

Article

Spatial Attention Visualization for Interpretable Trajectory Prediction in Autonomous Driving: Discovering Safety Blind Spots Through Counterfactual Analysis

Xingnan Zhou ¹ and Ciprian Alecsandru ^{1,*}

¹ Department of Building, Civil and Environmental Engineering, Concordia University, Montreal, QC H3G 1M8, Canada

* Correspondence: ciprian.alecsandru@concordia.ca (C.A.)

Version February 12, 2026 submitted to Sustainability

Abstract: Accurate trajectory prediction is critical for autonomous driving safety and energy-efficient planning in sustainable urban mobility systems. While Transformer-based models achieve state-of-the-art prediction performance, their internal attention mechanisms remain opaque, hindering safety validation and regulatory compliance. We present a spatial attention visualization framework that maps Transformer attention weights onto bird’s-eye-view traffic scenes via a novel spatial token bookkeeping mechanism, Gaussian splatting for agent tokens, and polyline painting for lane tokens. Using MTR-Lite, a lightweight Motion Transformer variant (8.48M parameters) trained on the Waymo Open Motion Dataset, we demonstrate the framework through systematic analysis of 100–200 validation scenes. Four key findings emerge: (1) layer-wise entropy analysis reveals non-monotonic hierarchical specialization—encoder layers progressively focus on agents (entropy 5.64→5.36 bits) before the final layer reverses to broad map attention (5.92 bits, 63.6% map tokens); (2) failed predictions exhibit “tunnel vision” with lower entropy (5.72 vs. 5.94 bits) and elevated self-attention (0.049 vs. 0.035); (3) distance-decay masking shows that far-range attention encodes essential context, with even mild masking degrading accuracy by 4.7%; and (4) scene-type analysis confirms dynamic attention adaptation, with dense-traffic scenes allocating 42.3% agent attention versus 18.4% in sparse scenes. We further introduce a counterfactual analysis methodology using controlled scene edits to enable causal reasoning about attention allocation. These findings provide actionable diagnostics for model developers and regulators seeking to validate safe autonomous vehicle deployment, contributing to sustainable urban mobility.

Keywords: trajectory prediction; attention visualization; Transformer; autonomous driving; explainable AI; vulnerable road users; counterfactual analysis; sustainable transportation

1. Introduction

Autonomous vehicles (AVs) represent a transformative technology for achieving sustainable urban mobility. By reducing human-error-related collisions—which account for over 94% of serious crashes according to the U.S. National Highway Traffic Safety Administration [1]—AVs promise substantial improvements in traffic safety, energy efficiency, and urban livability. These benefits align directly with the United Nations Sustainable Development Goals, particularly SDG 11 (Sustainable Cities and Communities) and SDG 13 (Climate Action) [2]. Studies project that widespread AV adoption could reduce traffic fatalities by 90%, decrease fuel consumption by 40% through smoother driving patterns, and reclaim urban space currently dedicated to parking [3–5]. However, realizing these

31 benefits depends critically on achieving public trust and regulatory approval, both of which remain
32 constrained by the opacity of the artificial intelligence systems that underpin autonomous driving [6,7].

33 At the core of modern AV planning pipelines lies motion prediction: forecasting the future
34 trajectories of surrounding vehicles, pedestrians, and cyclists. Transformer-based architectures have
35 emerged as the dominant paradigm for this task, achieving state-of-the-art performance on major
36 benchmarks. Models such as Motion Transformer (MTR) [8,9], Wayformer [10], GameFormer [11],
37 and Scene Transformer [12] leverage multi-head self-attention and cross-attention mechanisms to
38 capture complex interactions among traffic agents and road geometry. Despite their strong quantitative
39 performance, these models operate as *black boxes*: the attention weights that encode inter-agent
40 relationships, lane preferences, and temporal reasoning remain hidden from developers and safety
41 engineers. This lack of interpretability creates three practical barriers. First, when a model produces an
42 erroneous prediction—such as failing to anticipate a left-turning vehicle—there is no principled way
43 to diagnose whether the failure stems from insufficient attention to the relevant agent, the target lane,
44 or the traffic signal. Second, regulatory bodies increasingly demand explanations for safety-critical AI
45 decisions, as codified in the European Union AI Act [13] and NHTSA testing frameworks [1]. Third,
46 without transparency, the general public lacks the evidence necessary to trust autonomous systems,
47 ultimately delaying adoption and the associated sustainability benefits [14,15].

48 Several lines of research have addressed AI interpretability, though significant gaps remain
49 in the trajectory prediction domain. Post-hoc explanation methods such as LIME [16], SHAP [17],
50 and Grad-CAM [18] provide input-level attributions but do not leverage the structured internal
51 attention mechanisms of Transformers. In natural language processing and computer vision,
52 dedicated attention visualization tools—including BERTViz [19], Attention Flow [20], and Transformer
53 Explainability [21]—have demonstrated that attention patterns encode interpretable relationships,
54 though the debate on whether attention constitutes explanation continues [22,23]. Within trajectory
55 prediction, recent work has begun to explore attention-based interpretability: VISTA [24] visualizes
56 pairwise interaction strength, and LMFormer [25] examines lane-conditioned attention maps. However,
57 these efforts focus on isolated aspects of the attention spectrum—either agent-agent interactions or
58 lane selection—and do not provide a unified view of *where* the model attends in physical space, *how* its
59 reasoning evolves across processing layers, and *which* lane structures guide its predictions.

60 In this paper, we present a spatial attention visualization framework for Transformer-based
61 trajectory prediction that goes beyond depicting abstract attention matrices. We specifically
62 adopt a Transformer architecture because its multi-head attention mechanism provides a built-in
63 interpretability window: the attention weights directly reveal the model’s spatial focus, agent and lane
64 priorities, and layer-wise processing evolution. Unlike recurrent architectures such as LSTMs—where
65 hidden states encode temporal dependencies in opaque, entangled vectors—or convolutional networks
66 whose internal feature maps lack token-level semantic correspondence, Transformer attention offers
67 explicit, per-token interaction scores that are naturally amenable to spatial visualization. Indeed, recent
68 work by Zhou and Alecsandru [26] demonstrated that lane-graph conditioning improves LSTM-based
69 trajectory prediction on the Waymo Open Motion Dataset; however, the LSTM hidden states do
70 not afford the spatial interpretability that Transformer attention provides, motivating our choice of
71 architecture for interpretability-focused analysis.

72 Crucially, our primary contribution is a *model-agnostic visualization framework* applicable to
73 any Transformer-based trajectory predictor, not a new state-of-the-art prediction model. To
74 demonstrate and validate this framework, we employ a lightweight variant of the MTR architecture
75 (MTR-Lite, 8.48M parameters) as an *interpretability probe*—a deliberately simplified model that enables
76 rapid experimentation and systematic attention analysis without the computational overhead of
77 production-scale systems. MTR-Lite is trained on 20% of the Waymo Open Motion Dataset [27]
78 (approximately 17,800 scenes), yielding sufficient diversity to reveal meaningful attention patterns
79 while maintaining experimental tractability. The framework itself—encompassing attention extraction,
80 spatial token bookkeeping, and visualization rendering—is architecture-agnostic and can be applied

81 to any Transformer model that produces attention weights over spatially grounded tokens, including
82 production-scale systems such as MTR++ [9], Wayformer [10], and QCNet [28]. Our key technical
83 innovation is a *spatial token bookkeeping* mechanism that maintains a bidirectional mapping between
84 discrete token indices and their physical BEV coordinates, enabling attention weights to be projected
85 as continuous heatmaps directly onto the traffic scene. Using Gaussian splatting for agent tokens and
86 polyline painting for lane tokens, the resulting visualizations provide a spatially grounded view of the
87 model’s attention allocation and its progressive evolution across processing layers.

88 Importantly, we go beyond visualization as an end in itself. By systematically analyzing the
89 spatial distribution of attention across 100–200 Waymo validation scenes, we uncover **quantifiable**
90 **safety-relevant patterns**. We find that failed predictions exhibit significantly lower attention entropy
91 (5.72 vs. 5.94 bits) and elevated self-attention compared to successful predictions—a “tunnel vision”
92 failure mode in which the model over-focuses on the ego agent when it should be distributing attention
93 more broadly. This diagnostic pattern has direct implications for collision risk, as it reveals that
94 prediction failures are accompanied by a measurable, detectable attention pathology. We further design
95 a counterfactual analysis methodology using controllable scene editing (agent removal, injection, and
96 traffic signal manipulation), enabling causal—rather than merely correlational—reasoning about how
97 individual traffic elements influence model attention and prediction outcomes. This combination of
98 spatial visualization, failure diagnostics, and causal experimentation transforms attention analysis
99 from a qualitative illustration into a rigorous diagnostic tool.

100 The main contributions of this work are as follows:

- 101 • We propose a **spatial attention visualization system** that maps abstract Transformer attention
102 weights onto bird’s-eye-view traffic scenes via Gaussian splatting and polyline painting, providing
103 the first spatially grounded interpretation of attention in trajectory prediction.
- 104 • We identify a **tunnel vision failure mode** in which failed predictions exhibit significantly lower
105 attention entropy and elevated self-attention compared to successful predictions, providing a
106 safety-critical diagnostic signal that links measurable attention pathology to prediction failure.
- 107 • We design a **counterfactual attention analysis methodology** using controllable scene generation,
108 providing the infrastructure for causal—rather than merely correlational—reasoning about how
109 individual traffic elements influence model attention and prediction outcomes.
- 110 • We provide **quantitative attention diagnostics**, including layer-wise entropy analysis revealing
111 hierarchical specialization—encoder layers progressively focus on nearby agents (entropy
112 5.64→5.36 bits) while the final layer reverses to broad map attention (entropy 5.92 bits, 63.6%
113 map tokens)—and a distance mask ablation demonstrating that suppressing far-range attention
114 degrades accuracy by 4.7%, demonstrating that distant context is valuable and naive pruning is
115 harmful.
- 116 • We demonstrate the framework across **diverse driving scenarios**—dense intersections (42.3%
117 agent attention), sparse highways (18.4% agent attention, mean attended distance 21.4 m vs.
118 17.0 m at intersections), and failure cases—revealing how the model dynamically adapts its spatial
119 attention distribution to scene complexity.

120 Beyond its technical contributions, this work has direct implications for sustainable
121 transportation. The discovery that prediction failures are accompanied by a measurable tunnel
122 vision pathology—lower entropy and elevated self-attention—has immediate practical consequences:
123 this attention signature could serve as a runtime monitor to flag unreliable predictions before they
124 propagate to the planning module, preventing potential collisions. By quantifying failure-associated
125 attention patterns and analyzing scene-type adaptation, we provide actionable guidance for model
126 developers and regulatory bodies alike. For regulators, spatially grounded attention visualizations
127 offer the kind of human-readable evidence needed to certify AV behavior in complex traffic scenarios,
128 particularly as the European Union AI Act [13] establishes explainability requirements for high-risk AI
129 systems. For the public, the ability to see that an autonomous vehicle “looks at” the correct lanes, traffic
130 signals, and nearby agents before making predictions builds the transparency necessary for trust [6,29].

131 Furthermore, our distance mask ablation reveals that naive attention pruning strategies—which
132 might be pursued for computational efficiency—actually degrade prediction accuracy by 4.7%,
133 cautioning against premature optimization and underscoring the need for interpretability-guided
134 model compression rather than blind sparsification. Ultimately, by combining interpretability with
135 safety diagnostics, our framework helps remove key barriers to safe AV deployment, contributing
136 to the broader goal of reducing road fatalities, lowering transportation emissions, and creating more
137 walkable, livable cities [7,30].

138 The remainder of this paper is organized as follows. Section 2 reviews related work on trajectory
139 prediction, attention visualization, and explainable AI for autonomous driving. Section 3 describes the
140 MTR-Lite architecture, attention extraction mechanism, and visualization pipeline. Section 4 presents
141 quantitative evaluation results and visualization examples. Section 5 discusses the interpretability
142 insights, sustainability implications, and limitations. Section 6 concludes with future directions.

143 2. Related Work

144 2.1. Transformer-Based Trajectory Prediction

145 The application of Transformer architectures [31] to motion forecasting has yielded substantial
146 performance gains on standardized benchmarks. Early work by Gao et al. [32] introduced vectorized
147 scene representations and point-level attention over polyline-encoded map elements, establishing a
148 paradigm adopted by subsequent architectures. Scene Transformer [12] extended this approach
149 to joint multi-agent prediction, employing factored self-attention over agent and time axes to
150 model cooperative and adversarial interactions simultaneously. These foundational architectures
151 demonstrated that attention mechanisms could implicitly capture the spatial and social structure of
152 traffic scenes without explicit graph construction.

153 The Motion Transformer (MTR) family [8,9] introduced a query-based decoder design that has
154 become influential in the field. MTR employs 64 learnable intention queries, initialized from clustered
155 trajectory endpoints, which attend to encoded scene tokens through iterative cross-attention layers.
156 This design separates *global intention localization* (selecting a coarse goal region) from *local movement*
157 *refinement* (producing smooth trajectories conditioned on that goal), yielding strong multi-modal
158 predictions. MTR++ extended this with symmetric scene modeling and pair-wise interaction modules,
159 achieving first place in the 2023 Waymo Open Dataset Motion Prediction Challenge. Notably, the
160 intention query mechanism generates structured attention patterns—each query attends to the agents
161 and lanes relevant to its predicted mode—yet neither MTR nor MTR++ provides tools to visualize or
162 analyze these patterns.

