

# 人間のようには運転する：大規模言語モデルによる自律走行の再考

Daocheng Fu<sup>1\*</sup> Xin Li<sup>1,2\*</sup> Licheng Wen<sup>1\*</sup> Min Dou<sup>1</sup> Pinlong Cai<sup>1</sup>  
Botian Shi<sup>1†</sup> Yu Qiao<sup>1</sup>

<sup>1</sup>Shanghai AI Lab, <sup>2</sup>East China Normal University  
{fudaocheng, lixin, wenlicheng, doumin, caipinlong, shibotian, qiaoyu}@pjlab.org.cn

## Abstract

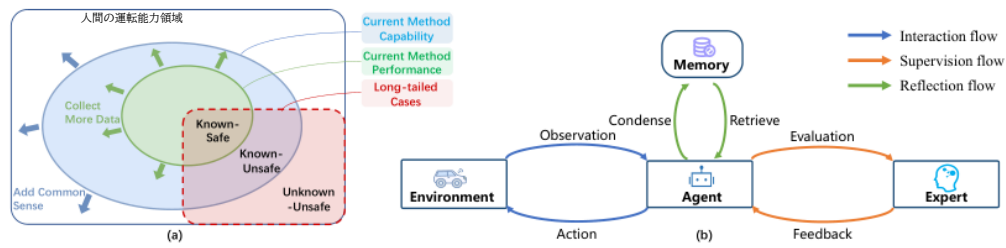
本論文では、大規模言語モデル(LLM)を用いて、人間のようには運転環境を理解する可能性を探り、複雑なシナリオに直面したときに推論、解釈、記憶する能力を分析する。我々は、従来の最適化ベースのモジュール型自律走行(AD)システムは、ロングテールのコーナーケースを扱う際に、固有の性能制限に直面すると主張する。この問題に対処するために、我々は理想的なADシステムが人間のようには運転し、継続的な運転を通じて経験を蓄積し、問題を解決するために常識を使用することを提案する。この目標を達成するために、我々はADシステムに必要な3つの重要な能力を特定する:推論、解釈、記憶。我々は、LLMの理解能力と環境との相互作用能力を示すために、閉ループシステムを構築することで、運転シナリオにLLMを採用することの実現可能性を実証する。我々の広範な実験により、LLMはロングテールのケースを推論し解決する素晴らしい能力を示し、人間のようには自律走行の開発に貴重な洞察を提供することが示された。関連コードは <https://github.com/PJLab-ADG/DriveLikeAHuman> で公開されている。

## 1 Introduction

ストップサインの前に緑色の光を待っている車輪の後ろに座っている場合を想像してください。一方、交通コーンを運ぶピックアップトラックは、交差点を前方に移動している。人間のドライバーとして、あなたの常識的な知識を活用して、これらの交通コーンがピックアップトラックの貨物であり、道路が建設中であることを意味しないことを推論することができる。しかし、人間のドライバーが扱いやすいこれらのシナリオは、既存の多くの自律走行(AD)システムにとってロングテールのコーナーケースである[6, 7, 8]。自動運転の開発者は、突然のブレーキを防ぐために、地面ではなく、車両の交通コーンについてルールを作成したり、より多くのデータを特別に収集したりすることで、このケースに対処することができるが、ノーゴーゾーンを示す地面に遭遇した場合、アルゴリズムは逆のケースで失敗する。ある問題を解決するのは、別の問題を見つけるためだけのようなもので、特に、現実世界では想像を絶するほど無限の稀なケースがある。このため、従来の最適化ベースのモジュール型ADシステムは、本質的に性能のボトルネックに直面していると考えています[6, 17]。

自律走行のストーリーボードを再考し、図1(a)の困難なオープンワールドに苦戦する従来の最適化ベースのADシステムがなぜなのかを明らかにする。最適化理論を基礎としたシステムは、複雑な自律走行タスクをサブタスクの集合に簡単に分割することができる。損失関数を最適化する目的は、複雑なシナリオに直面したとき、局所最適化に捕らわれる傾向があり、そのため汎化能力が制限される。

\* 均等な貢献。姓でアルファベット順にソート。著者名:† Corresponding author.



ルールベースとRLベースのパイプライン図1: (a)人間の運転と既存の自律走行システムとの関係、ルール/RLは、現在のアプローチの限界と、なぜロングテールのケースをすべて解決できないのかを特に強調している。(b) システムCarのスキーマ:  $(x \text{ can}^{z1, y1}, (x \text{ can}^{z1 2, w y1 2}, z^h 2^1, \text{drive} w^1 2), h 2, l 2)$  人間のように。<sup>123</sup> エージェントの行動は、環境を探索し、相互作用し、専門家のフィードバックに基づいて自己反省を行い、最終的に経験を蓄積することができる。AGIベースのパイプライン

より多くのデータ(グラフの緑の矢印)を取り入れることで、現在のモデルに対する<sup>xxxx</sup>間の<sup>xxxxxx</sup>性能<sup>xxxxxx</sup>ギャップ(緑の楕円)と最適化ベース(c)手法の最大容量(青の楕円)を削減できるだけである。これは主に、最適化プロセスがデータ内の支配的なパターンを学習することに重点を置いているためで、しばしば頻度の低いロングテールのコーナーケースを見落としている。常識(青矢印)を取り込まなければ、モデルの能力(青楕円)を促すことはできない。

さらに、連続的なデータ収集の間、常に無尽蔵の未知のロングテールのケースが存在する。このようなロングテールのコーナーケースに対処するのに苦労している現在の解決策と比較すると、人間は常に経験や常識によってスキルと楽しさで解決できる。簡単なアイデアが浮かび上がってくる: 限られた訓練コーパスに適合するのではなく、継続的な運転によって経験を蓄積できる人間のように運転できるようなシステムを作ることは可能なのだろうか?

