehicles may lead to a potential accident.

Now, let us make it more complicated: an autonomous truck comes to a sudden stop when it sees the "STOP" sign. There is another car, driven by a human or another AV, right behind it. This sudden stop could easily cause a crash. Now imagine a human driver seeing a partially blocked road and choosing to carefully go through the open part. As they proceed, they notice an autonomous vehicle (AV) approaching in the opposite direction. Using common sense, the human driver would slow down to avoid any risk of collision. However, the AV does not recognize the approaching car as a possible obstacle. It might not slow down and could keep going at its normal speed, which could lead to a crash.

However, if the AV had better reasoning abilities, it might recognize that the road is only partially closed and proceed safely through the open part. And if the AV could communicate with the vehicle behind it or if the lane was suddenly closed, it could alert the other car about its actions, preventing a possible accident. A fundamental challenge in traditional

autonomous vehicles (AVs) is to prepare for all truly unknown situations [2]. Any scenario we create is technically "known" to us, even if it feels unfamiliar to the AV, making it difficult to test how AVs handle truly unknown-unsafe situations. In addition, it is difficult to recreate complex real-world scenarios in simulation environments. For example, it is difficult to accurately simulate unpredictable human actions or complicated environmental conditions. In real life, there are countless rare or unexpected situations that are nearly impossible to include in training data sets for traditional autonomous driving systems (ADS). Human drivers use common sense to deal with these situations, so we need to explore whether LLMs can also show human-like reasoning to handle them effectively. To address these challenges, we propose a possible approach on integrating Large Language Models (LLMs) into a multi-agent framework for autonomous driving control system. LLMs have the potential to provide common sense or reasoning ability that traditional systems lack, helping AVs better understand the context of complex real-world scenarios. Incorporating Vehicle-to-Vehicle communication using multi-agent, which allows vehicles to exchange information with other vehicles. This capability is also essential for creating realistic and complex driving scenarios.

## III. METHODOLOGY

## IV. CONCLUSION

Reasoning is a fundamental aspect of human intelligence, essential for problem solving, decision making, and critical thinking. In recent years, Large Language Models (LLMs) have demonstrated emergent abilities [3], such as in-context learning [4], role play [5] and analogical reasoning [6]. These abilities allow LLMs to go beyond natural language processing problems to facilitate a wider range of tasks, such as code generation [7], robotic control [8], and autonomous agents [9]. Among these abilities, human-like reasoning has garnered significant attention from both academia and industry, since it demonstrates great potential for LLMs to generalize to complex real-world problems through abstract and logical reasoning. A notable breakthrough in this area is the 'chain of thought' prompting technique [6], which can elicit step-by-step human-like reasoning processes at test time without any additional training.

Our goal is to contribute to the development of safer and more reliable Level 5 [10] autonomous vehicles. We aim to contribute for level 5 (complete automation [11]). There are numerous instances in which traditional autonomous vehicles are capable of making effective decisions in familiar scenarios using pre-trained models. However, they may have trouble facing new or unusual situations where their programmed knowledge is not enough. Human drivers can use common sense and past experiences to handle unexpected events, such as knowing that traffic cones on a moving truck are not dangerous. However, current AV systems, even those that use rule-based approaches or reinforcement learning (RL [12]), will struggle in these scenarios. This limitation is especially

evident in what we call "unknown-unsafe" region situations where the correct action is not immediately clear. [13]

During this investigation with various LLM, we proceed under two assumptions: 1) Certain AIs can transform real-world situations into text-based explanations, which are subsequently used as input for the LLMs, and 2) The outputs generated by the LLMs can be translated into actual decisions made by the decision agent of AV by interpreting the responses from the LLMs. We are investigating the effectiveness of various Large Language Models (LLMs) for use as decision-making agents in autonomous driving systems. Our initial objective is to answer the research question, **Can an agent make decisions with the help of an LLM without undergoing a specific learning process for each new situation?**

| Recommendation | Model Name | Benchmark Performance |
|---|---|---|
| LM Studio Recommended Models | Mistral 7B | 74.6% (MMLU), 83.5% (GSM8K) |
| | OpenHermes 2.5 | 72.1% (MMLU), 79.2% (BIG-bench) |
| | DeepSeek 7B | 76.3% (MMLU), 85.4% (GSM8K) |
| Hugging Face Open-Source Models | LLaMA 2 13B | 75.8% (MMLU), 80.7% (BIG-bench) |
| | Mixtral 8x7B | 78.2% (MMLU), 88.1% (GSM8K) |
| Top Commercial Services | Claude 3 Sonnet | 82.3% (MMLU), 90.2% (BIG-bench) |
| | GPT-4 Turbo | 85.6% (MMLU), 92.4% (GSM8K) |
| | Falcon 7B/40B | 73.4% (MMLU), 78.9% (BIG-bench) |