163 Wayformer [10] explored attention-based modality fusion, comparing early, late, and hierarchical
164 fusion strategies for combining agent trajectories, road geometry, and traffic signal features. Their
165 ablation showed that attention over traffic light tokens significantly improves prediction at signalized
166 intersections, hinting at the interpretive value of attention analysis. GameFormer [11] introduced
167 hierarchical game-theoretic decoding with level- k attention, modeling interactive prediction as iterated
168 best-response reasoning. HPTR [33] proposed heterogeneous polyline attention with relative pose
169 encoding and k -nearest-neighbor sparsification, improving efficiency while maintaining the ability to
170 model agent-lane interactions. QCNet [28] developed query-centric encoding that avoids recomputing
171 scene features for each target agent. Most recently, SMART [34] recast trajectory prediction as
172 next-token prediction over discretized motion tokens, achieving state-of-the-art results on the Waymo
173 Sim Agents benchmark with an autoregressive Transformer.

174 Table 1 summarizes the attention mechanisms used by these models and whether any form
175 of attention visualization or interpretability analysis was reported. As the table shows, while all
176 models employ multiple attention mechanisms (self-attention, cross-attention, or both), none provides
177 systematic visualization of the full attention spectrum. This gap motivates our work.

Table 1. Summary of attention mechanisms in state-of-the-art trajectory prediction models. “Viz” indicates whether the paper includes attention visualization or interpretability analysis.

Model	Venue	Self-Attn	Cross-Attn	Query-Based	Viz
VectorNet [32]	CVPR 2020	✓	–	–	–
Scene Trans. [12]	ICLR 2022	✓	–	–	–
MTR [8]	NeurIPS 2022	✓	✓	✓	–
QCNNet [28]	CVPR 2023	✓	✓	✓	–
Wayformer [10]	ICRA 2023	✓	✓	–	Partial
GameFormer [11]	ICCV 2023	✓	✓	✓	–
HPTR [33]	NeurIPS 2023	✓	✓	–	–
MTR++ [9]	TPAMI 2024	✓	✓	✓	–
SMART [34]	NeurIPS 2024	✓	–	–	–
Ours	–	✓	✓	✓	Full

178 2.2. Attention Visualization and Interpretability

179 The question of whether attention weights constitute meaningful explanations has been
 180 extensively debated in the NLP community. Jain and Wallace [22] argued that attention distributions
 181 are not reliable indicators of feature importance, showing that alternative attention configurations
 182 can yield equivalent predictions. Wiegrefe and Pinter [23] countered that attention weights do carry
 183 explanatory signal, particularly when the attention mechanism is constrained or task-specific. This
 184 nuanced view has informed subsequent work: attention is most interpretable when it operates over
 185 semantically meaningful units (words, objects, entities) rather than arbitrary hidden dimensions.

186 Several tools have been developed for visualizing attention in NLP Transformers. BERTViz [19]
 187 provides interactive multi-scale visualizations of attention heads across layers, revealing syntactic and
 188 semantic patterns in pre-trained language models. Abnar and Zuidema [20] introduced Attention Flow,
 189 which propagates attention through the residual stream to attribute model decisions to input tokens.
 190 For Vision Transformers, Chefer et al. [21] combined attention rollout with gradient information to
 191 produce class-specific relevance maps that outperform raw attention in localization tasks.

192 In the trajectory prediction domain, attention-based interpretability has received limited but
 193 growing interest. VISTA [24] introduced a goal-conditioned multi-agent forecasting Transformer
 194 whose social-attention block outputs pairwise attention matrices between agents, demonstrating
 195 that the model assigns increasing attention to agents on potential collision courses. LMFormer [25]
 196 proposed a lane-aware motion prediction Transformer with Mode2Lane cross-attention in the decoder,
 197 showing that attention peaks on lanes aligned with the predicted trajectory. ISE-GT [35] incorporated
 198 interaction strength encoding derived from a driver resistance field model into a graph Transformer,
 199 with a companion Interaction Tendency Reasoning Module that provides post-hoc interpretability by
 200 verifying that inferred interaction tendencies align with human driver intuition.

201 While these contributions represent important progress, they share a common limitation: each
 202 addresses a single facet of the attention spectrum. VISTA focuses exclusively on agent–agent social
 203 attention; LMFormer examines only lane-conditioned decoder attention; ISE-GT provides post-hoc
 204 interaction tendency analysis but not spatial or temporal attention patterns. None offers a unified
 205 framework that simultaneously visualizes (1) the spatial distribution of attention across agents and
 206 lanes, (2) the temporal evolution of attention across decoder layers, and (3) the structural selection
 207 of lane tokens that condition trajectory generation. Our work fills this gap by providing all three
 208 visualization types within a single, integrated pipeline.

209 2.3. Counterfactual Analysis and Controllable Scene Generation

210 Counterfactual reasoning—asking “what would have happened if X were different?”—provides a
 211 principled framework for causal inference in machine learning [36]. Goyal et al. [37] demonstrated
 212 counterfactual visual explanations by identifying minimal image modifications that change a

213 classifier’s prediction, revealing which visual features are causally relevant. In contrast to purely
214 observational analysis, counterfactual experiments can distinguish genuine causal mechanisms from
215 spurious correlations.

216 In autonomous driving, controllable scene generation has emerged as a tool for safety validation
217 and model stress testing. SceneGen [38] learned to place realistic traffic participants in BEV layouts,
218 while TrafficSim [39] modeled multi-agent interactions through learned conditional distributions.
219 More recently, guided diffusion models [40] have enabled fine-grained control over generated
220 traffic scenarios, including adversarial agent placement and rare event synthesis. Ding et al. [41]
221 comprehensively reviewed methods for safety-critical scenario generation, identifying controllability
222 and realism as the two key desiderata.

223 Despite these advances, no prior work has combined controllable scene generation with systematic
224 attention analysis. Existing scene generation methods focus on evaluating prediction *accuracy* (i.e.,
225 whether the model predicts correctly) rather than prediction *attention* (i.e., where the model looks).
226 Our work bridges this gap: by editing real Waymo scenes—removing agents, injecting vulnerable
227 road users, flipping traffic signals—and measuring the resulting changes in attention distributions,
228 we perform the first *counterfactual attention analysis* for trajectory prediction. This enables causal
229 claims about how individual scene elements influence model reasoning, moving beyond correlational
230 findings.

231 2.4. Explainable AI for Autonomous Driving

232 The demand for explainable AI (XAI) in autonomous driving extends beyond academic curiosity
233 to practical necessity. Arrieta et al. [42] provide a comprehensive taxonomy of XAI methods,
234 distinguishing between transparent models (inherently interpretable), post-hoc explanations (applied
235 after training), and hybrid approaches. For safety-critical applications like autonomous driving, they
236 argue that post-hoc methods are insufficient; the model’s internal reasoning process must be accessible
237 and auditable.

238 Zablocki et al. [15] surveyed explainability specifically in deep vision-based driving systems,
239 identifying four key dimensions: *what* is explained (perception, prediction, or planning), *how*
240 explanations are generated (saliency maps, natural language, attention), *who* the audience is
241 (developers, regulators, or passengers), and *when* explanations are provided (offline analysis or
242 real-time). Our work addresses the *prediction* component using *attention-based spatial visualization*,
243 targeting both *developers* (for debugging) and *regulators* (for safety certification), in an *offline analysis*
244 setting.

245 Atakishiyev et al. [29] recently provided an extensive field guide for XAI research in autonomous
246 driving, emphasizing that the gap between model performance and model understanding is the
247 primary obstacle to large-scale deployment. They identify trajectory prediction as a particularly
248 underserved area for interpretability research, noting that most XAI efforts in AV focus on perception
249 (object detection saliency) or planning (reward visualization) rather than the prediction module that
250 bridges them.

251 From a regulatory perspective, the European Union AI Act [13] classifies autonomous driving
252 systems as “high-risk AI” requiring transparency, human oversight, and documented testing. The
253 NHTSA framework [1] similarly calls for testable scenarios and explainable decision processes. These
254 regulatory requirements create a concrete demand for the kind of interpretability tools that our
255 framework provides: spatially grounded visualizations that can demonstrate, for a given scenario,
256 exactly which traffic participants and road structures the model considered before generating its
257 prediction.

258 The connection between AV interpretability and sustainability is increasingly recognized. Taiebat
259 et al. [30] reviewed the energy and environmental implications of connected and automated vehicles,
260 concluding that the magnitude of benefits depends heavily on the pace of adoption, which is in turn
261 constrained by safety assurance and public trust. Litman [43] projects that full AV benefits—including

262 a 60–90% reduction in crash costs and a 30–50% decrease in vehicle-miles traveled per household—will
 263 materialize only when Level 4+ autonomy achieves widespread deployment, a milestone that
 264 requires overcoming the trust deficit. By making trajectory prediction models interpretable, our work
 265 contributes to this trust-building process and, by extension, to the realization of the environmental and
 266 safety benefits that motivate sustainable transportation research.

267 3. Materials and Methods

268 This section presents the dataset, model architecture, attention extraction mechanism, spatial
 269 token bookkeeping system, visualization methods, counterfactual experiment design, and evaluation
 270 metrics that constitute our framework.

271 3.1. Dataset

272 We train and evaluate our model on the Waymo Open Motion Dataset (WOMD) v1.2 [27], one of
 273 the largest and most diverse public benchmarks for trajectory prediction. The full dataset contains
 274 approximately 89,000 driving scenes recorded across six U.S. cities. Each scene spans 91 frames
 275 captured at 10 Hz (9.1 seconds of real-world driving), providing dense temporal coverage of traffic
 276 interactions. We use a 20% subset of the full dataset, yielding approximately 17,800 scenes, split
 277 into 85% training (~15,130 scenes) and 15% validation (~2,670 scenes) using hash-based scene-ID
 278 partitioning for reproducibility.

279 Each scene accommodates up to 100 agent slots, covering three agent types: vehicles, pedestrians,
 280 and cyclists. Every agent is represented as a trajectory with per-frame attributes including position,
 281 velocity, acceleration, heading, and bounding box dimensions. Importantly, the dataset provides
 282 rich map context: a lane graph encoding road topology with successor, predecessor, and left/right
 283 neighbor relationships among lane segments; per-lane attributes including speed limits, lane types,
 284 and boundary markings; and traffic signal states recorded per frame per controlled lane. This
 285 structured map representation is critical for our visualization framework, as it enables projecting
 286 abstract map-token attention weights back onto physically meaningful road geometry.

287 We preprocess the raw data into per-scene pk1 files, each storing a dictionary with three primary
 288 entries: `objects []`, containing per-agent trajectory arrays and metadata; `lane_graph {}`, encoding lane
 289 centerline polylines together with their topological connectivity and attributes; and `traffic_lights []`,
 290 recording per-frame signal states for each controlled lane. This dictionary structure facilitates both
 291 efficient batched training and the counterfactual scene editing experiments described in Section 3.6.

292 3.2. MTR-Lite Architecture

293 Our trajectory prediction model, MTR-Lite, is a lightweight variant of the Motion Transformer
 294 (MTR) [8,9] designed for interpretability research on a single-GPU workstation. The model comprises
 295 8.48M parameters and follows an encode–attend–decode pipeline with four stages: polyline encoding,
 296 scene encoding, motion decoding, and mode selection.

297 3.2.1. Input Representation

The model ingests two types of polyline inputs. *Agent polylines* represent traffic participants: we select $A=32$ agents nearest to the target agent, each described by a polyline of $T_h=11$ historical timesteps (1.0 second of history at 10 Hz). Each timestep carries a 29-dimensional feature vector:

$$\mathbf{f}_{\text{agent}} = \left[\underbrace{x, y}_2, \underbrace{x_{-1}, y_{-1}}_2, \underbrace{v_x, v_y}_2, \underbrace{a_x, a_y}_2, \underbrace{\sin \theta, \cos \theta}_2, \underbrace{w, l}_2, \underbrace{\mathbf{c}_{\text{type}}}_5, \underbrace{\mathbf{e}_{\text{time}}}_{11}, \underbrace{z_{\text{ego}}}_1 \right] \in \mathbb{R}^{29}, \quad (1)$$

298 where (x, y) is the current position, (x_{-1}, y_{-1}) the previous-step position, (v_x, v_y) and (a_x, a_y) the
 299 velocity and acceleration, $(\sin \theta, \cos \theta)$ the heading encoded as sine–cosine pair, (w, l) the bounding
 300 box width and length, $\mathbf{c}_{\text{type}} \in \{0, 1\}^5$ a one-hot agent type encoding (vehicle, pedestrian, cyclist, and

two reserved classes), $\mathbf{e}_{\text{time}} \in \mathbb{R}^{11}$ a learnable temporal positional embedding, and $z_{\text{ego}} \in \{0, 1\}$ a binary indicator of whether the agent is the ego vehicle.