最近の研究[14, 32, 46, 39, 47]によると、これまでのモジュラーADシステムは、推論、解釈、自己反省のような高度な知能を持たないタスク固有のコーパスで訓練されたインターネットAI[7, 8, 23, 17]とみなすことができると考えている。経験豊富な人間のドライバーのように車を運転できるエージェントを得たいのであれば、Embodied AI [31, 35]の研究からアイデアを借りる必要があると主張する。人間は実環境とのインタラクションから運転することを学ぶとともに、様々なシナリオとそれに対応する操作の記憶を説明し、推論し、凝縮することで、道路感覚を洗練させるフィードバックを得る。さらに、人間のドライバは、その論理的推論能力により、ルールを要約し、より一般的なシナリオに適用するために、常識を利用することができる(帰納的推論)[16]。一方、予測不可能なシナリオを扱うために、これまでの経験を無意識に喚起することができる(演繹的推論)[20]。

人間のように運転するという目標に向けて、我々は必要な3つの能力を特定する: 1) 推論: 特定の運転シナリオが与えられたとき、モデルは常識と経験を介した推論によって意思決定を行うことができるはずである。2) 解釈: エージェントが下した決定は解釈できるはずである。これは、内観の能力と宣言的記憶の存在を示している。3) 記憶: シナリオを推論し解釈した後、記憶メカニズムは以前の経験を記憶し、エージェントが同じような状況に直面したときに、同じような決定をすることを可能にすることが要求される。

以上の3つの性質に基づき、人間が運転する学習のパラダイムを参照し、図1(b)に描かれた運転システムの正準形式を凝縮する。このスキーマには4つのモジュールが含まれる: (1) 環境はエージェントが相互作用の流れによって相互作用できる段階を作り出す、(2) エージェントは環境を認識し、その記憶と専門家の助言からの学習を利用して意思決定を行うことができるドライバを表す、(3) 記憶はエージェントが経験を蓄積し、リフレクションフローを介してそれを使って行動を実行することを可能にする、(4) 専門家はエージェントの訓練について助言を提供し、エージェントが一貫性のない行動をとったときにフィードバックを与え、これが監督フローを形成する。具体的には、普遍的な運転フレームワークとして、環境、エージェント、エキスパートは、それぞれ実世界またはシミュレータ、人間のドライバまたは運転アルゴリズム、シミュレータまたはインストラクタのフィードバックによって表現することができる。

最近の研究に触発され、大規模言語モデル(LLM)は、その顕著な創発能力[40]と、Instruct Following[27]やIn-Context Learning(ICL)[3]のような新しい技術により、人工知能(AGI)[4, 39, 47, 1, 34]の初期バージョンと考えることができる。広範な

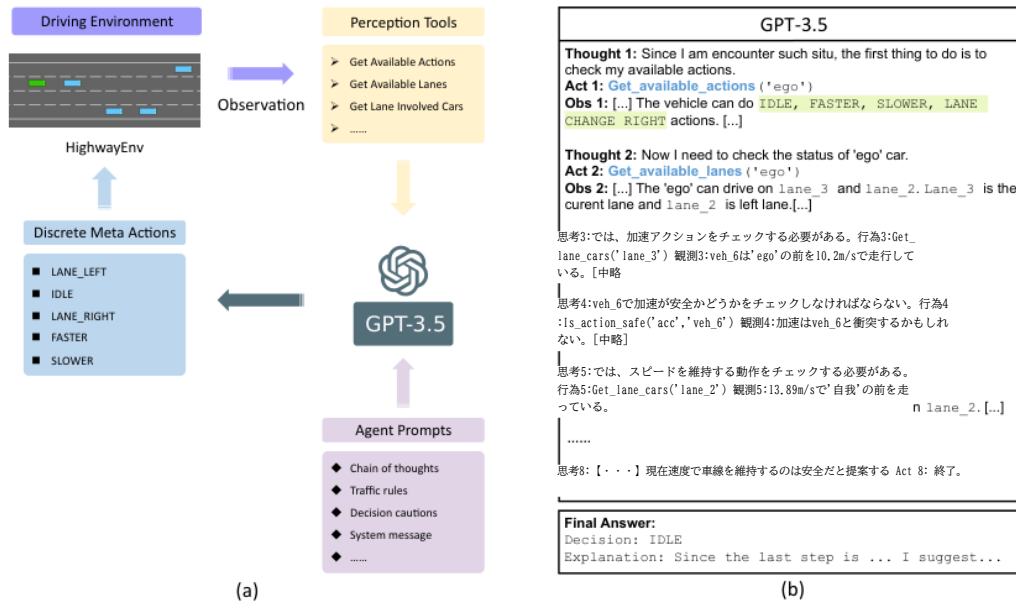


図2:HighwayEnvにおけるGPT-3.5のクローズドループ走行:(a)GPT-3.5はHighwayEnvにおいて知覚ツールを用いて環境を観察し、車両を制御するための意思決定を行い、クローズドループを形成している。(b) GPT-3.5は、思考、行動、観察のサイクルを通じて周囲を知覚しながら、行動を計画し、道具を使用するためにReAct戦略を採用している。

ChatGPT[26]のような最近リリースされたLLMの実験結果は、推論、解釈、暗記[41]の能力を実証している。そこで本稿では、まず人間のような運転交通シーンを理解するLLMの能力を探求し、ロングテールコーナーケースのようなシナリオを扱う際の推論、解釈、記憶のLLMの能力を、一連の定性的な実験を通して分析することを試みる。具体的には、まず、運転シナリオにおけるLLM(GPT-3.5)の理解能力と環境インタラクション能力を実証するために、閉ループシステムを構築する。次に、モジュラーADシステムにとって扱いにくい、人間のドライバにとって扱いやすいいくつかの典型的なロングテールのケースを解くことで、推論と記憶能力を実証する。

