

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

Abstract—In autonomous vehicles, making the right decision in unknown traffic situations remains a challenge. Traditional AVs rely on pre-trained models, which often struggle in these cases where human-like reasoning is required. In this research, we examine how LLM as the decision agent responds to scenario-specific text prompts derived from real-world traffic situations. In addition, we propose a Communication Agent that enables vehicle-to-vehicle (V2V) information sharing, such as speed, obstacles, and future intentions, in a structured format to provide contextual support for the decision agent. This paper proposes a preliminary approach to assess whether LLMs within a multi-agent framework, when provided with structured prompts and contextual data (unknown traffic situation), can support consistent and human-like decision-making.

Index Terms—autonomous vehicles, large language models (LLMs), multi-agent framework, vehicle-to-vehicle (V2V) communication, decision making agent, communication agent.

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 uncertain traffic situations. Our goal is to contribute to the development of safer and more reliable Level 5 [2] autonomous vehicles. We aim to contribute for level 5 (complete automation [3]). 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 directly covered in such situations. 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 [4]), 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 [5]. During this investigation with various LLM, we proceed our investigation 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 (Research Question 1) Can an agent make decisions with the help of an LLM without undergoing a specific learning process for each new situation? In our research, we will assume that vehicle-to-vehicle (V2V) communication has been achieved through a standardized networking protocol, such as LAN-based direct communication or a low-latency wireless system. Under this assumption, each vehicle can communicate with other vehicles within a defined radius (e.g., 10 meters). When a vehicle enters the information sharing radius of another vehicle, they exchange data regarding their current status, such as obstacle detection, traffic conditions, speed, direction, and future intended actions. To evaluate the impact of V2V communication on autonomous decision-making, we will investigate how Large Language Models (LLMs) process and utilize this information. Specifically, we will examine whether the incorporation of V2V exchanged data improves decision accuracy and enables AVs to make more precise and context-aware decisions. However, V2V communication introduces a critical question, (Research Question 2) Is a single centralized decision-making agent sufficient to process all incoming V2V data and make autonomous driving decisions, or is a multi-agent system required for improved efficiency and scalability? In this paper, we investigate the answers to these two research questions.

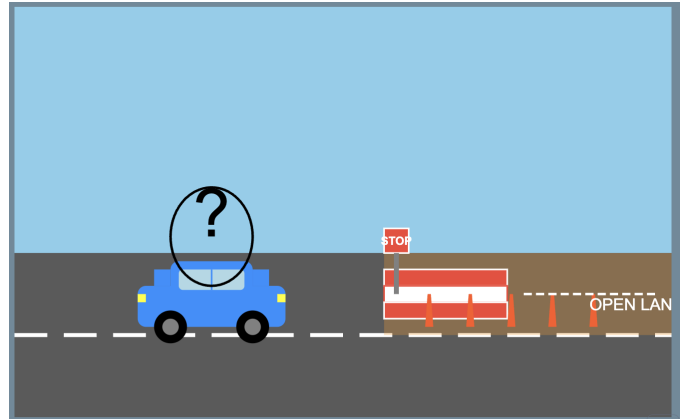


Fig. 1: What current AV might do in this situation?

II. MOTIVATION AND BACKGROUND

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.

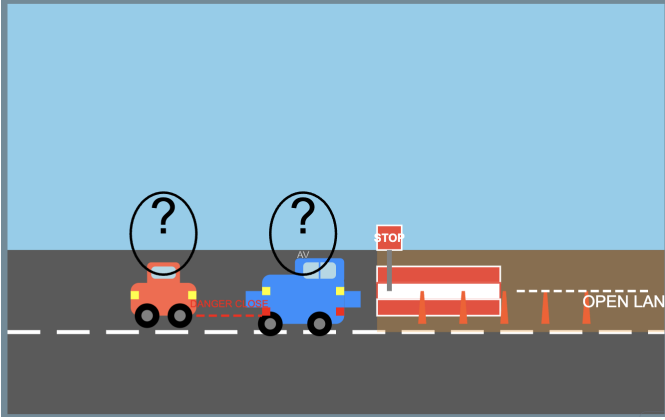


Fig. 2: Lack of communication between vehicles may lead to a potential accident.

Now, let us make it more complicated that 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 encountering a partially blocked road. To navigate safely, the driver shifts into the adjacent open lane, relying on instinct and situational awareness. Anticipating the possibility of an oncoming vehicle from the opposite direction or from behind, the driver proactively slows down, assessing the risk and taking the necessary precautions to avoid a collision. In contrast, a traditional AV, which lacks human-like reasoning and predictive thinking, may not recognize the potential danger. Without contextual understanding, it can continue at its normal speed, assuming the lane is clear, increasing the risk of a head-on collision.