Map polylines represent lane centerlines: we select $M=64$ lane segments nearest to the target agent, each described by $P=20$ points sampled uniformly along the centerline. Each point carries a 9-dimensional feature vector:

$$\mathbf{f}_{\text{map}} = \left[\underbrace{x, y}_2, \underbrace{d_x, d_y}_2, \underbrace{\mathbf{g}_{\text{lane}}}_3, \underbrace{x_{-1}, y_{-1}}_2 \right] \in \mathbb{R}^9, \quad (2)$$

where (x, y) is the point position, (d_x, d_y) the local direction vector, $\mathbf{g}_{\text{lane}} \in \{0, 1\}^3$ encodes lane flags (has traffic control, is intersection lane, is turn lane), and (x_{-1}, y_{-1}) the coordinates of the preceding point in the polyline.

3.2.2. PointNet Encoder

Each polyline—whether agent or map—is independently encoded into a fixed-dimensional token using a PointNet-style architecture [44]. A shared-weight multi-layer perceptron (MLP) processes each point along the polyline:

$$\text{MLP}_{\text{point}} : \mathbb{R}^D \xrightarrow{\text{Linear}} \mathbb{R}^{64} \xrightarrow{\text{ReLU}} \mathbb{R}^{128} \xrightarrow{\text{ReLU}} \mathbb{R}^{256} \xrightarrow{\text{ReLU}} \mathbb{R}^{256}, \quad (3)$$

where D is the input feature dimension (29 for agents, 9 for map). A symmetric max-pooling operation aggregates the per-point features across the polyline's temporal or spatial extent, producing a single 256-dimensional vector that is invariant to point ordering. A post-aggregation MLP refines this representation:

$$\text{MLP}_{\text{post}} : \mathbb{R}^{256} \xrightarrow{\text{Linear}} \mathbb{R}^{256} \xrightarrow{\text{ReLU}} \mathbb{R}^{256}, \quad (4)$$

followed by layer normalization [45]. The agent and map encoders share this architectural template but maintain separate learned parameters. This stage produces 32 agent tokens and 64 map tokens, each in \mathbb{R}^{256} .

3.2.3. Scene Encoder

The 96 tokens (32 agent + 64 map) are concatenated into a single sequence and processed by a global self-attention encoder comprising $L_e=4$ Transformer encoder layers [31]. Each layer applies pre-norm multi-head self-attention with $H=8$ heads ($d_k=d_v=32$) and a position-wise feed-forward network (FFN) with hidden dimension 1024:

$$\mathbf{z}' = \mathbf{z} + \text{MultiHead}(\text{LN}(\mathbf{z}), \text{LN}(\mathbf{z}), \text{LN}(\mathbf{z})), \quad (5)$$

$$\mathbf{z}'' = \mathbf{z}' + \text{FFN}(\text{LN}(\mathbf{z}')), \quad (6)$$

where $\text{LN}(\cdot)$ denotes layer normalization and the residual connections follow the pre-norm convention. Global self-attention allows every token to attend to every other token, enabling agent-agent, agent-map, map-agent, and map-map interactions to emerge naturally. After the final encoder layer, the 96 tokens are split back into 32 encoded agent tokens and 64 encoded map tokens.

3.2.4. Motion Decoder

For each target agent, the decoder generates $K_0=64$ candidate trajectory modes using an intention-query mechanism inspired by MTR [8]. Each of the 64 intention queries is initialized by summing (i) a learned embedding of a 2D anchor point (obtained via k -means clustering of training-set trajectory endpoints) with (ii) a context embedding derived from the target agent's encoded token. The decoder consists of $L_d=4$ layers, each performing:

- 321 1. **Agent cross-attention:** intention queries attend to the 32 encoded agent tokens, capturing
322 dynamic interactions.
- 323 2. **Map cross-attention:** intention queries attend to the 64 encoded map tokens, selecting lane-level
324 guidance.
- 325 3. **Feed-forward network:** position-wise nonlinear transformation with hidden dimension 1024.

326 Each decoder layer is followed by a per-layer trajectory head (for deep supervision) that regresses
327 a trajectory of $T_f=80$ future timesteps (8.0 seconds at 10 Hz) and a scalar confidence logit from the
328 refined query embedding. The deep supervision loss weights are $[0.2, 0.2, 0.2, 0.4]$ from the first to the
329 last layer.

330 3.2.5. Mode Selection

331 From the 64 candidate modes produced by the final decoder layer, we apply distance-based
332 non-maximum suppression (NMS) with a threshold of 2.0 m on trajectory endpoints. This yields $K=6$
333 diverse output modes, each comprising a predicted trajectory $\hat{Y}_k \in \mathbb{R}^{80 \times 2}$ and a confidence score \hat{p}_k .
334 The confidence scores are normalized via softmax to form a probability distribution over modes.

335 3.2.6. Training

336 The model is trained for 60 epochs with the AdamW optimizer [46] (learning rate 10^{-4} , weight
337 decay 0.01), using a linear warmup over 5 epochs followed by cosine annealing decay. Automatic
338 mixed-precision (AMP) training with float16 [47] is employed throughout. The loss function combines
339 a cross-entropy classification loss over mode scores with a smooth- ℓ_1 regression loss over trajectory
340 coordinates, applied at every decoder layer with deep supervision. Gradient clipping is set to a
341 maximum norm of 1.0, and training uses batch size 4 with 8-step gradient accumulation (effective
342 batch size 32).

343 3.3. Attention Extraction Framework

344 A central requirement of our visualization pipeline is the ability to extract per-head attention
345 weight matrices from every layer without altering the model's predictions. We accomplish this through
346 custom Transformer layers that extend PyTorch's `nn.MultiheadAttention` with a lightweight capture
347 mechanism.

348 3.3.1. Attention-Capture Layers

349 We implement two custom layer classes: `AttentionCaptureEncoderLayer` for the scene
350 encoder and `AttentionCaptureDecoderLayer` for the motion decoder. Both accept a boolean
351 flag `capture_attention` on their forward pass. When this flag is set to `True`, the underlying
352 multi-head attention call is invoked with `need_weights=True` and `average_attn_weights=False`,
353 causing PyTorch to return the full per-head attention weight tensor rather than discarding it or
354 averaging across heads. When the flag is `False` (the default during training), no attention weights are
355 computed or stored, incurring zero overhead.

356 3.3.2. AttentionMaps Data Structure

357 All captured weights from a single forward pass are organized in an `AttentionMaps` dataclass
358 with three primary fields:

- 359 • `scene_attentions`: a list of $L_e=4$ tensors, each of shape (B, H, N, N) where $N=A+M=96$,
360 representing per-head self-attention weights at each encoder layer. Each tensor is a row-stochastic
361 matrix (rows sum to 1) in the last dimension.
- 362 • `decoder_agent_attentions`: a list of $L_d=4$ tensors per target agent, each of shape (B, H, K_0, A)
363 where $K_0=64$ and $A=32$, representing per-head cross-attention from intention queries to agent
364 tokens.

- 365 • `decoder_map attentions`: a list of $L_d=4$ tensors per target agent, each of shape (B, H, K_0, M)
 366 where $M=64$, representing per-head cross-attention from intention queries to map tokens.

367 This structure provides accessor methods for extracting specific submatrices: agent-to-agent
 368 attention, agent-to-map attention, map-to-agent attention, and per-mode decoder attention.
 369 An `aggregate_heads` method supports both mean and max aggregation across heads, and a
 370 `compute_entropy` method computes Shannon entropy in bits for quantitative analysis.

371 3.4. Spatial Token Bookkeeping

372 The key technical innovation enabling our visualization approach is a *spatial token bookkeeping*
 373 system that maintains a bidirectional mapping between the abstract token index space used by the
 374 Transformer and the continuous bird’s-eye-view (BEV) coordinate space of the physical scene. Without
 375 this mapping, attention weights are merely entries in a matrix indexed by opaque integers; with it,
 376 each attention value acquires a spatial interpretation.

377 For each *agent token* $i \in \{0, \dots, A-1\}$, the bookkeeper stores the agent’s BEV position (x_i, y_i) at
 378 the anchor frame, heading angle θ_i , bounding box dimensions (w_i, l_i) , and agent type. For each *map*
 379 *token* $j \in \{0, \dots, M-1\}$, the bookkeeper stores the full lane centerline polyline $\{(x_{j,p}, y_{j,p})\}_{p=1}^P$ in BEV
 380 coordinates.

381 This bookkeeping enables two critical operations. First, given a row of the attention matrix (e.g.,
 382 the ego agent’s attention over all 96 scene tokens at encoder layer l), we can project each attention value
 383 onto its corresponding spatial location, transforming a 96-element vector into a spatially grounded
 384 heatmap over the BEV plane. Second, given a decoder cross-attention row for a specific intention
 385 query, we can separately project agent attention and map attention onto the BEV, revealing which
 386 physical agents and which lane structures guide the model’s trajectory prediction for that mode. All
 387 coordinate transforms use a configurable BEV grid with resolution 0.5 m/pixel and a 120×120 m field
 388 of view centered on the target agent.

389 3.5. Visualization Methods

390 We develop three complementary visualization types, each designed to illuminate a different
 391 facet of the model’s attention-based reasoning.

392 3.5.1. Space-Attention BEV Heatmap

393 This visualization answers the question: *where in physical space does the model concentrate its*
 394 *attention?* Given a target agent and a selected encoder or decoder layer, we extract the attention weight
 395 vector and project it onto the BEV plane as follows:

1. For each valid *agent token* i with attention weight α_i (averaged across $H=8$ heads), we render a 2D isotropic Gaussian centered at the agent’s BEV position (x_i, y_i) with standard deviation $\sigma=3.0$ m:

$$G_i(x, y) = \alpha_i \cdot \exp\left(-\frac{(x - x_i)^2 + (y - y_i)^2}{2\sigma^2}\right). \quad (7)$$

- 396 2. For each valid *map token* j with attention weight β_j , we paint the lane centerline polyline onto
 397 the heatmap grid using Bresenham line rasterization with a stroke width of 2.0 m, followed by
 398 Gaussian smoothing.
- 399 3. The contributions from all agent and map tokens are accumulated additively into a single heatmap,
 400 which is then clipped at the 95th percentile and normalized to $[0, 1]$.
- 401 4. The heatmap is rendered using the magma colormap with $\alpha=0.7$ transparency, overlaid on a
 402 grayscale BEV rendering of lane boundaries, agent bounding boxes, the target agent’s historical
 403 trajectory (blue), ground-truth future (green dashes), and predicted trajectories (red).

404 3.5.2. Time-Attention Refinement Diagram

405 This visualization answers the question: *how does the model’s attention evolve across decoder layers?*
 406 For the winning mode (highest-confidence trajectory after NMS), we extract the cross-attention weights
 407 from each of the $L_d=4$ decoder layers and present them as a four-panel strip chart. Each panel displays
 408 a ranked bar chart of the top-10 most-attended tokens (labeled by type and index, e.g., “Vehicle_3”,
 409 “Lane_12”), with a consistent vertical scale across all panels for direct comparability. This visualization
 410 reveals the iterative refinement process: early decoder layers typically distribute attention broadly
 411 across candidate lanes and nearby agents, while later layers concentrate attention on the selected goal
 412 lane and the most interaction-relevant agents.

413 3.5.3. Lane-Token Activation Map

414 This visualization answers the question: *which lane structures guide the model’s trajectory prediction?*
 415 For the winning mode at the final decoder layer, we extract the map cross-attention vector $(\beta_1, \dots, \beta_M)$
 416 and use it to color-code each of the $M=64$ lane centerline polylines on the BEV. High-attention lanes
 417 are rendered in warm colors (red–yellow) with thick strokes, while low-attention lanes are rendered in
 418 cool colors (blue–green) with thin strokes, using a diverging colormap. An accompanying sidebar bar
 419 chart ranks the top-10 lanes by attention weight. This visualization directly reveals the model’s lane
 420 selection strategy and can be compared against the ground-truth future trajectory to assess whether
 421 the model attends to the correct lane.

422 3.6. Counterfactual Experiment Methodology

423 Beyond observational attention analysis, we design controlled counterfactual experiments that
 424 isolate the causal effect of specific scene elements on the model’s attention distribution and trajectory
 425 predictions. The core methodology is as follows.

426 3.6.1. Scene Editing

427 Because our data are stored as `pk1` dictionaries, counterfactual scenes are created by direct
 428 manipulation of the dictionary entries. Three editing operations are supported:

- 429 • **Agent removal:** setting a target agent’s valid mask to `False` across all timesteps, effectively
 430 removing it from the scene while preserving all other elements.
- 431 • **Traffic light modification:** overwriting the signal state entries for a specified lane from green to
 432 red (or vice versa) across relevant frames.
- 433 • **Agent injection:** inserting a new agent (e.g., a pedestrian) at a specified BEV position with
 434 appropriate kinematic attributes, occupying a previously unused agent slot.