本論文の主な貢献は以下の通りである：

1. ロングテールコーナーに直面したときに、既存のADシステムの壊滅的な忘却を防ぐために、自律走行システムを人間のように運転させる方法を深く掘り下げ、人間のように運転する3つの重要な能力に要約する：推論、解釈、記憶。
2. 運転シナリオにLLMを採用することの実現可能性を実証し、模擬運転環境における意思決定能力を活用したのは我々が初めてである。
3. 本研究における広範な実験は、印象的な理解とロングテールのケースを解決する能力を表現している。これらの洞察が、学界や産業界に、人間のような自律走行の発展に貢献することを期待している。

## 2 運転シナリオにおける閉ループ相互作用能力

解釈能力は、LLMが運転環境を理解し、環境との相互作用の基礎を形成し、推論能力、記憶能力を強化することを可能にする。LLMの解釈と環境とのインタラクション能力を検証するため、GPT-3.5を用いてHighwayEnv\*のクローズドループ運転実験を行った。テキストのみの大規模言語モデルであるGPT-3.5はHighwayEnvと直接対話できないため、GPT-3.5の観察と意思決定を支援する知覚ツールとエージェントプロンプトを提供した。

\* <https://github.com/ファラマ財団/ハイウェイエンブ>

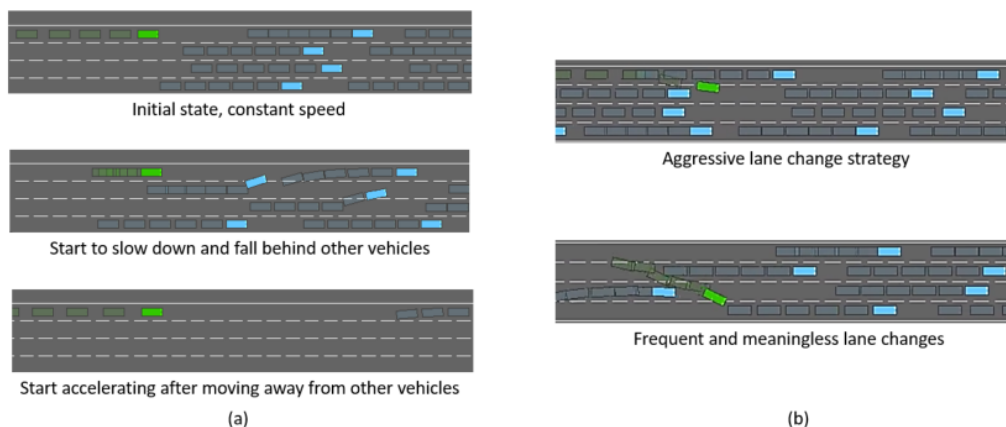


図3:RLベースとサーチベースの手法の運転行動:(a)RLベースのエージェントは、中間ステップを無視して、最終的な報酬を達成することのみに焦点を当てる。これにより、他の車両の後方に落ちるのを遅くしたり、衝突を避けるためにオープンロードを走行したりといった、型破りな行動をとることができる。(b)目的関数を最適化することで、探索に基づく決定を行う。安全を確保しつつ、最大限の効率性を求めて攻撃的な行動をとることができる。

図2に示すように、エージェントプロンプトはGPT-3.5に現在の行動、運転ルール、注意事項に関する情報を提供する。GPT-3.5は、思考、行動、観察のサイクルを通じて、周囲の環境を認識し、分析するために、ReAct戦略[42]を採用している。この情報に基づいて、GPT-3.5はHighwayEnvで車両を決定し制御し、閉ループ走行システムを形成する。

人間と同様に、GPT-3.5は運転中の行動の潜在的な結果を評価し、最も賢明な判断を下すために結果を重み付けする。広く使われている強化学習(RL)ベースや探索ベースのアプローチとは異なり、GPT-3.5はシナリオや行動を解釈するだけでなく、意思決定プロセスを最適化するために常識を利用する。

GPT-3.5は、RLベースのアプローチと比較して、HighwayEnvにおいて、微調整なしで60%以上のゼロショット合格率を達成した。対照的に、RLベースのアプローチは、競争力のある性能を達成するために、多数の反復に大きく依存する。例えば、図3(a)が示すように、衝突に対するペナルティが大きいため、RLベースのエージェントは、衝突を防ぐために、最初に減速して、その後の加速のための広範な空間を作り出すという方針を学習した。RLベースのアプローチは、しばしばこのような予期せぬ解を生成することを示している。

探索に基づくアプローチは、目的関数を最適化することで、その関数に記載されていない未定義部分を見逃して意思決定を行う。図3(b)が示すように、探索型エージェントは、高い走行効率を達成するために積極的な車線変更を行う可能性があり、それによって衝突のリスクが高まる。さらに、サーチベースのアプローチは、他の車両が先行していない場合でも、無意味な車線変更を行う可能性がある。これは、探索型エージェントの場合、安全性を前提とした目的関数において、車線変更と速度維持の優先順位が等しいためと考えられる。その結果、エージェントは行動の1つをランダムに選択する。

要約すると、RLベースのアプローチも探索ベースのアプローチも、常識、シナリオを解釈する能力、長所と短所を天秤にかける能力を欠いているため、人間のように真に考え、駆り立てることはできない。これに対して、GPT-3.5は各行動の結果を説明することができ、プロンプトを提供することで、GPT-3.5を価値志向にし、より人間に近い判断を可能にすることができる。

HighwayEnv環境におけるGPT-3.5の解釈能力とインタラクション能力、およびクローズドループプロセスにおける意思決定の一貫性を示す2つの例を示す。最初のケースは図4に描かれており、一番右の車線lane\_3を走行する緑色のエゴ・カーが関与している。エゴ・カーは、lane\_3を走行する先行車veh4にしばらく追従しているが、左車線lane\_2のveh1はveh4よりも速い速度で走行している。その後、GPT-3.5はReAct処理を開始する。まず、エゴ・カーが現在実行できるアクション(現在の車線の速度の加速、減速、維持、左車線の変更など)を決定する。しかし、エゴ・カーが一番右の車線にいるため、右の車線変更を続けることはできない。次に