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 [6]. 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 could be countless unexpected situations that are nearly impossible to include in training data sets for traditional autonomous driving systems (ADS) to avoid accidents. Human drivers use common sense to deal with these situations, so through this research we will explore whether LLMs can also show human-like reasoning to handle them effectively or not.

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. PRELIMINARY

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 the Experimentation

Hugging Face¹ is a platform where various open LLM models are available. LM Studio² is a platform to test and integrate LLM available on Hugging Face a locally available system on ordinary computing devices, even without powerful GPUs. Now, the question is 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 [7] 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 non-quantized versions across most benchmarks. This approach achieves memory reduction of 50 to 75 % while preserving precision in complex tasks such as reasoning, decision making,

¹<https://huggingface.co/>

²<https://lmstudio.ai/>

```

AV = Autonomous Vechiels,
S = Stop Sign of a partial road,
Other lane is open.

AV      S
-----

Initial Position:
AV      S
-----

Position Right Before STOP sign:
      AV S
-----

Simulation ended.

```

(a) Traffic Light

```

H = Human Car,
AV = Autonomous Vechiels,
TL = traffic light.

Initial Position:
AV      TL:green
-----

H
-----

Position Right Before TL:
      AV TL:green
-----

H
-----

Simulation ended.

```

(b) Stop sign of Partial road

```

H = Human,
AV = Autonomous Vechiels,

Initial Position:

AV
-----

Position Right Before H:

      AV      H
-----

Simulation ended.

```

(c) Human suddenly crossing road

Fig. 4: Illustration of real life traffic scenarios

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.

To test these selected models, we have prepared a set of simple, text-based simulation scenarios using Python, as illustrated in Fig. 4. These scenarios are designed to represent situations that cover situations that require logical reasoning or the application of common sense knowledge, such as recognizing a red traffic light and stopping accordingly. 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.3.

IV. OUR APPROACH

The main phase of this research focuses on using the shared data of the communication agent within our proposed multi-agent framework. The Communication Agent ensures that vehicles share critical information, including obstacle detection, traffic conditions, speed, and intended actions, in a standardized format within a defined radius, the expected format is demonstrated in Fig. 5. This structured data exchange is expected to significantly improve the accuracy of Decision Agent by providing a more comprehensive contextual understanding of the driving environment. For initial testing, we selected a basic traffic scenario involving an autonomous vehicle (AV) approaching a green traffic light, as illustrated in Fig. 4(a), and presented it to various language models to evaluate their decision-making behavior. The scenario was communicated through the prompt: 'What would the AV on the above do? Please just answer STOP or FORWARD.' The selected models included Claude³ 3.7 Sonnet, GPT⁴-4 Turbo, Falcon⁵ 3 7B and open-source alternatives such as OpenHermes-2.5-Mistral-7B and google/flan-t5-large. The responses of these models, sourced from Hugging Face, LM Studio, and commercial providers, are shown in Figure 6. GPT-4 and Falcon provided direct and expected responses 'FORWARD', which aligns with the logic of the situation. However, Claude deviated from the expected format by including an explanatory response, despite the prompt requesting a one-word answer. OpenHermes-2.5

