

MITEE: TOWARDS HUMAN-LIKE AUTONOMOUS DRIVING VIA MULTI-MODAL TRAJECTORY GENERATION AND EVALUATION

Anonymous authors

Paper under double-blind review

ABSTRACT

End-to-end autonomous driving systems integrate perception, decision-making, and control, enhancing vehicle adaptability by transforming sensory inputs directly into driving actions. Traditionally reliant on imitation learning, these systems face challenges like causal confusion and distribution shifts, limiting their effectiveness in diverse conditions. In this paper, we introduce MITEE, the first Human-Like-Centric end-to-end system, which employs a Sparse Perception Module, a Diffusion-Based Motion Planner using Denoising Diffusion Probabilistic Models (DDPM), and a GPT-4o guided Trajectory Scoring System. This innovative framework addresses the limitations of previous systems by capturing the variability in human driving and aligning trajectory evaluations with human judgment. MITEE enhances trajectory prediction robustness and adapts scoring based on real-time contexts, setting a new standard for autonomous vehicle systems with its advanced, human-centric approach. We will open-source the code.

1 INTRODUCTION

End-to-end autonomous driving systems represent a significant shift from traditional modular approaches, where discrete components separately handle perception, decision-making, and control tasks. This paradigm shift towards integrated processing promises more fluid and adaptive vehicle behaviour by directly mapping sensory inputs to driving actions. Historically, such systems have heavily relied on imitation learning techniques, where models are trained to emulate human driving patterns from large-scale driving datasets. These methods aim to simplify the complex multi-stage processing pipeline, reducing latency and potentially increasing the robustness of the decision-making process under varied driving conditions.

Unfortunately, while imitation learning provides a straightforward framework for teaching vehicles to drive, it introduces inherent limitations in trajectory prediction that often manifest as causal confusion and distribution shifts. Causal confusion occurs when models fail to discern the actual causal relationships from correlations within the data, leading to decisions that may not generalize well across different driving scenarios. Moreover, distribution shifts between the training environments and real-world conditions can degrade the model’s performance unpredictably. Existing trajectory scoring systems, which typically employ rule-based, prediction-based, or hybrid methods, also struggle to achieve human-like driving adaptability. These systems often fail to consider the nuanced and stochastic nature of human driving, leading to evaluations that do not capture the full spectrum of plausible and safe driving behaviours.

Our key observation is that diffusion models, particularly Denoising Diffusion Probabilistic Models (DDPM), can effectively address the aforementioned challenges in imitation learning. These models excel in generating diverse and realistic outputs, which makes them particularly suited for predicting multi-modal trajectories that reflect the true variability in human driving. Furthermore, the advent of advanced language models like GPT-4 offers unprecedented opportunities to enhance decision-making processes. By integrating these models, we can develop a trajectory scoring system that not only evaluates the safety and feasibility of predicted trajectories but also aligns closely with human judgment and contextual appropriateness. This approach not only mitigates the issues of tra-

ditional scoring systems but also paves the way for more nuanced and human-like decision-making in autonomous driving.

To this end, we introduce MITEE, the first Human-Like-Centric end-to-end autonomous driving system, as illustrated in Fig. 1b. MITEE consists of a sparse perception module, a diffusion-based motion planner, and a GPT-4o guided trajectory scoring system. Utilizing decoupled instance features and geometric anchors for a complete representation of each entity (be it a dynamic road agent or a static map element), the Sparse Perception module integrates detection, tracking, and online mapping into a unified framework with a symmetric architecture, achieving a fully sparse scene representation. The Motion Planner leverages DDPM to generate multi-modal trajectories by considering both the ego vehicle and surrounding agent instances processed through sparse perception, ensuring a thorough and adaptable planning approach. Subsequently, the Trajectory Scoring system employs a rule-based method for initial trajectory evaluation, then dynamically incorporates GPT-4o to refine scores based on the prevailing environmental context, promoting human-like decision-making. This layered strategy meticulously selects the most suitable trajectory, prioritizing both safety and adaptability.