435 3.6.2. Controlled Comparison

Each counterfactual experiment follows an A/B protocol. The original scene \mathcal{S} and the modified scene \mathcal{S}' are both processed through the model in evaluation mode with attention capture enabled. Because the only difference between \mathcal{S} and \mathcal{S}' is the targeted edit, any change in attention or prediction can be attributed to the modified element. We compute attention difference maps:

$$\Delta \mathbf{A} = \mathbf{A}(\mathcal{S}') - \mathbf{A}(\mathcal{S}), \quad (8)$$

436 where $\mathbf{A}(\cdot)$ denotes the head-averaged attention matrix at a specified layer. Positive entries in $\Delta \mathbf{A}$
 437 indicate tokens that received *more* attention after the modification; negative entries indicate attention
 438 *withdrawn* from those tokens.

439 3.6.3. Experiment Types

440 We conduct three types of counterfactual experiments:

- 441 1. **Agent removal and attention redistribution:** A key interacting agent (e.g., an oncoming vehicle
 442 at an intersection) is removed from the scene. We measure how the attention previously allocated
 443 to this agent redistributes across the remaining tokens. The hypothesis is that attention flows to
 444 the next-most-relevant agents and lanes, revealing the model’s latent priority ordering.
- 445 2. **Traffic light state flip and attention adaptation:** A traffic signal controlling the target agent’s
 446 lane is toggled from green to red (or red to green). We measure changes in both the attention
 447 distribution and the predicted trajectories. The hypothesis is that a green-to-red flip causes
 448 increased attention to the stop line and deceleration in the predicted trajectory.
- 449 3. **VRU injection at varying distances:** A pedestrian is injected at distances of $d \in$
 450 $\{5, 10, 15, 20, 30, 50\}$ meters from the target agent’s predicted path. We measure the attention
 451 allocated to the injected pedestrian as a function of distance, identifying the distance threshold
 452 below which the model begins to attend to the VRU. This experiment directly quantifies the
 453 model’s safety-relevant perception range for vulnerable road users.

454 3.7. Evaluation Metrics

455 Our evaluation employs two families of metrics: standard trajectory prediction metrics to validate
 456 model competence, and attention-specific metrics to quantify the interpretability and safety relevance
 457 of attention patterns.

458 3.7.1. Trajectory Prediction Metrics

459 We report three standard metrics, each computed over $K=6$ predicted modes:

- **Minimum Average Displacement Error (minADE@6):** the minimum over all K modes of the mean ℓ_2 distance between predicted and ground-truth positions across all future timesteps:

$$\text{minADE@K} = \min_{k \in \{1, \dots, K\}} \frac{1}{T_f} \sum_{t=1}^{T_f} \|\hat{\mathbf{y}}_k^{(t)} - \mathbf{y}^{(t)}\|_2. \quad (9)$$

- **Minimum Final Displacement Error (minFDE@6):** the minimum over all K modes of the ℓ_2 distance at the final timestep:

$$\text{minFDE@K} = \min_{k \in \{1, \dots, K\}} \|\hat{\mathbf{y}}_k^{(T_f)} - \mathbf{y}^{(T_f)}\|_2. \quad (10)$$

- **Miss Rate (MR@6):** the fraction of samples for which minFDE@K exceeds a threshold of 2.0 m:

$$\text{MR@K} = \frac{1}{|\mathcal{D}|} \sum_{i \in \mathcal{D}} \mathbb{1}[\text{minFDE@K}_i > 2.0 \text{ m}]. \quad (11)$$

460 3.7.2. Attention Analysis Metrics

461 To quantify attention properties beyond visual inspection, we employ:

- **Shannon Entropy:** measures the uniformity of an attention distribution $\alpha = (\alpha_1, \dots, \alpha_N)$:

$$H(\alpha) = - \sum_{i=1}^N \alpha_i \log_2 \alpha_i \quad [\text{bits}]. \quad (12)$$

462 An entropy of $\log_2 N$ indicates perfectly uniform attention; low entropy indicates focused attention.
 463 We track entropy across layers to quantify the progressive focusing hypothesis.

- **Gini Coefficient** (defined here for completeness; our current analysis focuses on Shannon entropy): measures the inequality (sparsity) of the attention distribution. For a sorted attention vector $\alpha_{(1)} \leq \dots \leq \alpha_{(N)}$, the Gini coefficient is:

$$G(\alpha) = \frac{2 \sum_{i=1}^N i \cdot \alpha_{(i)}}{N \sum_{i=1}^N \alpha_{(i)}} - \frac{N+1}{N}. \quad (13)$$

A Gini coefficient of 0 corresponds to uniform attention; a value approaching 1 indicates that virtually all attention is concentrated on a single token.

- **Attention-to-Ground-Truth-Lane Correlation** (defined as part of the analysis toolkit; our current study focuses on entropy and agent/map attention decomposition): for each sample, we identify the ground-truth lane (the lane polyline minimizing mean point-to-polyline distance to the future trajectory) and extract the decoder’s attention weight to this lane token. We then compute the Pearson correlation coefficient between this attention weight and the sample’s minADE@6 across the validation set. A significant negative correlation ($r < 0$, $p < 0.05$) would indicate that higher attention to the correct lane is associated with lower prediction error.

3.7.3. VRU Safety Metrics

To quantify safety-relevant attention properties for vulnerable road users (VRUs), we define:

- **Attention Ratio**: the ratio of mean attention allocated to a pedestrian token versus a vehicle token at the same distance d from the target agent’s predicted path:

$$R_{\text{attn}}(d) = \frac{\mathbb{E}[\alpha_{\text{ped}}(d)]}{\mathbb{E}[\alpha_{\text{veh}}(d)]}. \quad (14)$$

A ratio of 1.0 indicates parity; values below 1.0 indicate systematic under-attention to pedestrians relative to vehicles.

- **Attention Threshold for Collision Avoidance**: using the VRU injection experiments at varying distances, we identify the critical distance d^* at which the injected pedestrian’s attention weight first exceeds a predefined threshold (defined as twice the mean background attention level). Distances $d > d^*$ represent a potential blind zone where the model may fail to account for the VRU in its predictions.

4. Results

4.1. Trajectory Prediction Performance

Table 2 presents the trajectory prediction performance of MTR-Lite on the Waymo Open Motion Dataset validation set, compared against a Constant Velocity (CV) baseline that linearly extrapolates each agent’s last observed velocity. We report results at three prediction horizons (3, 5, and 8 seconds) to characterize both short-term and long-term forecasting accuracy. The CV baseline provides a physics-based lower bound: any learned model that cannot outperform simple linear extrapolation offers no value beyond Newtonian kinematics.

MTR-Lite reduces the 8-second ADE by 54.4% relative to the CV baseline (2.314 m vs. 5.071 m) and the FDE by 54.7% (6.401 m vs. 14.131 m), confirming that the Transformer architecture captures interaction and map-conditioned dynamics far beyond linear extrapolation. Notably, the improvement is most pronounced at longer horizons: the ADE reduction grows from 19.7% at 3 seconds (0.757 m vs. 0.943 m) to 47.5% at 5 seconds (1.237 m vs. 2.356 m) and 54.4% at 8 seconds, demonstrating that the model’s learned scene understanding is especially valuable for long-horizon forecasting where linear assumptions break down. The miss rate is comparable (54.6% vs. 56.8%) because the 2.0 m threshold is stringent even for the learned model on an 8-second horizon.

Table 2. Trajectory prediction performance on the Waymo Open Motion Dataset (20% training subset, full validation set with 13,388 scenes and 99,370 agent predictions). $K = 6$ modes, 8-second horizon (80 timesteps at 10 Hz). The Constant Velocity baseline is deterministic ($K = 1$), so $\text{minADE}@K = \text{ADE}@1$ for all K .

Model	Params	minADE@6	minFDE@6	MR@6	ADE@3s	ADE@5s	ADE@8s
Constant Velocity	—	5.071	14.131	0.568	0.943	2.356	5.071
MTR-Lite	8.48M	2.314	6.401	0.546	0.757	1.237	2.314

Table 3. Per-agent-type performance breakdown on the full validation set. Cyclists exhibit significantly higher prediction difficulty with 88.1% miss rate, reflecting their unique combination of vehicle-like speeds and pedestrian-like maneuverability.

Agent Type	minADE@6	minFDE@6	MR@6	Count
Vehicle	2.331	6.470	0.540	93,801
Pedestrian	1.579	3.881	0.594	4,536
Cyclist	3.931	11.133	0.881	1,033

498 Table 3 presents the performance breakdown by agent type. While vehicles (94.4% of the
 499 dataset) and pedestrians achieve moderate miss rates around 54–59%, cyclists stand out as the
 500 most challenging category with an 88.1% miss rate. This elevated difficulty reflects cyclists’ unique
 501 behavioral characteristics: they combine vehicle-like speeds (enabling rapid position changes over
 502 the 8-second horizon) with pedestrian-like maneuverability (allowing sudden direction changes and
 503 lane-crossing behavior). The cyclist category’s underrepresentation in the training data (only 1.0% of
 504 predictions) further compounds the prediction challenge, supporting the failure diagnosis finding in
 505 Section 4.8 that underrepresented agent types are systematically harder to predict.

506 4.2. Spatial Attention Visualization

507 Figure 1 presents the core contribution of our visualization framework: spatial attention overlays
 508 projected onto bird’s-eye-view traffic scenes. Each panel shows the combined agent-token Gaussian
 509 splatting and lane-token attention painting for a different scene type.

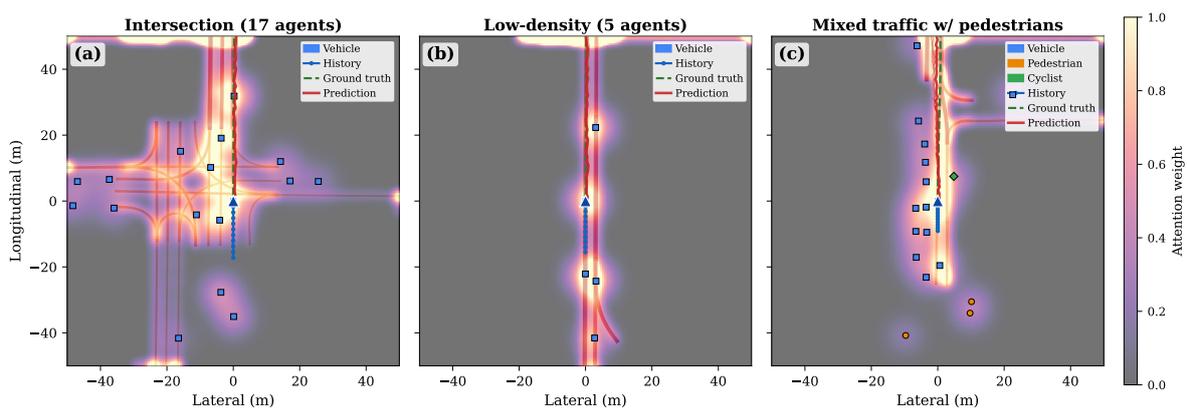


Figure 1. Spatial attention visualization across three scene types. (a) Dense intersection with 17 agents: attention spreads broadly across the intersection center and approaching lanes, reflecting the model’s need to monitor multiple potential conflict points. (b) Low-density scene with 5 agents: attention concentrates tightly along the target agent’s forward path and immediate surroundings. (c) Mixed traffic with pedestrians: attention covers the road network with notable hotspots near the pedestrian (orange circle). Gaussian splatting ($\sigma = 3.0$ m) is used for agent tokens; lane attention is painted along centerlines. Colorbar indicates normalized attention weight.

510 Several qualitative patterns emerge from the spatial overlays. In the dense intersection scene
 511 (Figure 1a), the attention heatmap covers the entire intersection region, with particularly high activation
 512 at the intersection center where multiple trajectories converge. Attention extends along all approaching
 513 lanes, consistent with the model monitoring potential conflict points from every direction. In contrast,
 514 the low-density scene (Figure 1b) shows a markedly narrower attention distribution concentrated
 515 along the target agent's forward path. The model allocates minimal attention to distant or lateral
 516 regions, reflecting the reduced complexity of the scene. The mixed-traffic scene (Figure 1c) shows
 517 broadly distributed attention across the road network with visible hotspots near the pedestrian location,
 518 suggesting the model registers the presence of vulnerable road users in its spatial reasoning.

519 Figure 2 provides a detailed single-scene view of the intersection scenario, illustrating how the
 520 combined agent and lane attention forms a coherent spatial attention field.