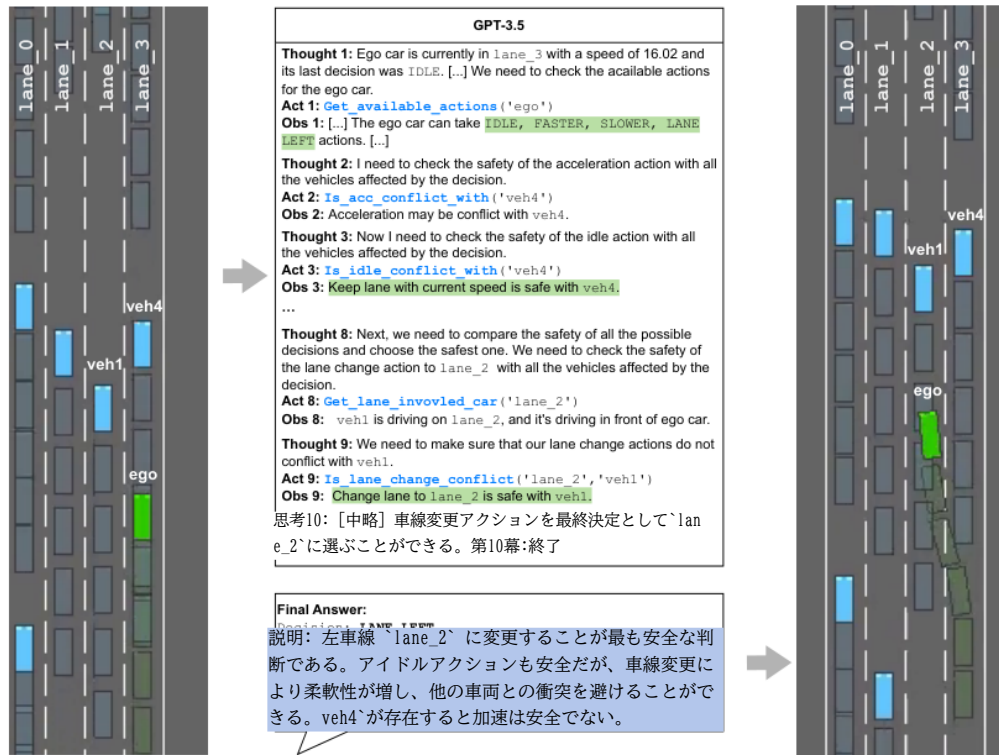


図4:GPT-3.5による車線変更の意思決定プロセス

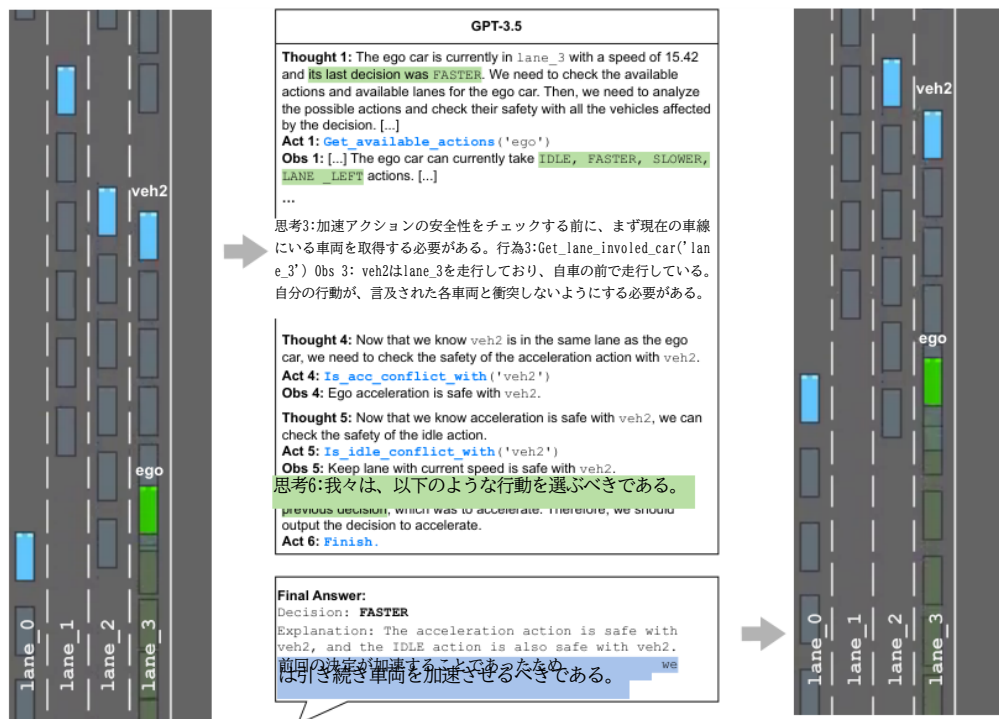


図5:アクセラレーションプロセスにおける意思決定の一貫性

GPT-3.5は、利用可能な各アクションの安全性をチェックする。知覚ツールは、加速が前車veh4と衝突する可能性があり、速度維持が非常に安全であることを示している。左車線変更アクションをチェックするとき、GPT-3.5はまず、lane\_2上のそのようなアクションによって影響を受ける車両を決定し、次に、左車線変更がveh1で安全であることを学習する。この時点で、GPT-3.5は各アクションを検討し、最終的に左車線変更の決定を行った。アイドル時と車線変更時の動作はどちらも安全であるが、車線\_2に変更した方がエゴ車に柔軟性を与えるため、より良い動きであると述べ、成熟した合理的な説明を提供する。veh1の方が高速であることを考慮すると、この決定はより良いパフォーマンスにつながる。

HighwayEnvにおけるクローズドループ走行は、LLMが各時間ステップで安全な意思決定を行うだけでなく、意思決定間の一貫性を要求し、頻繁な加速・減速や無意味な車線変更などの行動を回避する。我々のフレームワークでは、前のフレームからの決定結果と説明は、エージェントプロンプトの一部として含まれ、GPT-3.5に入力される。図5の2番目の例を用いて、GPT-3.5がこのような意思決定の一貫性を持つことを示す。