The main contribution of our work are summarized as follows:

- **Diffusion-based Motion Planner.**
- **Plug-and-Play Trajectory Scoring Module.**
- **Excellent Planning Results in Closed-loop Benchmarks.**

2 GENERAL FORMATTING INSTRUCTIONS

The text must be confined within a rectangle 5.5 inches (33 picas) wide and 9 inches (54 picas) long. The left margin is 1.5 inch (9 picas). Use 10 point type with a vertical spacing of 11 points. Times New Roman is the preferred typeface throughout. Paragraphs are separated by 1/2 line space, with no indentation.

Paper title is 17 point, in small caps and left-aligned. All pages should start at 1 inch (6 picas) from the top of the page.

Authors' names are set in boldface, and each name is placed above its corresponding address. The lead author's name is to be listed first, and the co-authors' names are set to follow. Authors sharing the same address can be on the same line.

Please pay special attention to the instructions in section 4 regarding figures, tables, acknowledgments, and references.

There will be a strict upper limit of 10 pages for the main text of the initial submission, with unlimited additional pages for citations.

3 HEADINGS: FIRST LEVEL

First level headings are in small caps, flush left and in point size 12. One line space before the first level heading and 1/2 line space after the first level heading.

3.1 HEADINGS: SECOND LEVEL

Second level headings are in small caps, flush left and in point size 10. One line space before the second level heading and 1/2 line space after the second level heading.

3.1.1 HEADINGS: THIRD LEVEL

Third level headings are in small caps, flush left and in point size 10. One line space before the third level heading and 1/2 line space after the third level heading.

4 CITATIONS, FIGURES, TABLES, REFERENCES

These instructions apply to everyone, regardless of the formatter being used.

4.1 CITATIONS WITHIN THE TEXT

Citations within the text should be based on the `natbib` package and include the authors' last names and year (with the "et al." construct for more than two authors). When the authors or the publication are included in the sentence, the citation should not be in parenthesis using `\citet{}` (as in "See Hinton et al. (2006) for more information.>"). Otherwise, the citation should be in parenthesis using `\citep{}` (as in "Deep learning shows promise to make progress towards AI (Bengio & LeCun, 2007).>").

The corresponding references are to be listed in alphabetical order of authors, in the REFERENCES section. As to the format of the references themselves, any style is acceptable as long as it is used consistently.

4.2 FOOTNOTES

Indicate footnotes with a number¹ in the text. Place the footnotes at the bottom of the page on which they appear. Precede the footnote with a horizontal rule of 2 inches (12 picas).²

4.3 FIGURES

All artwork must be neat, clean, and legible. Lines should be dark enough for purposes of reproduction; art work should not be hand-drawn. The figure number and caption always appear after the figure. Place one line space before the figure caption, and one line space after the figure. The figure caption is lower case (except for first word and proper nouns); figures are numbered consecutively.

Make sure the figure caption does not get separated from the figure. Leave sufficient space to avoid splitting the figure and figure caption.

You may use color figures. However, it is best for the figure captions and the paper body to make sense if the paper is printed either in black/white or in color.

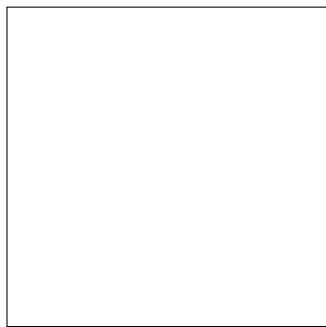


Figure 1: Sample figure caption.

4.4 TABLES

All tables must be centered, neat, clean and legible. Do not use hand-drawn tables. The table number and title always appear before the table. See Table 1.

¹Sample of the first footnote

²Sample of the second footnote

Table 1: Sample table title

PART	DESCRIPTION
Dendrite	Input terminal
Axon	Output terminal
Soma	Cell body (contains cell nucleus)

Place one line space before the table title, one line space after the table title, and one line space after the table. The table title must be lower case (except for first word and proper nouns); tables are numbered consecutively.

5 DEFAULT NOTATION

In an attempt to encourage standardized notation, we have included the notation file from the textbook, *Deep Learning* Goodfellow et al. (2016) available at https://github.com/goodfeli/dlbook_notation/. Use of this style is not required and can be disabled by commenting out `math_commands.tex`.