³<https://claude.ai/chats>

⁴<https://chat.openai.com/>

⁵<https://chat.falconllm.tii.ae/>

```
1 import json
2 from dataclasses import dataclass, asdict
3 from typing import Dict, List, Any
4 import random
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65
66
67
68
69
70
71
72
73
74
75
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94
95
96
97
98
99
100
101
102
103
104
105
106
107
108
109
110
111
112
113
114
115
116
117
118
119
120
121
122
123
124
125
126
127
128
129
130
131
132
133
134
135
136
137
138
139
140
141
142
143
144
145
146
147
148
149
150
151
152
153
154
155
156
157
158
159
160
161
162
163
164
165
166
167
168
169
170
171
172
173
174
175
176
177
178
179
180
181
182
183
184
185
186
187
188
189
190
191
192
193
194
195
196
197
198
199
200
201
202
203
204
205
206
207
208
209
210
211
212
213
214
215
216
217
218
219
220
221
222
223
224
225
226
227
228
229
230
231
232
233
234
235
236
237
238
239
240
241
242
243
244
245
246
247
248
249
250
251
252
253
254
255
256
257
258
259
260
261
262
263
264
265
266
267
268
269
270
271
272
273
274
275
276
277
278
279
280
281
282
283
284
285
286
287
288
289
290
291
292
293
294
295
296
297
298
299
300
301
302
303
304
305
306
307
308
309
310
311
312
313
314
315
316
317
318
319
320
321
322
323
324
325
326
327
328
329
330
331
332
333
334
335
336
337
338
339
340
341
342
343
344
345
346
347
348
349
350
351
352
353
354
355
356
357
358
359
360
361
362
363
364
365
366
367
368
369
370
371
372
373
374
375
376
377
378
379
380
381
382
383
384
385
386
387
388
389
390
391
392
393
394
395
396
397
398
399
400
401
402
403
404
405
406
407
408
409
410
411
412
413
414
415
416
417
418
419
420
421
422
423
424
425
426
427
428
429
430
431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497
498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546
547
548
549
550
551
552
553
554
555
556
557
558
559
560
561
562
563
564
565
566
567
568
569
570
571
572
573
574
575
576
577
578
579
580
581
582
583
584
585
586
587
588
589
590
591
592
593
594
595
596
597
598
599
600
601
602
603
604
605
606
607
608
609
610
611
612
613
614
615
616
617
618
619
620
621
622
623
624
625
626
627
628
629
630
631
632
633
634
635
636
637
638
639
640
641
642
643
644
645
646
647
648
649
650
651
652
653
654
655
656
657
658
659
660
661
662
663
664
665
666
667
668
669
670
671
672
673
674
675
676
677
678
679
680
681
682
683
684
685
686
687
688
689
690
691
692
693
694
695
696
697
698
699
700
701
702
703
704
705
706
707
708
709
710
711
712
713
714
715
716
717
718
719
720
721
722
723
724
725
726
727
728
729
730
731
732
733
734
735
736
737
738
739
740
741
742
743
744
745
746
747
748
749
750
751
752
753
754
755
756
757
758
759
760
761
762
763
764
765
766
767
768
769
770
771
772
773
774
775
776
777
778
779
780
781
782
783
784
785
786
787
788
789
790
791
792
793
794
795
796
797
798
799
800
801
802
803
804
805
806
807
808
809
810
811
812
813
814
815
816
817
818
819
820
821
822
823
824
825
826
827
828
829
830
831
832
833
834
835
836
837
838
839
840
841
842
843
844
845
846
847
848
849
850
851
852
853
854
855
856
857
858
859
860
861
862
863
864
865
866
867
868
869
870
871
872
873
874
875
876
877
878
879
880
881
882
883
884
885
886
887
888
889
890
891
892
893
894
895
896
897
898
899
900
901
902
903
904
905
906
907
908
909
910
911
912
913
914
915
916
917
918
919
920
921
922
923
924
925
926
927
928
929
930
931
932
933
934
935
936
937
938
939
940
941
942
943
944
945
946
947
948
949
950
951
952
953
954
955
956
957
958
959
960
961
962
963
964
965
966
967
968
969
970
971
972
973
974
975
976
977
978
979
980
981
982
983
984
985
986
987
988
989
990
991
992
993
994
995
996
997
998
999
1000
```

Fig. 5: An illustration of Communication Agent

responded with a clear ‘FORWARD’, demonstrating alignment with the input instructions. Surprisingly, the google/flan-t5-large output differed from all others by responding ‘STOP’, contradicting the intended logic of the 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, each scenario was submitted to the models multiple times. This approach allowed us to observe whether the models produced stable and repeatable responses or if their decisions varied unpredictably. Consistency is especially critical in autonomous vehicle (AV) applications, where uncertain output can cause serious safety concerns. For this analysis, we ran the prompt at least 20 times on OpenHermes-2.5-Mistral-7B, and the consistency results are presented in Fig. 6. Consistency was calculated by mapping ‘FORWARD’ = 1 and ‘STOP’ = 0. A mean closer

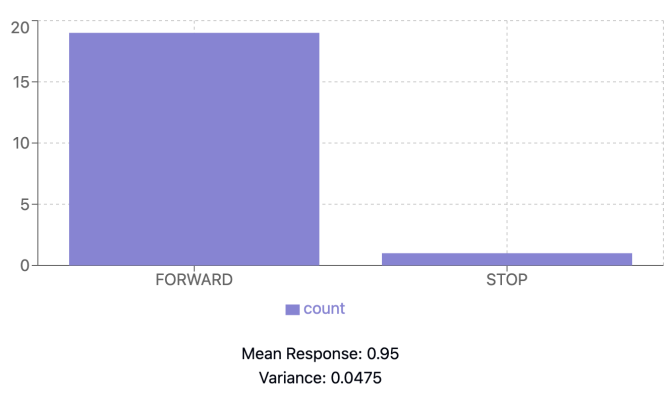


Fig. 6: LLM Model Consistency on the Tested Scenario

to 1 indicates greater consistency with the expected response.

V. CONCLUSION

We presented an approach for leveraging Large Language Models (LLMs) to enhance reasoning capabilities in uncertain traffic situations within a multi-agent framework for autonomous vehicles. We have investigated LLMs from Hugging Face, LM Studio, as well as major commercial services such as Claude, GPT-4, and Falcon. Their responses to basic traffic scenarios were analyzed, revealing varying levels of consistency and accuracy. The preliminary evaluation reveals both the promise and the current limitations of this approach, highlighting the need for improved consistency in the model responses. The proposed approach aims at the realization of adaptive autonomous systems capable of human-like reasoning in unpredictable situations. By integrating contextual data derived from V2V communication with LLM reasoning, our goal is to bridge the critical gap in decision-making accuracy for autonomous vehicles.