Fig. 3: Selected Models for Experimentation

Hugging Face is a platform where various open LLM models are available. LM Studio is a platform to test and integrate LLMs available at Hugging Face into a locally available system on ordinary computing devices even without powerful GPUs. Now, what factors were considered in the selection of LLMs for this study? One key factor in our model selection was quantization, which optimizes model performance while reducing computational requirements. This research [14] indicates that 8-bit quantization enables the majority of LLMs to maintain a performance level comparable to their nonquantized equivalents, regardless of model size (e.g. 7B to 70B parameters). Moreover, LLMs that are quantized to 4 bits can also up hold similar performance to their nonquantized versions across most benchmarks. This approach achieves memory reduction of 50 to 75 % while preserving precision in complex tasks such as reasoning, decision making, and domain-specific applications. The models chosen for this investigation were selected using a structured approach based on three key factors: popularity and performance in text-to-text generation, LM Studio recommendations for optimized accuracy, and comparative evaluations of the top commercial services. Specifically, 1) Some models were chosen based on the highest number of downloads from Hugging Face,

ensuring widespread adoption and benchmark effectiveness. 2) Models suggested by LM Studio were included due to their strong performance and compatibility with local inference environments. 3) The best commercial services were selected based on their comparative performance with open-source counterparts, prioritizing accuracy, interpret ability, and real-time inference. The selected models are shown in Fig. 3. However, not all Hugging Face models are fully compatible with the LM Studio runtime. As a result, some models were tested directly on the Hugging Face interface to avoid compatibility issues. Furthermore, while models such as OpenAI's o1, DeepSeek R1, and Llama 3.1 405B have demonstrated strong benchmark performance, they were not included in this study due to limited quantization support, lack of local deployment feasibility, and restricted open-source availability. These limitations make them less practical for our research. To test this model, we will create simple text-based scenarios. Some of these scenarios will involve situations that can only be solved through logical reasoning, while others will require recognizing and applying common or well-known information (e.g. stopping when seeing a red traffic light). These scenarios will be used as input for various LLMs to assess how effectively they handle both types of challenges. The experiments will be conducted using the LLM selected in Fig.4. The purpose was to assess whether the models could accurately recommend actions based on the complexity and ambiguity of each scenario. A major consideration was consistency. **Does repeated querying of LLMs on the same scenario lead to variations in decision-making outcomes?** To test this, we posed each scenario to each model at least 20 times. This helped us to check whether the models could provide stable and repeatable outputs or whether their decisions varied. This consistency is crucial for tasks like AV decision-making, where unpredictable output could lead to safety risks.

Fig. 4: Text-base scenario

Fig. 5: LLM Model Response

We presented a text-base scenario showing the autonomous vehicle (AV) and human-driven car (H) approaching a green traffic light similar to the one in Fig. 4, to the language models. We used the following prompt: What would the AV do in this situation? Please, just answer STOP or FORWARD. We ran this same prompt at least 20 times in OpenHermes-2.5-Mistral-7B and analyzed the results.

```python
import numpy as np

# Example responses from the LLM (1 for FORWARD, 0 for STOP)
responses = [1] * 19 + [0]  # 19 FORWARD + 1 FORWARD (counted as 1)

# Calculate mean and variance
mean_response = np.mean(responses)
variance_response = np.var(responses)

print("Mean:", mean_response)
print("Variance:", variance_response)
```

PROBLEMS    OUTPUT    PORTS    TERMINAL

```
rafimdashifujjman@Rafis-Laptop code % python3 mean.py
Mean: 0.95
Variance: 0.0475
```
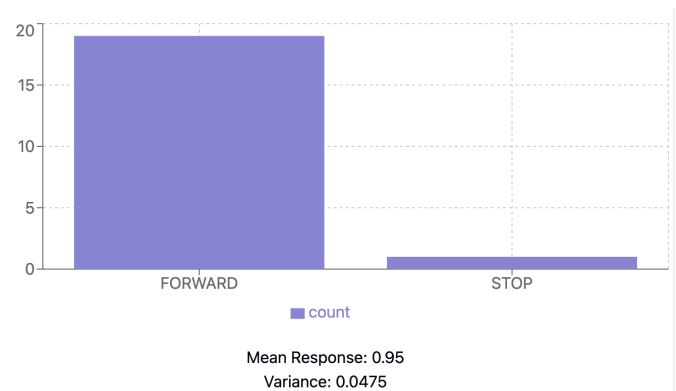
Fig. 6: LLM Model Response Analysis
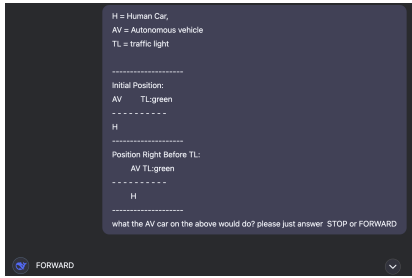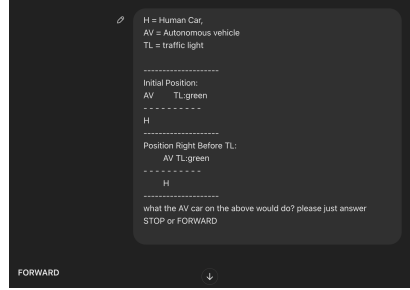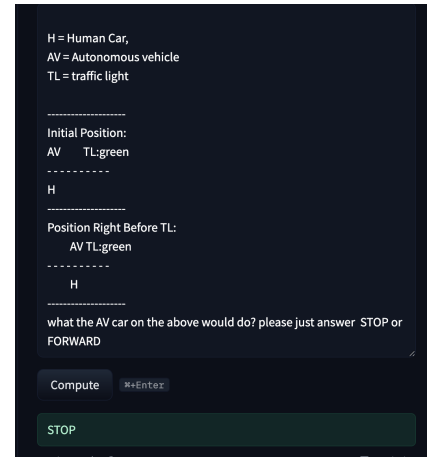
Mean Response: 0.95
Variance: 0.0475
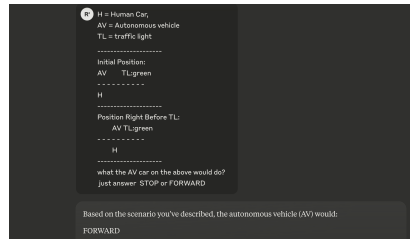
Fig. 7: LLM Model Response Analysis