Numbers and Arrays

a	A scalar (integer or real)
\mathbf{a}	A vector
\mathbf{A}	A matrix
\mathbf{A}	A tensor
\mathbf{I}_n	Identity matrix with n rows and n columns
\mathbf{I}	Identity matrix with dimensionality implied by context
$\mathbf{e}^{(i)}$	Standard basis vector $[0, \dots, 0, 1, 0, \dots, 0]$ with a 1 at position i
$\text{diag}(\mathbf{a})$	A square, diagonal matrix with diagonal entries given by \mathbf{a}
\mathbf{a}	A scalar random variable
\mathbf{a}	A vector-valued random variable
\mathbf{A}	A matrix-valued random variable

Sets and Graphs

\mathbb{A}	A set
\mathbb{R}	The set of real numbers
$\{0, 1\}$	The set containing 0 and 1
$\{0, 1, \dots, n\}$	The set of all integers between 0 and n
$[a, b]$	The real interval including a and b
$(a, b]$	The real interval excluding a but including b
$\mathbb{A} \setminus \mathbb{B}$	Set subtraction, i.e., the set containing the elements of \mathbb{A} that are not in \mathbb{B}
\mathcal{G}	A graph
$\text{Pa}_{\mathcal{G}}(\mathbf{x}_i)$	The parents of \mathbf{x}_i in \mathcal{G}

Indexing

216	a_i	Element i of vector \mathbf{a} , with indexing starting at 1
217	\mathbf{a}_{-i}	All elements of vector \mathbf{a} except for element i
218	$A_{i,j}$	Element i, j of matrix \mathbf{A}
219	$\mathbf{A}_{i,:}$	Row i of matrix \mathbf{A}
220	$\mathbf{A}_{:,i}$	Column i of matrix \mathbf{A}
221	$\mathbf{A}_{i,j,k}$	Element (i, j, k) of a 3-D tensor \mathbf{A}
222	$\mathbf{A}_{:,:,i}$	2-D slice of a 3-D tensor
223	\mathbf{a}_i	Element i of the random vector \mathbf{a}
224		
225		Calculus
226	$\frac{dy}{dx}$	Derivative of y with respect to x
227	$\frac{\partial y}{\partial x}$	Partial derivative of y with respect to x
228	$\nabla_{\mathbf{x}} y$	Gradient of y with respect to \mathbf{x}
229	$\nabla_{\mathbf{X}} y$	Matrix derivatives of y with respect to \mathbf{X}
230	$\nabla_{\mathbf{X}} y$	Tensor containing derivatives of y with respect to \mathbf{X}
231	$\frac{\partial f}{\partial \mathbf{x}}$	Jacobian matrix $\mathbf{J} \in \mathbb{R}^{m \times n}$ of $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$
232	$\nabla_{\mathbf{x}}^2 f(\mathbf{x})$ or $\mathbf{H}(f)(\mathbf{x})$	The Hessian matrix of f at input point \mathbf{x}
233	$\int f(\mathbf{x}) d\mathbf{x}$	Definite integral over the entire domain of \mathbf{x}
234	$\int_{\mathbb{S}} f(\mathbf{x}) d\mathbf{x}$	Definite integral with respect to \mathbf{x} over the set \mathbb{S}
235		
236		Probability and Information Theory
237	$P(\mathbf{a})$	A probability distribution over a discrete variable
238	$p(\mathbf{a})$	A probability distribution over a continuous variable, or over a variable whose type has not been specified
239	$\mathbf{a} \sim P$	Random variable \mathbf{a} has distribution P
240	$\mathbb{E}_{\mathbf{x} \sim P}[f(\mathbf{x})]$ or $\mathbb{E}f(\mathbf{x})$	Expectation of $f(\mathbf{x})$ with respect to $P(\mathbf{x})$
241	$\text{Var}(f(\mathbf{x}))$	Variance of $f(\mathbf{x})$ under $P(\mathbf{x})$
242	$\text{Cov}(f(\mathbf{x}), g(\mathbf{x}))$	Covariance of $f(\mathbf{x})$ and $g(\mathbf{x})$ under $P(\mathbf{x})$
243	$H(\mathbf{x})$	Shannon entropy of the random variable \mathbf{x}
244	$D_{\text{KL}}(P \ Q)$	Kullback-Leibler divergence of P and Q
245	$\mathcal{N}(\mathbf{x}; \boldsymbol{\mu}, \boldsymbol{\Sigma})$	Gaussian distribution over \mathbf{x} with mean $\boldsymbol{\mu}$ and covariance $\boldsymbol{\Sigma}$
246		
247		Functions
248		
249		
250		
251		
252		
253		
254		
255		
256		
257		
258		
259		
260		
261		
262		
263		
264		
265		
266		
267		
268		
269		