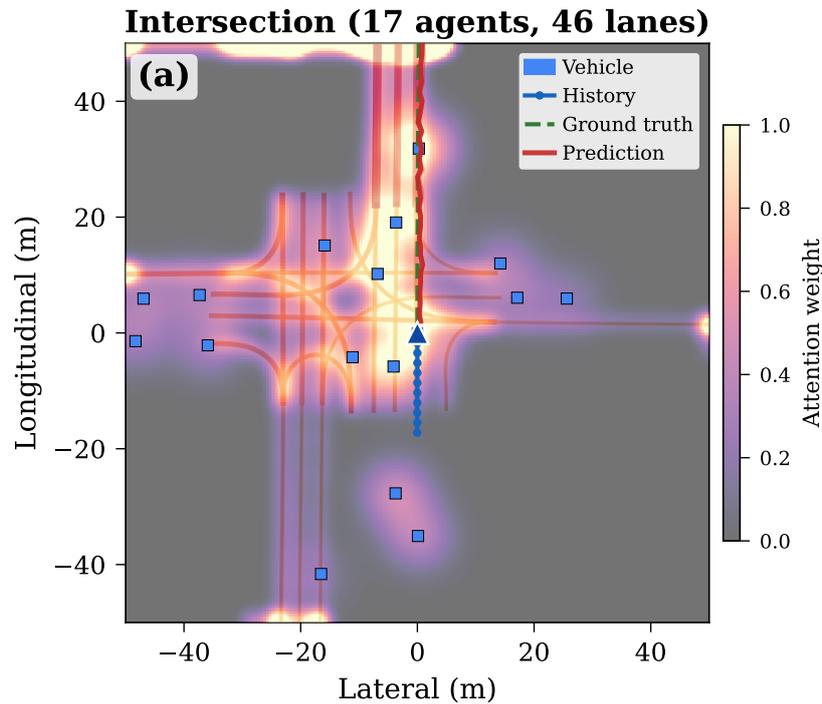


Figure 2. Detailed spatial attention overlay for an intersection scenario (17 agents, 46 lanes). The ego vehicle (blue triangle) is located at the center. Attention is highest along the forward trajectory path and at the intersection center, with secondary peaks at nearby vehicles and approaching lanes. Blue squares denote vehicles; the green dashed line shows ground truth; the red solid line shows the best predicted trajectory.

521 4.3. Layer-Wise Attention Evolution

522 To quantify how attention evolves across processing layers, we compute the Shannon entropy
 523 $H = -\sum_j w_j \log_2 w_j$ of the attention distribution for each encoder layer, averaged across 100–200
 524 validation scenes. Figure 3 presents the results.

525 The key finding is a *non-monotonic* entropy pattern that contradicts the naive expectation of simple
 526 progressive focusing. Layers 0 through 2 progressively decrease entropy (5.64→5.50→5.36 bits) while
 527 increasing agent attention share (49.7%→55.1%→62.4%), consistent with the model narrowing its
 528 focus onto the most relevant traffic agents. However, Layer 3 reverses this trend: entropy increases to
 529 5.92 bits and the attention composition flips to 63.6% map tokens. This pattern suggests a two-phase
 530 processing strategy: *agent interaction modeling* (Layers 0–2) followed by *map-conditioned trajectory*
 531 *planning* (Layer 3), where the final layer broadens attention to incorporate the road geometry needed
 532 for generating lane-following trajectories.

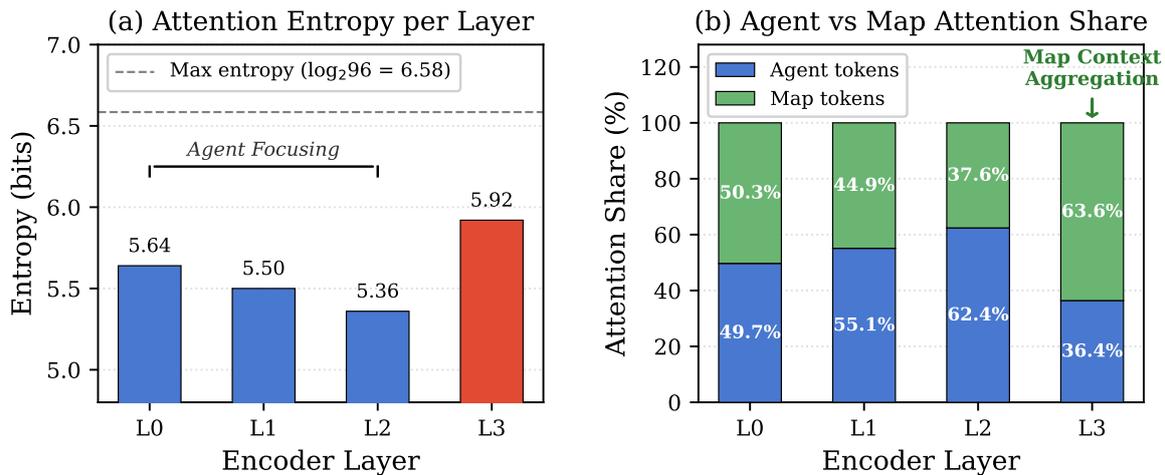


Figure 3. Layer-wise attention analysis across the four encoder layers. (a) Shannon entropy reveals a non-monotonic pattern: entropy decreases from Layer 0 (5.64 bits) through Layer 2 (5.36 bits) as the model focuses on relevant agents, but Layer 3 reverses to 5.92 bits. The dashed line marks maximum entropy for 96 tokens ($\log_2 96 = 6.58$ bits). (b) Agent vs. map attention share explains the reversal: Layers 0–2 progressively increase agent attention (49.7%→62.4%), while Layer 3 pivots to 63.6% map attention, broadening its scope to incorporate road geometry for trajectory generation.

533 4.4. Attention Head Specialization

534 While the layer-wise analysis reveals aggregate attention patterns, it obscures within-layer
 535 heterogeneity across attention heads. To investigate whether individual heads specialize in different
 536 aspects of scene understanding, we compute per-head agent-to-map attention ratios and visualize
 537 their spatial attention patterns. Figure 4 presents the results.

538 The head-wise analysis reveals functional specialization that is invisible in aggregate layer
 539 statistics. While Section 4.3 showed that Layer 3 shifts to 63.6% map attention overall, this transition
 540 is not uniform across heads. Head 5 in Layer 3 allocates 93.3% of its attention to map tokens,
 541 strongly focusing on lane geometry and road boundaries. In contrast, Head 3 retains 58.8% agent
 542 attention—acting as an “agent sentinel” that preserves social context even as other heads pivot
 543 toward spatial planning. The spread between the most agent-focused and most map-focused heads in
 544 Layer 3 reaches 52.1 percentage points, confirming that the layer’s aggregate map-dominance conceals
 545 substantial functional diversity.

546 The spatial heatmaps in Figure 4b illustrate this disentanglement qualitatively. Map-focused
 547 heads produce attention that aligns tightly with lane centerlines and extends along the road network,
 548 while agent-focused heads concentrate on vehicle clusters and intersection conflicts. This head-level
 549 specialization suggests that the model learns complementary representations within each layer: some
 550 heads track dynamic agents, others encode static geometry, and the final decoder aggregates both
 551 sources of information. The persistence of agent-specialized heads in Layer 3 contradicts a naive
 552 interpretation of the layer as purely map-focused, revealing instead a collaborative division of labor
 553 across attention heads.

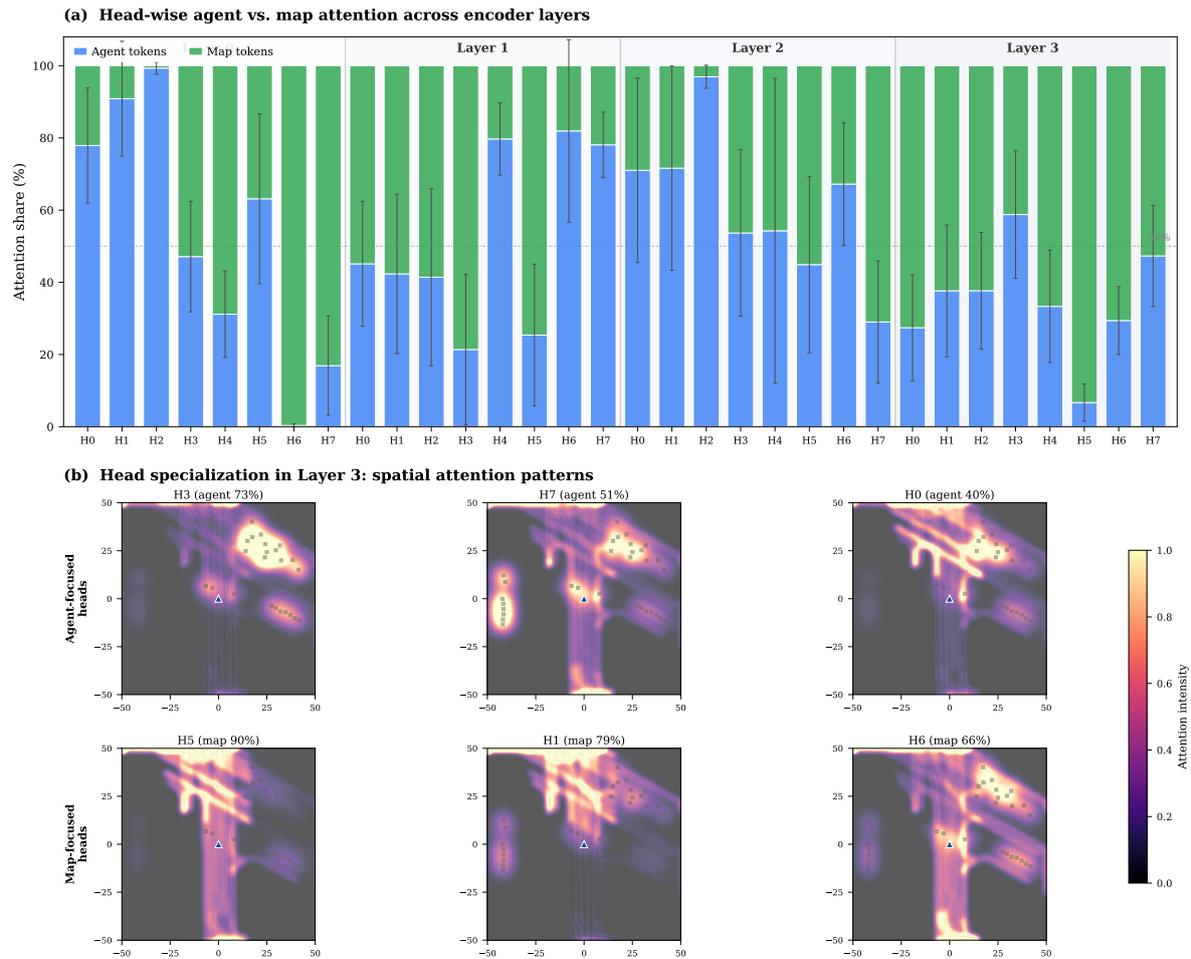


Figure 4. Attention head specialization analysis. (a) Per-head agent-to-map attention ratio across all four encoder layers (8 heads per layer, 32 bars total). Layer 3 exhibits the strongest head-wise disentanglement, with Head 5 allocating 93.3% of its attention to map tokens while Head 3 maintains 58.8% agent attention. (b) BEV spatial attention heatmaps for the three most agent-focused heads versus the three most map-focused heads in Layer 3, demonstrating qualitatively distinct attention patterns.

554 4.5. Lane-Token Activation Analysis

555 Figure 5 visualizes the cumulative decoder map-attention projected onto the lane topology,
 556 revealing which road structures the model prioritizes when generating trajectory predictions.

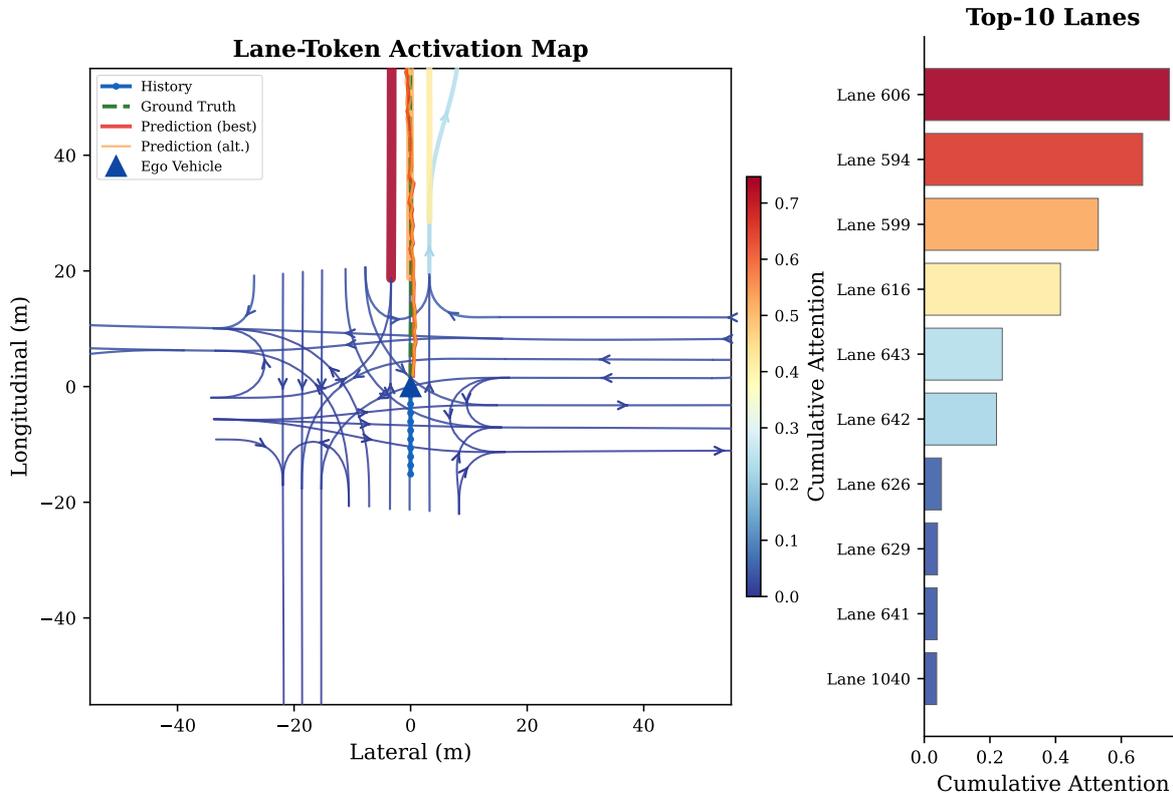


Figure 5. Lane-token activation map for an intersection scenario. **Left:** BEV with lanes colored by cumulative decoder map-attention (warm = high, cool = low), with line width proportional to attention weight. The two highest-attended lanes (606 and 594) align closely with the ego vehicle's forward trajectory (green dashed line). **Right:** Bar chart ranking the top-10 most-attended lane tokens. Lane 606 (cumulative attention 0.68) and Lane 594 (0.39) dominate, both corresponding to the northbound road segment. Alternative mode predictions (shown in orange) attend to adjacent lanes.