(a) Output from DeepQ

(b) Output from ChatGPT

(c) Output from google/flan-t5-large

(d) Output from Claude

Fig. 8: Comparison of responses from different LLMs on the task

## REFERENCES

[1] R. Zhang, X. Guo, W. Zheng, C. Zhang, K. Keutzer, and L. Chen, "Instruct large language models to drive like humans," *arXiv preprint arXiv:2406.07296*, 2024.

[2] T. Singh, E. van Hassel, A. Sheorey, and M. Alirezaei, "A systematic approach for creation of sotif's unknown unsafe scenarios: An optimization based method," tech. rep., SAE Technical Paper, 2024.

[3] J. Wei, Y. Tay, R. Bommasani, C. Raffel, B. Zoph, S. Borgeaud, D. Yogatama, M. Bosma, D. Zhou, D. Metzler, *et al.*, "Emergent abilities of large language models," *arXiv preprint arXiv:2206.07682*, 2022.

[4] Q. Dong, L. Li, D. Dai, C. Zheng, J. Ma, R. Li, H. Xia, J. Xu, Z. Wu, T. Liu, *et al.*, "A survey on in-context learning," *arXiv preprint arXiv:2301.00234*, 2022.

[5] M. Shanahan, K. McDonell, and L. Reynolds, "Role play with large language models," *Nature*, vol. 623, no. 7987, pp. 493–498, 2023.

[6] J. Wei, X. Wang, D. Schuurmans, M. Bosma, F. Xia, E. Chi, Q. V. Le, D. Zhou, *et al.*, "Chain-of-thought prompting elicits reasoning in large language models," *Advances in neural information processing systems*, vol. 35, pp. 24824–24837, 2022.

[7] J. Gehring, K. Zheng, J. Copet, V. Mella, Q. Carbonneaux, T. Cohen, and G. Synnaeve, "Rlef: Grounding code llms in execution feedback with reinforcement learning," *arXiv preprint arXiv:2410.02089*, 2024.

[8] M. Ahn, A. Brohan, N. Brown, Y. Chebotar, O. Cortes, B. David, C. Finn, C. Fu, K. Gopalakrishnan, K. Hausman, *et al.*, "Do as i can, not as i say: Grounding language in robotic affordances," *arXiv preprint arXiv:2204.01691*, 2022.

[9] X. Deng, Y. Gu, B. Zheng, S. Chen, S. Stevens, B. Wang, H. Sun, and Y. Su, "Mind2web: Towards a generalist agent for the web," *Advances in Neural Information Processing Systems*, vol. 36, pp. 28091–28114, 2023.

[10] A. Balasubramaniam and S. Pasricha, "Object detection in autonomous vehicles: Status and open challenges," *arXiv preprint arXiv:2201.07706*, 2022.

[11] M. Raza, "Autonomous vehicles: levels, technologies, impacts and concerns," *International Journal of Applied Engineering Research*, vol. 13, no. 16, pp. 12710–12714, 2018.

[12] L. P. Kaelbling, M. L. Littman, and A. W. Moore, "Reinforcement learning: A survey," *Journal of artificial intelligence research*, vol. 4, pp. 237–285, 1996.

[13] D. Fu, X. Li, L. Wen, M. Dou, P. Cai, B. Shi, and Y. Qiao, "Drive like a human: Rethinking autonomous driving with large language models," in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 910–919, 2024.

[14] R. Jin, J. Du, W. Huang, W. Liu, J. Luan, B. Wang, and D. Xiong, "A comprehensive evaluation of quantization strategies for large language models," *ArXiv*, vol. abs/2402.16775, 2024.