270	$f : \mathbb{A} \rightarrow \mathbb{B}$	The function f with domain \mathbb{A} and range \mathbb{B}
271		
272	$f \circ g$	Composition of the functions f and g
273	$f(x; \theta)$	A function of x parametrized by θ . (Sometimes we write
274		$f(x)$ and omit the argument θ to lighten notation)
275	$\log x$	Natural logarithm of x
276		
277	$\sigma(x)$	Logistic sigmoid, $\frac{1}{1 + \exp(-x)}$
278		
279	$\zeta(x)$	Softplus, $\log(1 + \exp(x))$
280	$\ \mathbf{x}\ _p$	L^p norm of \mathbf{x}
281	$\ \mathbf{x}\ $	L^2 norm of \mathbf{x}
282		
283	x^+	Positive part of x , i.e., $\max(0, x)$
284	$\mathbf{1}_{\text{condition}}$	is 1 if the condition is true, 0 otherwise
285		

286 6 FINAL INSTRUCTIONS

287
288
289
290 Do not change any aspects of the formatting parameters in the style files. In particular, do not modify
291 the width or length of the rectangle the text should fit into, and do not change font sizes (except
292 perhaps in the REFERENCES section; see below). Please note that pages should be numbered.

293 7 PREPARING POSTSCRIPT OR PDF FILES

294
295
296 Please prepare PostScript or PDF files with paper size “US Letter”, and not, for example, “A4”. The
297 -t letter option on dvips will produce US Letter files.

298
299 Consider directly generating PDF files using pdf_lat_ex (especially if you are a MiKTeX user).
300 PDF figures must be substituted for EPS figures, however.

301 Otherwise, please generate your PostScript and PDF files with the following commands:

302
303 `dvips mypaper.dvi -t letter -Ppdf -G0 -o mypaper.ps`
304 `ps2pdf mypaper.ps mypaper.pdf`

305 7.1 MARGINS IN L^AT_EX

306
307
308 Most of the margin problems come from figures positioned by hand using `\special` or other
309 commands. We suggest using the command `\includegraphics` from the `graphicx` package.
310 Always specify the figure width as a multiple of the line width as in the example below using `.eps`
311 `graphics`

312
313 `\usepackage[dvips]{graphicx} ...`
314 `\includegraphics[width=0.8\linewidth]{myfile.eps}`

315
316 or

317
318 `\usepackage[pdftex]{graphicx} ...`
319 `\includegraphics[width=0.8\linewidth]{myfile.pdf}`

320 for `.pdf` graphics. See section 4.4 in the graphics bundle documentation (<http://www.ctan.org/tex-archive/macros/latex/required/graphics/grfguide.ps>)

321
322 A number of width problems arise when LaTeX cannot properly hyphenate a line. Please give
323 LaTeX hyphenation hints using the `\-` command.

AUTHOR CONTRIBUTIONS

If you'd like to, you may include a section for author contributions as is done in many journals. This is optional and at the discretion of the authors.

ACKNOWLEDGMENTS

Use unnumbered third level headings for the acknowledgments. All acknowledgments, including those to funding agencies, go at the end of the paper.

REFERENCES

- Yoshua Bengio and Yann LeCun. Scaling learning algorithms towards AI. In *Large Scale Kernel Machines*. MIT Press, 2007.
- Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. *Deep learning*, volume 1. MIT Press, 2016.
- Geoffrey E. Hinton, Simon Osindero, and Yee Whye Teh. A fast learning algorithm for deep belief nets. *Neural Computation*, 18:1527–1554, 2006.

A APPENDIX

You may include other additional sections here.