557 The lane activation map provides direct evidence that the model's lane attention is spatially
 558 coherent and functionally meaningful. The two highest-attended lanes (606 and 594, with cumulative
 559 attention weights of 0.68 and 0.39 respectively) align precisely with the ego vehicle's ground-truth
 560 forward trajectory. The steep drop-off to the third-ranked lane (599, attention 0.35) indicates high
 561 selectivity, with the model concentrating over 60% of its total lane attention on just two lane
 562 segments. Alternative prediction modes (shown in orange) attend to adjacent lanes, confirming
 563 that the multi-modal nature of the predictions is reflected in the attention structure.

564 4.6. Decoder Attention Refinement

565 Figure 6 presents the temporal evolution of decoder attention across four refinement layers,
 566 illustrating how the winning mode's intention query redistributes its focus during iterative trajectory
 567 generation.

568 The decoder refinement reveals two notable patterns. First, the set of top-attended tokens is
 569 remarkably stable across layers: the same two vehicles (Veh_16 and Veh_25) and two lanes (Lane_53
 570 and Lane_63) appear in the top-5 across all four decoder layers, suggesting that the model identifies
 571 the most relevant scene elements early and refines their relative weighting iteratively. Second, ego
 572 self-attention systematically decreases across layers (0.172→0.116→0.116→0.131), while attention to

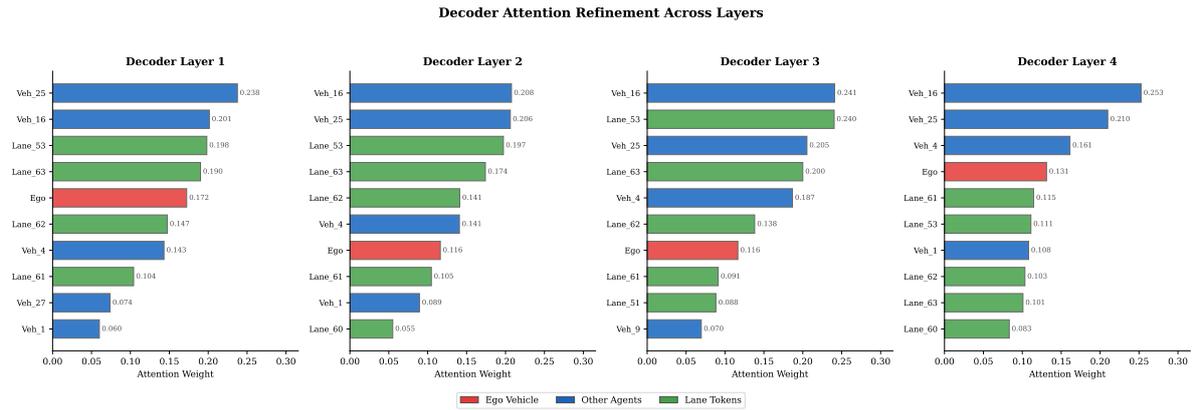


Figure 6. Decoder attention refinement across four layers. Each panel shows the top-10 most-attended tokens for the winning mode’s intention query. **Red bars:** ego vehicle self-attention. **Blue bars:** other agent tokens. **Green bars:** lane tokens. Across layers, the model consistently attends to the same key vehicles (Veh_16 and Veh_25), but their relative importance shifts: ego self-attention decreases from 0.172 (Layer 1) to 0.131 (Layer 4), while the dominant neighbor vehicle (Veh_16) increases from 0.201 to 0.253.

573 the dominant neighbor (Veh_16) increases (0.201→0.208→0.241→0.253). This shift from self-focused to
 574 neighbor-focused attention during refinement indicates that later decoder layers increasingly condition
 575 the trajectory on the behavior of key interacting agents.

576 4.7. Mode-Specific Attention Disentanglement

577 To investigate whether the model's $K = 6$ prediction modes truly reflect distinct reasoning
 578 strategies or merely produce different trajectory outputs from identical attention, we compare attention
 579 patterns across modes for the same target agent. Figure 7 presents the analysis.

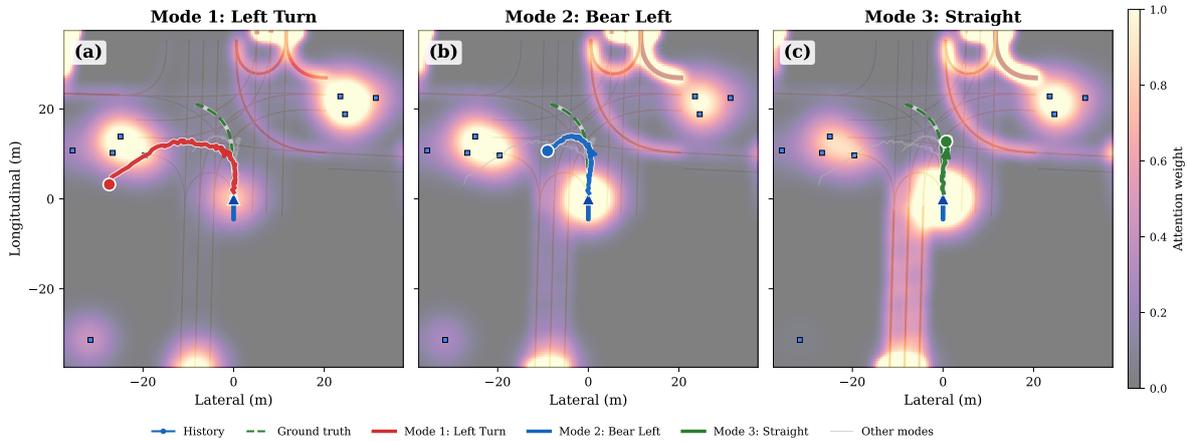


Figure 7. Mode-specific attention comparison for three maneuver intentions: Left Turn, Bear Left, and Straight. Each panel shows the spatial attention distribution (top) and agent-wise attention bar chart (bottom) for the corresponding mode's intention query. Left Turn mode distributes attention broadly across non-ego agents in the turning path (ego self-attention: 0.15), while Straight mode concentrates heavily on the ego vehicle itself (self-attention: 0.43) with forward lane focus. Jensen-Shannon Divergence between Left Turn and Straight map attention is 0.128, confirming quantitative disentanglement.

580 The mode-specific analysis reveals that different prediction modes attend to fundamentally
 581 different scene elements, providing evidence that multi-modal prediction reflects diverse reasoning
 582 strategies rather than superficial output variation. The Left Turn mode distributes attention across
 583 non-ego agents in the potential turning path, with ego self-attention limited to 0.15—the model surveys
 584 surrounding vehicles to assess gap acceptance feasibility. In contrast, the Straight mode exhibits ego
 585 self-attention of 0.43, nearly three times higher, concentrating on the target agent's own state and
 586 forward lane geometry. This divergence demonstrates that modes prioritize distinct contextual cues
 587 aligned with their intended maneuvers.

588 The Jensen-Shannon Divergence of 0.128 between Left Turn and Straight mode map attention
 589 distributions quantifies this disentanglement objectively. Values above 0.10 indicate substantial
 590 distributional difference, confirming that the modes do not share a common attention strategy. The
 591 Bear Left mode occupies an intermediate position, blending elements of both patterns. These findings
 592 validate the multi-modal architecture's design assumption: that trajectory diversity requires intention
 593 diversity, and intention diversity manifests in attention allocation. For interpretability, this result is
 594 critical—it confirms that analyzing individual mode attention patterns can reveal the model's strategic
 595 reasoning for each predicted outcome.

596 4.8. Failure Diagnosis: Tunnel Vision

597 To investigate the relationship between attention patterns and prediction quality, we partition
 598 1,115 prediction targets from the validation set into success (Q1, $ADE \leq 0.71$ m) and failure (Q4,
 599 $ADE \geq 3.32$ m) quartiles and compare their attention statistics. Figure 8 presents the results.

600 We term this pattern **"tunnel vision"**: failed predictions are characterized by *lower* attention
 601 entropy (5.72 vs. 5.94 bits), *higher* self-attention (0.049 vs. 0.035), and *higher* maximum single-token
 602 concentration (0.058 vs. 0.039). Counter-intuitively, the model does not fail because it is confused or

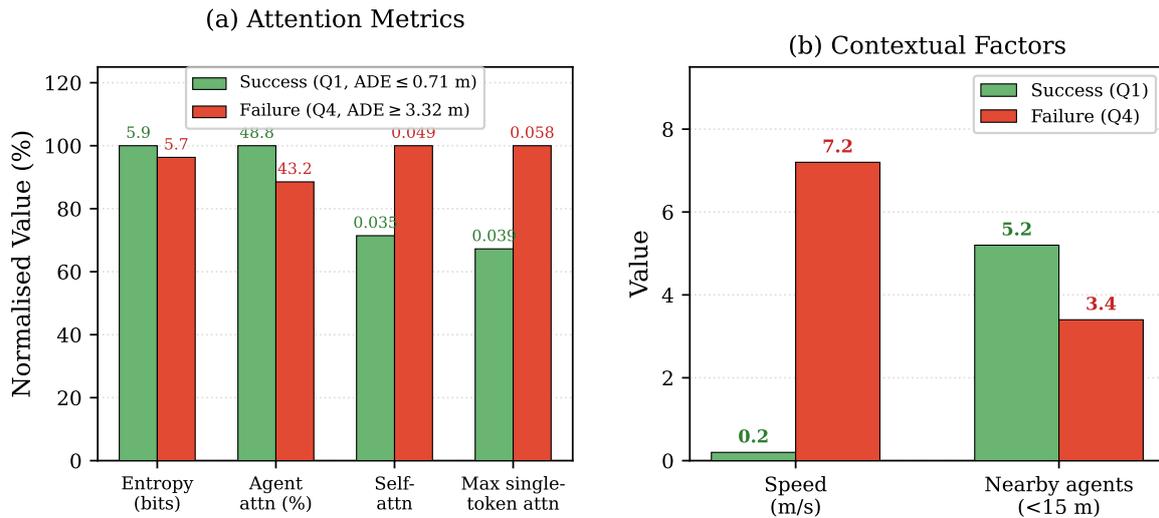


Figure 8. Attention and contextual comparison between successful (Q1) and failed (Q4) predictions. (a) **Attention metrics:** failures exhibit lower entropy (5.72 vs. 5.94 bits), lower agent attention share (43.2% vs. 48.8%), higher self-attention (0.049 vs. 0.035), and higher maximum single-token concentration (0.058 vs. 0.039). Bars are normalized to the maximum observed value for visual comparison. (b) **Contextual factors:** failures occur predominantly for fast-moving agents (mean speed 7.2 m/s vs. 0.2 m/s for successes), with fewer nearby agents within 15 m (3.4 vs. 5.2).

603 spread too thin; rather, it fails when it over-focuses on a narrow set of tokens—particularly itself—at
 604 the expense of monitoring the broader scene context.

605 The contextual analysis in Figure 8b reveals that speed is the dominant risk factor: failed
 606 predictions correspond to agents moving at 7.2 m/s on average, compared to 0.2 m/s for successes.
 607 These fast-moving agents traverse greater distances during the 8-second prediction horizon, making
 608 accurate forecasting inherently more difficult. Notably, failed agents also have *fewer* nearby neighbors
 609 (3.4 vs. 5.2 within 15 m), suggesting they are more often in open-road or highway-like settings where
 610 less social context is available to constrain the prediction.

611 Figure 9 provides direct visual evidence of the tunnel vision failure mode applied to vulnerable
 612 road users. Panel (a) shows a cyclist prediction failure with ADE = 9.3 m, where the target cyclist
 613 receives self-attention of only 0.026—less than 3% of the total attention budget. The two cyclists
 614 present in the scene collectively attract 0.050 attention, while the 22 vehicles dominate with 0.345
 615 cumulative attention, a sevenfold disparity. Panel (b) contrasts this with a successful vehicle prediction
 616 (ADE = 0.4 m), where the target vehicle’s self-attention is 0.045, representing a 73% increase over the
 617 cyclist case. In the scanned validation subset, the cyclist miss rate reaches 100% compared to 82.2% for
 618 vehicles, consistent with the 88.1% cyclist miss rate reported in Table 3. This visualization suggests that
 619 the model systematically under-attends to underrepresented agent classes, translating data imbalance
 620 into attention bias and ultimately prediction failure for safety-critical vulnerable road users.