この例では、緑色のエゴ・カーは一番右の車線にあり、比較的長い距離を保ちながらveh2に従っている。前回の決定では、GPT-3.5は先行車からの距離が遠すぎると判断したため、veh2に追いつくためにスピードアップすることを決定した。ReActプロセスの開始時、GPT-3.5は依然としてGet\_available\_actionツールを使用して、現在の時間ステップで利用可能な4つのアクションをすべて取得している。そして、veh2はまだエゴ・カーの前を走行しており、アイドル動作と加速動作の両方が先行車と安全であることがわかった。GPT-3.5の最終決定は、その最終回答で説明されているように、「前の決定と一致する行動を選択する」ので、加速し続けることである。その結果、エゴ・カーはフロント・ビークルとの距離を縮め、全体的な交通の流れをより助長する。最初の例と比較して、前の決定結果を参照しているため、GPT-3.5によって呼び出されるツールの数と推論コストが大幅に削減されている。

### 3 常識に基づく推論能力

人間のドライバーも、これまでの最適化ベースのADシステムも、基本的な運転スキルを持っているが、両者の根本的な違いは、人間が世界について常識的な理解を持っていることである。常識とは、日常生活から蓄積された、私たちの周りに起こっていることに対する健全で実践的な判断である[11]。運転に貢献する常識は、日常生活のあらゆる側面から導き出すことができる。新しい運転状況が提示されると、人間のドライバーは常識に基づいてシナリオを素早く評価し、合理的な判断を下すことができる。対照的に、従来のADシステムは運転領域での経験があるかもしれないが、常識に欠けるため、そのような状況に対処することができない。

GPT-3.5のようなLLMは、膨大な量の自然言語データで学習され、常識に精通している[2]。これは従来のAD手法とは大きく異なる点であり、LLMが人間のドライバーのように常識的な複雑な運転シナリオを推論する力を与えている。本節では、自律走行システムにおける2つの典型的なロングテールケースを評価する。これは、セクション1の冒頭で説明したように、交通コーンを運ぶピックアップトラックを含む。

図6に示すように、類似しているが異なる2枚の写真をLLMに入力する。最初の写真は、目的地まで移動中にトラックベッドに複数のトラフィックコーンを積んだピックアップトラックを描いている。2つ目のピックアップも、トラックのベッドにコーンが置かれ、しかも周囲に地面にコーンが散らばっているピックアップを描いている。GPT-3.5は画像を含むマルチモーダル入力を処理する機能がないため、画像処理フロントエンドとしてLLaMA-Adapter v2 [15]を採用した。LLaMA-Adapterに写真の詳細をできるだけ説明するよう指示する。この記述を観察として利用し、GPT-3.5にシナリオが潜在的に危険であるかどうかを評価し、トラックの後続であると想定される自転車について判断してもらう。

図6(a)に描かれた最初のケースでは、LLaMA-Adapterは、写真のピックアップトラックが複数の交通コーンを搭載していることを識別し、目的地まで輸送している可能性があるかと推測した。これらの観察結果から、GPT-3.5は運転シナリオの解析に成功した。GPT-3.5は、交通コーンの存在に惑わされる代わりに、貨物を目的地まで輸送するトラックはよくあることだという常識に基づき、このシナリオは危険ではないと判断した。GPT-3.5は、エゴ・カーの運転者に減速する必要がないことを助言し、不必要な減速は交通の流れにとって危険である可能性があることに注意を促した。





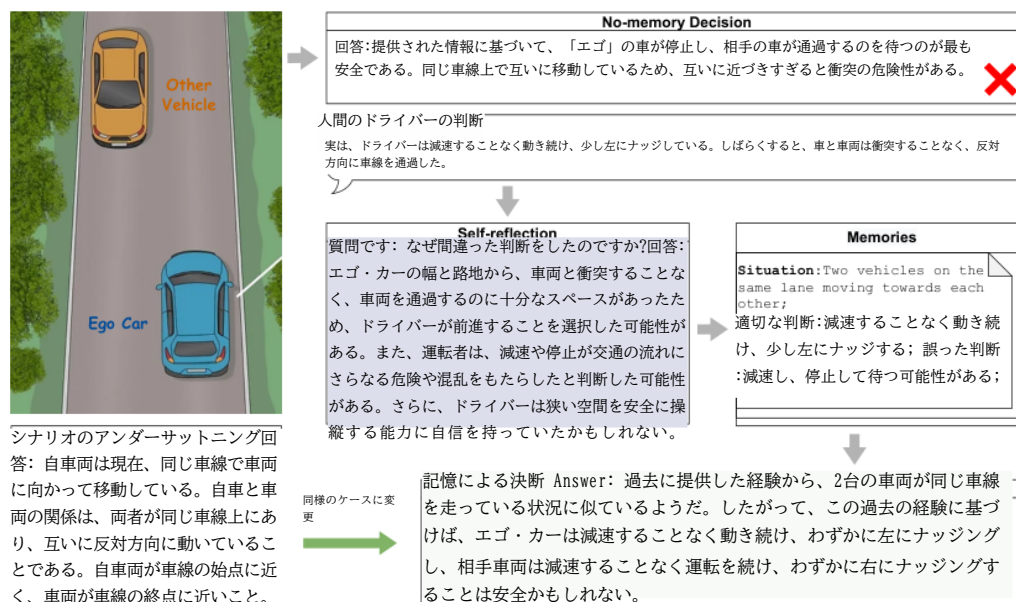


Figure 7: An example of the self-reflection and memorization ability of LLM

図1(b)に示すように、我々の提案するLLMベースのアプローチは、別のメモリモジュールを使用することで、人間に忠実に追従する。メモリモジュールは、“専門家から逸脱した”決定シナリオのみを記録する。専門家は、LLMの決定に対する開発者の評価、または実世界における人間のドライバーの決定のグランドトゥールズのいずれかである。専門家からのフィードバックが得られたら、LLMは自己反省プロセスを経て、なぜその決定が専門家の判断から外れるのかを判断する。そして、トラフィック状況を意思決定シナリオにまとめ、適切な意思決定とともに、新しいメモリエントリとしてメモリプールに追加する。次に同様のケースに遭遇したとき、LLMはこのメモリエントリを素早く取り出して参照し、十分な情報を得た上で決定することができる。