REFERENCES

[1] Y. Chang, X. Wang, J. Wang, Y. Wu, L. Yang, K. Zhu, H. Chen, X. Yi, C. Wang, Y. Wang, et al., “A survey on evaluation of large language models,” *ACM transactions on intelligent systems and technology*, vol. 15, no. 3, pp. 1–45, 2024.

[2] J. Cho and J. Heo, “The more you know, the more you can trust: Drivers’ understanding of the advanced driver assistance system,” in *HCI in Mobility, Transport, and Automotive Systems. Automated Driving and In-Vehicle Experience Design: Second International Conference, MobiTAS 2020, Held as Part of the 22nd HCI International Conference, HCII 2020, Copenhagen, Denmark, July 19–24, 2020, Proceedings, Part I 22*, pp. 230–248, Springer, 2020.

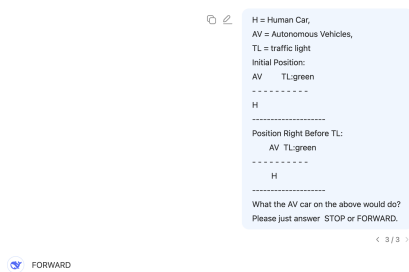
[3] M. Raza, “Autonomous vehicles: levels, technologies, impacts and concerns,” *International Journal of Applied Engineering Research*, vol. 13, no. 16, pp. 12710–12714, 2018.

[4] 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.

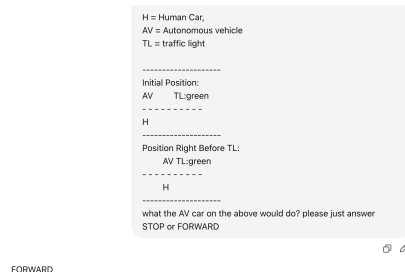
[5] D. Fu, X. Li, L. Wen, M. Dou, P. Cai, B. Shi, and Y. Qiao, “Drive like a human: Rethinking autonomous driving with large language models,” in *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 910–919, 2024.

[6] T. Singh, E. van Hassel, A. Sheorey, and M. Alirezai, “A systematic approach for creation of sotif’s unknown unsafe scenarios: An optimization based method,” tech. rep., SAE Technical Paper, 2024.

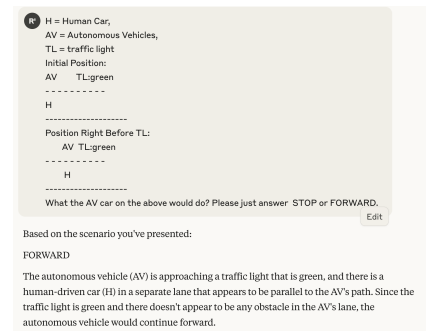
[7] R. Jin, J. Du, W. Huang, W. Liu, J. Luan, B. Wang, and D. Xiong, “A comprehensive evaluation of quantization strategies for large language models,” *ArXiv*, vol. abs/2402.16775, 2024.



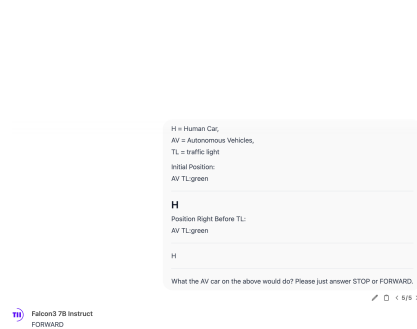
(a) Output from DeepSeek-V3



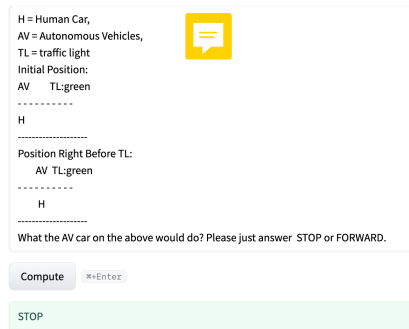
(b) Output from ChatGPT



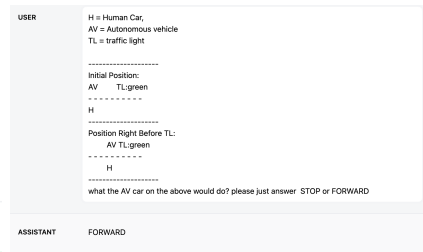
(c) Output from Claude 3.7 Sonnet



(d) Output from Falcon3 7B



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



(f) Output from OpenHermes-2.5-Mistral-7B

Fig. 7: Responses from different LLMs