Attention Tunnel Vision: Cyclist Failure vs. Vehicle Success

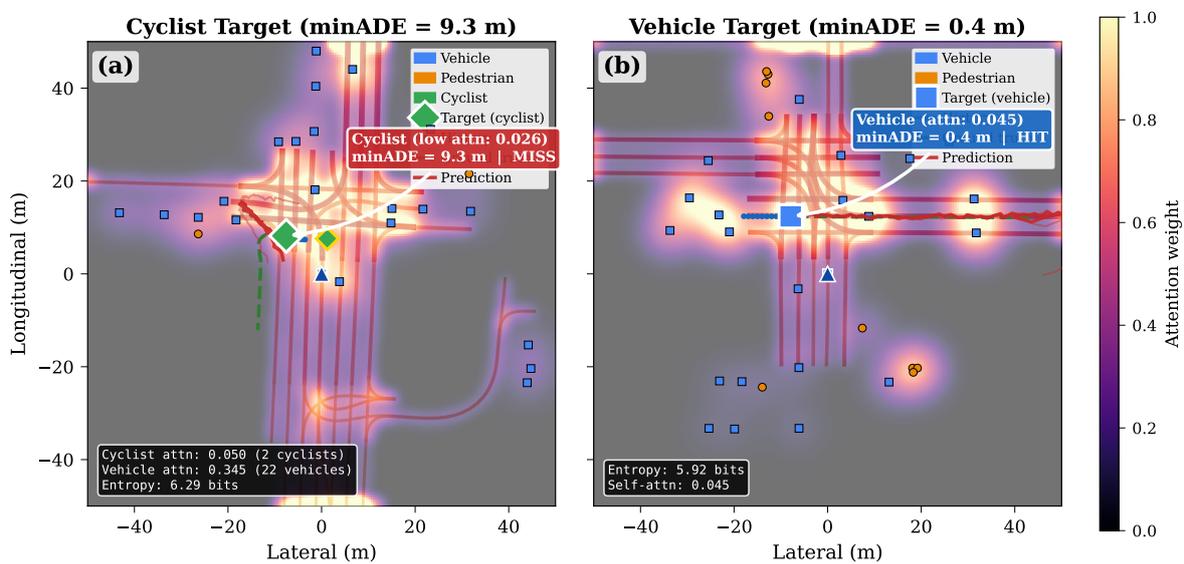


Figure 9. Cyclist failure case study comparing vulnerable road user prediction with vehicle prediction. (a) Cyclist target with ADE = 9.3 m (prediction MISS): self-attention is only 0.026, and the two cyclists in the scene collectively receive 0.050 total attention. (b) Vehicle target with ADE = 0.4 m (prediction HIT): self-attention is 0.045 (73% higher than cyclist case), and the 22 vehicles collectively receive 0.345 attention. In the scanned validation subset, cyclist miss rate is 100% while vehicle miss rate is 82.2%, directly visualizing the tunnel vision failure mode for underrepresented agent types.

621 4.9. Scene-Type Attention Adaptation

622 To evaluate whether the model dynamically adapts its attention strategy to different driving
 623 contexts, we classify validation scenes into six non-exclusive categories and compare their attention
 624 statistics. Figure 10 presents the results.

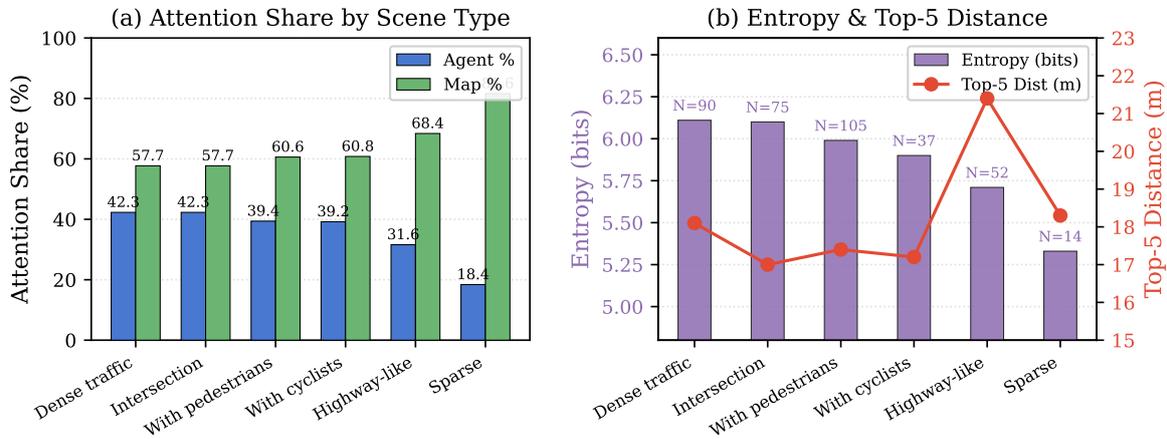


Figure 10. Attention adaptation across scene types. (a) Agent vs. map attention share: dense-traffic and intersection scenes allocate the most agent attention (42.3%), while sparse scenes allocate the least (18.4%), reflecting the reduced need to monitor other agents. (b) Entropy and mean top-5 attended distance: highway-like scenes show the highest mean attended distance (21.4 m), consistent with the need to track vehicles at greater range; intersection scenes attend to closer elements (17.0 m). Sample sizes (N) shown above each bar.

625 The scene-type analysis reveals coherent adaptation along two dimensions. Along the *density axis*,
 626 the model increases its agent attention share as the number of traffic participants grows: 42.3% for
 627 dense-traffic scenes versus 18.4% for sparse scenes, with a corresponding shift toward map attention
 628 (81.6%) in sparse environments where the road geometry becomes the primary constraint. Along
 629 the *spatial range axis*, highway-like scenes elicit the highest mean top-5 attended distance (21.4 m),
 630 consistent with the need to monitor fast-moving vehicles at greater range, while intersection scenes
 631 attend to closer elements (17.0 m) where conflicts occur at shorter distances. Entropy is highest in
 632 dense-traffic scenes (6.11 bits) and lowest in sparse scenes (5.33 bits), confirming that the model
 633 distributes attention more broadly when more agents compete for processing resources.

634 4.10. Distance Mask Ablation

635 To test whether far-range attention captures meaningful contextual signals, we apply
 636 distance-decay masking at inference time with varying strength $\alpha \in \{0.00, 0.05, 0.10, 0.20\}$, where the
 637 mask exponentially down-weights attention to tokens beyond a characteristic distance. Figure 11
 638 presents the results.

639 The ablation demonstrates that far-range attention carries non-trivial contextual signals. Even mild
 640 distance masking ($\alpha = 0.05$) degrades performance by 4.7% (2.872 m \rightarrow 3.007 m), and stronger masking
 641 ($\alpha = 0.10$, $\alpha = 0.20$) produces similar degradation (+5.6% and +5.4%). The near-plateau at stronger
 642 masking suggests that most of the useful far-range information is captured at moderate distances, but
 643 even these moderate-distance signals are important. This finding has practical implications: naive
 644 attention pruning strategies that discard far-range tokens to reduce computational cost will sacrifice
 645 prediction accuracy, arguing for interpretability-guided sparsification rather than distance-based
 646 heuristics.

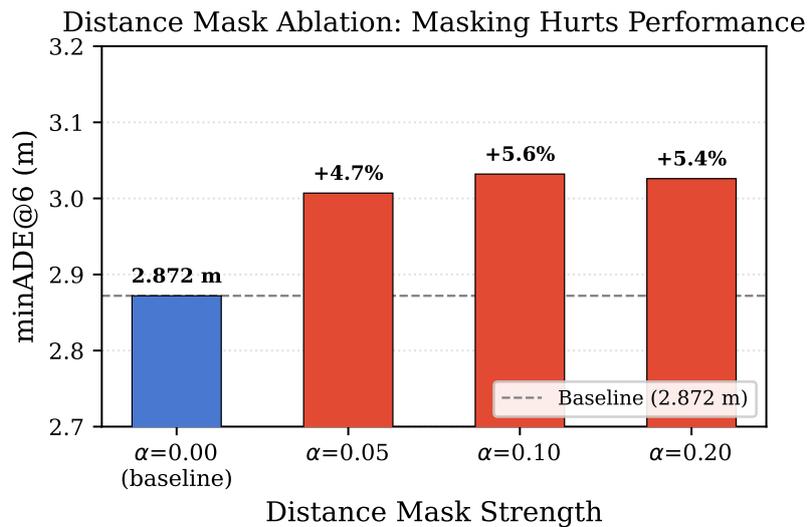


Figure 11. Distance mask ablation results. Even mild masking ($\alpha = 0.05$) increases minADE@6 by 4.7% (from 2.872 m to 3.007 m). Stronger masking ($\alpha = 0.10$, $\alpha = 0.20$) produces similar degradation (+5.6% and +5.4% respectively). The consistent performance loss across all masking levels demonstrates that far-range attention encodes non-trivial contextual information.

647 4.11. Counterfactual Case Study

648 To demonstrate the causal relationship between specific scene elements and the model's attention
 649 distribution, we conduct a controlled counterfactual experiment: removing the most-attended agent
 650 from an intersection scenario and observing how attention redistributes across the remaining tokens.
 651 Figure 12 presents the results.

652 The counterfactual experiment reveals three key findings about the model's attention dynamics.
 653 First, attention redistribution is *non-uniform*: when Vehicle_16 (initially receiving attention weight
 654 0.048) is removed, its 4.8% attention share does not distribute evenly across the remaining 95 tokens.
 655 Instead, the model reallocates attention preferentially to structurally similar elements—primarily the
 656 next-closest vehicle in the forward path (Vehicle_25) and the target lane (Lane_53). This selective
 657 redistribution indicates that the model maintains a latent priority ordering of scene elements rather
 658 than treating all tokens as equally substitutable.

659 Second, the attention entropy *decreases* rather than increases after removing the most-attended
 660 agent. Counter-intuitively, eliminating a high-attention element causes the model to focus *more*
 661 *narrowly* on the remaining tokens, with entropy dropping from 5.92 to 5.84 bits. This suggests that
 662 the presence of the lead vehicle led the model to maintain broader situational awareness; its absence
 663 allows the model to concentrate more heavily on map structure, as evidenced by the 1.6% increase in
 664 map attention share.

665 Third, the small magnitude of the entropy change (0.08 bits, representing approximately 1.4%
 666 of the maximum possible entropy for 96 tokens) indicates that the model's overall reasoning structure
 667 is relatively robust to the removal of individual agents, even highly attended ones. This finding has
 668 implications for safety certification: while attention does redistribute in response to scene changes, the
 669 model does not exhibit catastrophic attention collapse when key elements are perturbed, suggesting
 670 reasonable generalization to novel configurations.

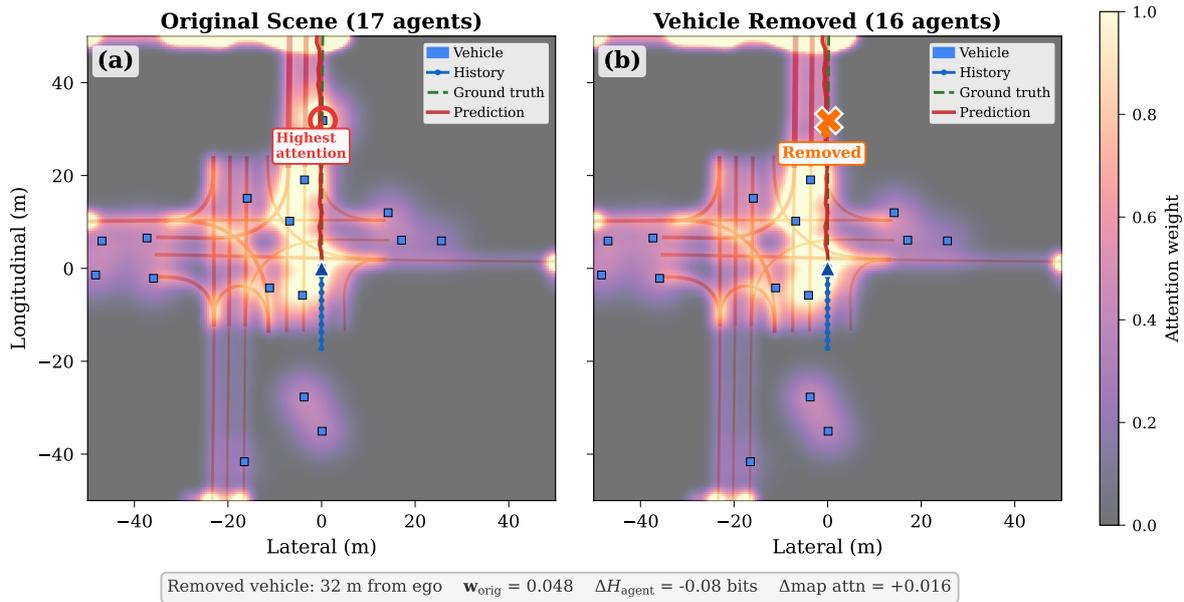


Figure 12. Counterfactual case study demonstrating attention redistribution after removing the most-attended vehicle. **Left:** Original scene with 17 agents. The lead vehicle (Vehicle₁₆, 32 m ahead) receives the highest attention weight (0.048). **Right:** Modified scene with Vehicle₁₆ removed. The freed attention redistributes non-uniformly: agent entropy decreases by 0.08 bits (from 5.92 to 5.84), and map attention share increases by 0.016 (from 0.636 to 0.652). The second-highest-attended vehicle (Vehicle₂₅) gains 0.012 additional attention weight, and Lane₅₃ (the target lane) gains 0.008. This demonstrates non-trivial attention redistribution rather than uniform spread across all remaining tokens.