図7に暗記処理の例を示す。このシナリオでは、青いエゴの車と、車の幅の2倍より少し広い狭い車線内を反対方向に移動する黄色い車との遭遇が行われる。GPT-3.5に入力するためにシーンを構造化されたテキストに変換した後、モデルは車両の状態、向き、目的地を含むシーンをよく理解していることがわかった。しかし、シナリオの決定を依頼したところ、GPT-3.5は、エゴ・カーが停止し、相手車両が先に通過するのを待つべきだという、安全だが過度に慎重なアドバイスをした。LLMの性能を向上させるために、専門家は人間のドライバーがどのように状況に対処するかについて実践的なアドバイスをを行う。これは、車の動きを維持し、少し左にナッジすることを含む。そしてLLMは、両車両が通過するのに十分なスペースがあり、減速すると交通の流れが混乱することを認識する。状況を「同じ車線にいる2台の車両が互いに移動している」とまとめ、適切な判断とともに記憶を記録する。これらの記憶を使って、狭い路地で異なる速度と位置で2台の車が合流するという別のシナリオを入力し、LLMに決定を求めた。LLMは、これが「同じ車線にいる2台の車が互いに向かっている」という決定シナリオの別の変形に過ぎないことを認識し、減速して待つのではなく、エゴの車が動き続けることは安全であると助言している。

記憶能力は、経験を積むための運転事例を継続的に収集し、既存の記憶を検索することで意思決定を支援するもので、LLMに自律走行分野における継続的な学習能力を付与する。さらに、これは類似シナリオにおけるLLMの決定コストを大幅に削減し、実用的な性能を向上させる。



自動運転の自律性。自律走行車は、モジュール式[37, 12]とエンドツーエンド[5, 17, 33, 36]の2つの主要なパラダイムを包含する。モジュラーアプローチは、知覚[23, 24, 43]、計画[21, 44]、制御[30, 19]などの様々なサブタスクを処理する、相互接続されたコンポーネントのスタックを含む。このアーキテクチャは、モジュール性や汎用性といった魅力的な機能を提供する。しかし、パイプラインのチューニングやエラー伝播の管理には課題がある。一方、エンドツーエンドの自律性は、センサー入力をプランナーやコントローラコマンドに直接マッピングする。これらの方法は一般的に開発が容易であるが、解釈可能性に欠けるため、エラーの診断、安全性の確保、交通ルールの組み込みが困難である。とはいえ、エンドツーエンドで学習可能なパイプラインの自律性における最近の進歩は、両方のパラダイムの長所を組み合わせることで、有望な結果を示している[17, 5]。自動運転のこれら2つのパラダイムは大きな進歩を遂げているにもかかわらず、実世界の環境で発生するロングテールデータや分布外のシナリオを扱う際に、脆弱になることがよく観察されており[22]、セーフティクリティカルな自律走行に課題を投げかけている。大規模言語モデルによる高度なタスク大規模言語モデル(Large Language Models: LLM)の成功は、機械が人間の知識をどの程度学習できるかを示すものであり、間違いなくエキサイティングである。LLMにおける最近の取り組みは、ゼロショットプロンプトと複雑な推論[2, 25, 9, 27, 10]、具現化エージェント研究[39, 47, 38, 13, 42]、主要な輸送問題への対処[45]において、印象的な性能を示している。PaLM-E[13]は、マルチモーダルなプロンプトをサポートするために、事前に訓練されたLLMを適応させるために、微調整技術を採用している。Reflexion[34]は、LLMを使用して推論トレースとタスク固有の行動の両方を生成するために、思考連鎖プロンプト[27]でエージェントの推論能力をさらに強化するために自己反省を組み込んでいる。VOYAGER [39]は、LLMに基づくプロンプト機構、スキルライブラリ、自己検証を用いた生涯学習を提示している。これら3つのモジュールは、エージェントのより複雑な行動の開発を強化することを目的としている。生成エージェント[29]は、エージェントの経験の完全な記録を保存するためにLLMを採用する。時間の経過とともに、これらの記憶はより高次の反射に合成され、動的に取り出され、行動を計画する。Instruct2Act[18]は、ロボット操作タスクのために、マルチモーダル命令を逐次アクションにマッピングするために、大規模言語モデルを利用するフレームワークを導入している。

## 6 Conclusion

本論文では、人間のように運転できるシステムを構築するという我々のアイデアを提示する。我々は、大域的最適化の壊滅的な忘却により、ロングテールのコーナーケースを扱う場合、これまでの最適化ベースの自律走行システムには限界があると考えている。そこで、ADシステムが不完全性を打ち負かすために必要な能力として、(1)推論、(2)解釈、(3)記憶の3つを挙げる。そして、この3つの信条に倣って、人間の運転学習のプロセスを模倣した新しいパラダイムを設計する。最後に、人工知能の期待に応え、GPT-3.5をLLMのテストベッドとして使用し、交通シナリオを理解する素晴らしい能力を示すことを試みる。予備的な研究として、我々はLLMを運転エージェントとして使用するのではなく、この技術を採用することの利点と機会を強調するために、閉ループ運転におけるLLMの可能性の表面をかすめただけである。我々の願望は、この研究が学界と産業界の双方にとって、人間のように運転できるAGIベースの自律走行システムを革新し、構築するための触媒として役立つことである。

## References

- [1] Michael Ahn, Anthony Brohan, Noah Brown, Yevgen Chebotar, Omar Cortes, Byron David, Chelsea Finn, Chuyuan Fu, Keerthana Gopalakrishnan, Karol Hausman, et al. Do as i can, not as i say: Grounding language in robotic affordances. *arXiv preprint arXiv:2204.01691*, 2022.
- [2] 寧天、韓向平、孫樂、林宏宇、呂耀傑、何博文。Chatgptは知識はあるが経験の浅いソルバーである：大規模言語モデルにおけるコモンセンス問題の調査, 2023.
- [3] Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33:1877–1901, 2020.
- [4] Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, Johannes Gehrkke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuanzhi Li, Scott Lundberg, et al. Sparks of artificial general intelligence: Early experiments with gpt-4. *arXiv preprint arXiv:2303.12712*, 2023.

- [5] Sergio Casas, Abbas Sadat, and Raquel Urtasun. Mp3: A unified model to map, perceive, predict and plan. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14403–14412, 2021.
- [6] Long Chen, Yuchen Li, Chao Huang, Bai Li, Yang Xing, Daxin Tian, Li Li, Zhongxu Hu, Xiaoxiang Na, Zixuan Li, et al. Milestones in autonomous driving and intelligent vehicles: Survey of surveys. *IEEE Transactions on Intelligent Vehicles*, 8(2):1046–1056, 2022.
- [7] Long Chen, Yuchen Li, Chao Huang, Yang Xing, Daxin Tian, Li Li, Zhongxu Hu, Siyu Teng, Chen Lv, Jinjun Wang, et al. Milestones in autonomous driving and intelligent vehicles—part 1: Control, computing system design, communication, hd map, testing, and human behaviors. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2023.
- [8] Long Chen, Siyu Teng, Bai Li, Xiaoxiang Na, Yuchen Li, Zixuan Li, Jinjun Wang, Dongpu Cao, Nanning Zheng, and Fei-Yue Wang. Milestones in autonomous driving and intelligent vehicles—part ii: Perception and planning. *IEEE Transactions on Systems, Man, and Cybernetics: Systems*, 2023.
- [9] Aakanksha Chowdhery, Sharan Narang, Jacob Devlin, Maarten Bosma, Gaurav Mishra, Adam Roberts, Paul Barham, Hyung Won Chung, Charles Sutton, Sebastian Gehrmann, et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.
- [10] Hyung Won Chung, Le Hou, Shayne Longpre, Barret Zoph, Yi Tay, William Fedus, Eric Li, Xuezhi Wang, Mostafa Dehghani, Siddhartha Brahma, et al. Scaling instruction-finetuned language models. *arXiv preprint arXiv:2210.11416*, 2022.
- [11] Common sense. Common sense — Wikipedia, the free encyclopedia, 2023. [Online; accessed 3-July-2023].
- [12] Jonathan Daudelin, Gangyuan Jing, Tarik Tosun, Mark Yim, Hadas Kress-Gazit, and Mark Campbell. An integrated system for perception-driven autonomy with modular robots. *Science Robotics*, 3(23):eaat4983, 2018.
- [13] Danny Driess, Fei Xia, Mehdi SM Sajjadi, Corey Lynch, Aakanksha Chowdhery, Brian Ichter, Ayzaan Wahid, Jonathan Tompson, Quan Vuong, Tianhe Yu, et al. Palm-e: An embodied multimodal language model. *arXiv preprint arXiv:2303.03378*, 2023.
- [14] Jiafei Duan, Samson Yu, Hui Li Tan, Hongyuan Zhu, and Cheston Tan. A survey of embodied ai: From simulators to research tasks. *IEEE Transactions on Emerging Topics in Computational Intelligence*, 6(2):230–244, 2022.
- [15] Peng Gao, Jiaming Han, Renrui Zhang, Ziyi Lin, Shijie Geng, Aojun Zhou, Wei Zhang, Pan Lu, Conghui He, Xiangyu Yue, Hongsheng Li, and Yu Qiao. Llama-adapter v2: Parameter-efficient visual instruction model, 2023.
- [16] Brett K Hayes, Evan Heit, and Haruka Swendsen. Inductive reasoning. *Wiley interdisciplinary reviews: Cognitive science*, 1(2):278–292, 2010.
- [17] Yihan Hu, Jiazhi Yang, Li Chen, Keyu Li, Chonghao Sima, Xizhou Zhu, Siqi Chai, Senyao Du, Tianwei Lin, Wenhai Wang, et al. Planning-oriented autonomous driving. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17853–17862, 2023.
- [18] Siyuan Huang, Zhengkai Jiang, Hao Dong, Yu Qiao, Peng Gao, and Hongsheng Li. Instruct2act: Mapping multi-modality instructions to robotic actions with large language model, 2023.
- [19] Michael A Johnson and Mohammad H Moradi. *PID control*. Springer, 2005.
- [20] Philip N Johnson-Laird. Deductive reasoning. *Annual review of psychology*, 50(1):109–135, 1999.
- [21] Alonzo Kelly and Bryan Nagy. Reactive nonholonomic trajectory generation via parametric optimal control. *The International Journal of Robotics Research*, 22(7-8):583–601, 2003.
- [22] Lingdong Kong, Youquan Liu, Xin Li, Runnan Chen, Wenwei Zhang, Jiawei Ren, Liang Pan, Kai Chen, and Ziwei Liu. Robo3d: Towards robust and reliable 3d perception against corruptions. *arXiv preprint arXiv:2303.17597*, 2023.
- [23] Xin Li, Tao Ma, Yuenan Hou, Botian Shi, Yuchen Yang, Youquan Liu, Xingjiao Wu, Qin Chen, Yikang Li, Yu Qiao, et al. Logonet: Towards accurate 3d object detection with local-to-global cross-modal fusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 17524–17534, 2023.
- [24] Zhiqi Li, Wenhai Wang, Hongyang Li, Enze Xie, Chonghao Sima, Tong Lu, Yu Qiao, and Jifeng Dai. Bevformer: Learning bird’s-eye-view representation from multi-camera images via spatiotemporal transformers. In *European conference on computer vision*, pages 1–18. Springer, 2022.
- [25] John J Nay. Law informs code: A legal informatics approach to aligning artificial intelligence with humans. *Nw. J. Tech. & Intell. Prop.*, 20:309, 2022.
- [26] OpenAI. Introducing chatgpt. <https://openai.com/blog/chatgpt/>, 2023.

- [27] Long Ouyang, Jeffrey Wu, Xu Jiang, Diogo Almeida, Carroll Wainwright, Pamela Mishkin, Chong Zhang, Sandhini Agarwal, Katarina Slama, Alex Ray, et al. Training language models to follow instructions with human feedback. *Advances in Neural Information Processing Systems*, 35:27730–27744, 2022.
- [28] German I Parisi, Ronald Kemker, Jose L Part, Christopher Kanan, and Stefan Wermter. Continual lifelong learning with neural networks: A review. *Neural networks*, 113:54–71, 2019.
- [29] Joon Sung Park, Joseph C O’Brien, Carrie J Cai, Meredith Ringel Morris, Percy Liang, and Michael S Bernstein. Generative agents: Interactive simulacra of human behavior. *arXiv preprint arXiv:2304.03442*, 2023.
- [30] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-real transfer of robotic control with dynamics randomization. In *2018 IEEE international conference on robotics and automation (ICRA)*, pages 3803–3810. IEEE, 2018.
- [31] Rolf Pfeifer and Fumiya Iida. Embodied artificial intelligence: Trends and challenges. In *Embodied Artificial Intelligence: International Seminar, Dagstuhl Castle, Germany, July 7-11, 2003. Revised Papers*, pages 1–26. Springer, 2004.
- [32] Santhosh K Ramakrishnan, Dinesh Jayaraman, and Kristen Grauman. An exploration of embodied visual exploration. *International Journal of Computer Vision*, 129:1616–1649, 2021.
- [33] Abbas Sadat, Sergio Casas, Mengye Ren, Xinyu Wu, Pranaab Dhawan, and Raquel Urtasun. Perceive, predict, and plan: Safe motion planning through interpretable semantic representations. In *Computer Vision—ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XXIII 16*, pages 414–430. Springer, 2020.
- [34] Noah Shinn, Beck Labash, and Ashwin Gopinath. Reflexion: an autonomous agent with dynamic memory and self-reflection. *arXiv preprint arXiv:2303.11366*, 2023.
- [35] Linda Smith and Michael Gasser. The development of embodied cognition: Six lessons from babies. *Artificial life*, 11(1-2):13–29, 2005.
- [36] Ardi Tampuu, Tambet Matiisen, Maksym Semikin, Dmytro Fishman, and Naveed Muhammad. A survey of end-to-end driving: Architectures and training methods. *IEEE Transactions on Neural Networks and Learning Systems*, 33(4):1364–1384, 2020.
- [37] Sebastian Thrun, Mike Montemerlo, Hendrik Dahlkamp, David Stavens, Andrei Aron, James Diebel, Philip Fong, John Gale, Morgan Halpenny, Gabriel Hoffmann, et al. Stanley: The robot that won the darpa grand challenge. *Journal of field Robotics*, 23(9):661–692, 2006.
- [38] Sai Vemprala, Rogerio Bonatti, Arthur Buckner, and Ashish Kapoor. Chatgpt for robotics: Design principles and model abilities. *Microsoft Auton. Syst. Robot. Res.*, 2:20, 2023.
- [39] Guanzhi Wang, Yuqi Xie, Yunfan Jiang, Ajay Mandlekar, Chaowei Xiao, Yuke Zhu, Linxi Fan, and Anima Anandkumar. Voyager: An open-ended embodied agent with large language models. *arXiv preprint arXiv:2305.16291*, 2023.
- [40] Jason Wei, Yi Tay, Rishi Bommasani, Colin Raffel, Barret Zoph, Sebastian Borgeaud, Dani Yogatama, Maarten Bosma, Denny Zhou, Donald Metzler, et al. Emergent abilities of large language models. *arXiv preprint arXiv:2206.07682*, 2022.
- [41] Lilian Weng. Llm-powered autonomous agents. *lilianweng.github.io*, Jun 2023.
- [42] Shunyu Yao, Jeffrey Zhao, Dian Yu, Nan Du, Izhak Shafran, Karthik Narasimhan, and Yuan Cao. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.
- [43] Tianwei Yin, Xingyi Zhou, and Philipp Krahenbuhl. Center-based 3d object detection and tracking. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 11784–11793, 2021.
- [44] Chris Zhang, Runsheng Guo, Wenyan Zeng, Yuwen Xiong, Binbin Dai, Rui Hu, Mengye Ren, and Raquel Urtasun. Rethinking closed-loop training for autonomous driving. In *European Conference on Computer Vision*, pages 264–282. Springer, 2022.
- [45] Ou Zheng, Mohamed Abdel-Aty, Dongdong Wang, Zijin Wang, and Shengxuan Ding. Chatgpt is on the horizon: Could a large language model be all we need for intelligent transportation?, 2023.
- [46] Hao Zhu, Raghav Kapoor, So Yeon Min, Winson Han, Jiatai Li, Kaiwen Geng, Graham Neubig, Yonatan Bisk, Aniruddha Kembhavi, and Luca Weihs. Excalibur: Encouraging and evaluating embodied exploration. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 14931–14942, 2023.
- [47] Xizhou Zhu, Yuntao Chen, Hao Tian, Chenxin Tao, Weijie Su, Chenyu Yang, Gao Huang, Bin Li, Lewei Lu, Xiaogang Wang, et al. Ghost in the minecraft: Generally capable agents for open-world environments via large language models with text-based knowledge and memory. *arXiv preprint arXiv:2305.17144*, 2023.