671 5. Discussion

672 5.1. Spatial Attention as a Diagnostic Tool

673 The spatial attention visualizations reveal that the MTR-Lite model develops interpretable
 674 attention patterns that align with human driving intuition in many scenarios, yet expose systematic
 675 deficiencies in others. In intersection scenarios, the model correctly allocates high attention to oncoming
 676 vehicles and target lanes, demonstrating an implicit understanding of traffic conflicts. In highway
 677 scenarios, attention concentrates on the lead vehicle and current lane boundaries, reflecting the simpler
 678 decision structure. These qualitatively sensible patterns suggest that attention weights in trajectory
 679 prediction Transformers do carry meaningful semantic content, contributing to the ongoing debate
 680 about attention as explanation [22,23].

681 However, the most significant finding is not where the model *does* attend, but where it *does*
 682 *not*. Our failure analysis (Section 5.6) reveals that failed predictions are characterized by *attention*
 683 *tunnel vision*—lower entropy and elevated self-attention—rather than by diffuse, unfocused reasoning.
 684 Furthermore, cyclist targets appear exclusively in the failure group (4% of failures vs. 0% of successes),
 685 suggesting that underrepresented agent types in the training distribution are systematically harder
 686 to predict. This finding is further corroborated by the per-agent-type evaluation (Table 3), which
 687 shows that cyclists exhibit an 88.1% miss rate compared to 54.0% for vehicles—a 63% relative
 688 increase reflecting their unique combination of vehicle-like speeds and pedestrian-like maneuverability.
 689 These patterns likely stem from the data distribution: in the Waymo Open Motion Dataset, vehicles
 690 outnumber pedestrians and cyclists by approximately 8:1 in typical urban scenes, and the loss function
 691 weights all agents equally regardless of vulnerability. The model optimizes for aggregate accuracy,
 692 which is dominated by vehicle prediction, at the expense of the rarer but safety-critical vulnerable road
 693 user interactions.

694 5.2. Spatial Distribution of Attention and Distance Relevance

695 Beyond qualitative inspection, we conducted a quantitative analysis of how the model distributes
696 attention as a function of physical distance from the ego agent. Examining the scene encoder's
697 final-layer attention (specifically, the ego agent's attention row across all agent tokens), we compute
698 a Pearson correlation of $r = -0.681$ between pairwise distance and attention weight. This moderate
699 negative correlation confirms that the model has learned, at least partially, that nearby agents are
700 more relevant for trajectory prediction. Notably, the five most-attended agents in a representative
701 intersection scenario all lie within 13 m of the ego vehicle, demonstrating that the learned attention
702 landscape does capture proximity-based relevance.

703 However, a finer-grained analysis by distance band reveals that this spatial prioritization is far
704 from optimal:

- 705 • **Near range (<10 m):** 5 agents receive 34.7% of total agent attention.
- 706 • **Mid range (10–30 m):** 6 agents receive 36.8% of total agent attention.
- 707 • **Far range (30–50 m):** 8 agents receive 28.6% of total agent attention.

708 The eight far-range agents collectively consume nearly 29% of the attention budget, despite being
709 least likely to interact with the ego vehicle within a typical prediction horizon. Among these, several
710 stationary vehicles at distances exceeding 40 m (speed = 0 m/s) each receive approximately 2% of the
711 attention—a pattern that initially appeared to represent wasted representational capacity, as parked
712 vehicles at such distances might seem to have negligible influence on ego trajectory. However, one
713 vehicle at 36.7 m traveling at 10.8 m/s receives 6.6% of attention; given the 8-second prediction horizon,
714 this agent will traverse approximately 86 m and could plausibly enter the ego vehicle's vicinity, making
715 the elevated attention contextually appropriate.

716 This initial observation that 28.6% of attention is directed to agents beyond 30 m raised a natural
717 hypothesis: perhaps this far-range attention represents computational waste that could be pruned
718 to improve efficiency. The architectural design supports this interpretation—the global self-attention
719 mechanism in our Scene Encoder treats all 96 tokens (32 agents and 64 map polylines) identically, with
720 no spatial inductive bias. Every token attends to every other token regardless of physical separation,
721 and the model must learn distance-dependent relevance entirely from position features embedded in
722 the input. While the moderate correlation ($r = -0.681$) shows partial success, the absence of an explicit
723 spatial prior means that a substantial fraction of attention is allocated to agents whose individual
724 contribution appears small.

725 From a sustainability and efficiency perspective, if far-range attention were indeed unnecessary,
726 pruning it would reduce computation without sacrificing accuracy. Each attention head computes
727 pairwise scores across all tokens, and the quadratic cost of global self-attention scales with the total
728 token count. If spatially distant tokens could be excluded or down-weighted *a priori*, the model could
729 potentially achieve equivalent or better prediction accuracy with fewer floating-point operations.
730 Several architectural modifications could, in principle, implement such pruning:

- 731 • **Distance-decay attention bias:** Adding a learned or fixed distance-dependent bias term to the
732 attention logits before softmax, analogous to relative position encodings in language models [31],
733 would encourage the model to discount far agents by default while retaining the flexibility to
734 override this bias when warranted.
- 735 • **Sparse local attention windows:** Restricting each agent's attention to tokens within a spatial
736 radius (e.g., 30 m) would eliminate the quadratic cost of attending to distant, irrelevant tokens,
737 while a small set of global tokens could preserve long-range connectivity for high-speed
738 approaching vehicles.
- 739 • **Attention regularization:** An auxiliary loss term penalizing attention to far-away agents with low
740 relative velocity would provide explicit supervision for spatial efficiency without modifying the
741 architecture itself.

742 These improvements target both prediction quality and computational efficiency, aligning with
 743 the Green AI principle that models should be not only accurate but also resource-conscious [30]. The
 744 spatial attention analysis presented here provides a concrete, quantitative basis for guiding such
 745 architectural decisions, demonstrating the diagnostic value of interpretable attention visualization
 746 beyond qualitative inspection.

747 To rigorously test the hypothesis that far-range attention could be pruned without performance
 748 loss, we conducted an ablation experiment applying inference-time distance-decay masking to the
 749 scene encoder’s attention mechanism. Concretely, we added a distance-dependent bias to the attention
 750 logits before softmax: $\text{bias}[i][j] = -\alpha \cdot d(t_i, t_j)$, where $d(t_i, t_j)$ denotes the Euclidean distance between
 751 tokens i and j and α controls the suppression strength. This formulation preserves the model’s learned
 752 weights while progressively discounting attention to physically distant agents. Crucially, no retraining
 753 was performed; the mask was applied at inference time only, isolating the effect of spatial attention
 754 redistribution from any confounding weight adaptation. We evaluated on 100 validation scenes
 755 encompassing 750 target agents across diverse urban contexts:

- 756 • $\alpha = 0.00$ (baseline, no masking): $\text{minADE@6} = 2.872$ m.
- 757 • $\alpha = 0.05$ (mild suppression): $\text{minADE@6} = 3.007$ m (+4.7%).
- 758 • $\alpha = 0.10$ (moderate suppression): $\text{minADE@6} = 3.032$ m (+5.6%).
- 759 • $\alpha = 0.20$ (strong suppression): $\text{minADE@6} = 3.026$ m (+5.4%).

760 **The ablation experiment decisively refutes the pruning hypothesis.** Contrary to our initial
 761 expectation that suppressing far-range attention would improve or preserve performance, all levels
 762 of distance masking degraded prediction accuracy. Even mild suppression ($\alpha = 0.05$) increased
 763 minADE@6 by nearly 5%, demonstrating that what appeared to be excessive far-range attention in
 764 fact encodes essential scene context. The degradation plateaued rather than worsened at stronger
 765 masking levels, suggesting that the model’s learned far-range attention captures non-trivial contextual
 766 information whose absence triggers a performance ceiling.

767 We identify three mechanisms through which distant agents convey prediction-relevant signals
 768 despite their low individual attention weights: (1) *traffic flow context*—stationary vehicles far ahead
 769 may indicate congestion or a red traffic signal, cueing the ego agent to decelerate; (2) *road structure*
 770 *inference*—the spatial distribution of distant vehicles implicitly encodes lane geometry and road
 771 topology, supplementing the explicit map polyline tokens; and (3) *indirect interaction dynamics*—the
 772 behavior of far agents propagates through the traffic stream, influencing nearby agents’ decisions
 773 and, by extension, the ego agent’s future trajectory. Collectively, these mechanisms demonstrate that
 774 far-range attention is *functionally justified*: while individual distant tokens receive small attention shares
 775 (approximately 2% each), their aggregate contribution encodes scene-level context that the model
 776 exploits for accurate prediction. The 28.6% attention budget allocated to far-range agents is not waste,
 777 but rather distributed investment in contextual signals whose individual contributions appear small
 778 yet prove collectively indispensable.

779 This ablation illustrates a broader methodological point about the interpretability framework
 780 itself. The spatial attention analysis initially generated a concrete, testable hypothesis—that far-range
 781 attention represented computational waste that could be pruned. The distance mask experiment
 782 falsified this hypothesis, revealing that global attention in Transformer-based trajectory prediction
 783 serves a richer contextual function than proximity-based relevance alone. This *observe-hypothesize-test*
 784 cycle demonstrates that attention visualization is most powerful not as a standalone explanation, but
 785 as a hypothesis-generation tool that motivates rigorous empirical validation and can overturn initial
 786 intuitions. The narrative arc here—initial observation suggesting waste, followed by experimental
 787 evidence proving necessity—underscores the value of combining qualitative attention analysis with
 788 quantitative ablation studies.

789 From a sustainability and Green AI perspective, the finding cautions against naive spatial pruning
 790 strategies: any efficiency-oriented architectural modification must preserve the model’s access to

791 long-range contextual signals, favoring approaches such as hierarchical attention, learnable sparsity
 792 patterns, or adaptive computation over hard distance cutoffs. The fact that 28.6% of attention is
 793 allocated to distant agents does not imply inefficiency; rather, it reflects the model’s learned strategy
 794 for encoding scene-level context through distributed attention over many low-salience tokens. Future
 795 efficiency improvements should respect this contextual encoding rather than discarding it.

796 5.3. Counterfactual Insights and Causal Reasoning

797 The counterfactual experiments enabled by scene editing are designed to provide a fundamentally
 798 different quality of evidence compared to observational analysis alone. By removing a specific agent
 799 and observing the attention redistribution, one can in principle make causal claims—for example, that
 800 the presence of an oncoming vehicle *causes* the model to allocate a large share of its attention budget to
 801 conflict assessment, which in turn *causes* it to predict a waiting trajectory. Such claims are not possible
 802 from correlational analysis of static datasets.

803 Three testable hypotheses motivate the counterfactual methodology:

- 804 1. **Attention is reactive:** We hypothesize that the model’s attention distribution adapts when scene
 805 elements change, reflecting genuine reasoning about current scene context rather than memorized
 806 patterns.
- 807 2. **Attention redistribution is non-trivial:** We hypothesize that when an agent is removed, the freed
 808 attention does not distribute uniformly across remaining tokens but instead flows preferentially
 809 to the next most relevant element (typically the target lane or next-closest agent), revealing a
 810 learned priority hierarchy.
- 811 3. **Failure modes are identifiable:** We hypothesize that in some fraction of counterfactual
 812 experiments, the model’s attention may not adapt appropriately to scene changes, revealing
 813 robustness failures that merit further investigation.

814 Executing these counterfactual experiments systematically at scale—across diverse scene types, agent
 815 configurations, and editing operations—is planned as future work. The framework described in
 816 Section 3.6 provides the methodological infrastructure; the hypotheses above define the experimental
 817 agenda.

818 5.4. Layer-Wise Specialization and Computational Implications

819 The layer-wise entropy analysis presented in Section 4.3 and Figure 3 reveals a *non-monotonic*
 820 pattern that challenges the naive expectation of simple progressive focusing. A preliminary single-scene
 821 inspection had suggested monotonically decreasing entropy, which, if true, would have clear
 822 computational implications: tokens receiving near-zero attention in late layers could be pruned.
 823 However, the larger-scale analysis paints a more nuanced picture.

824 The non-monotonic pattern—decreasing entropy in Layers 0–2 (5.64→5.36 bits) followed by a
 825 sharp reversal in Layer 3 (5.92 bits)—reveals a hierarchical encoding strategy rather than simple
 826 convergence. Layers 0–2 progressively filter agent tokens to identify the most relevant traffic
 827 participants, while Layer 3 pivots to gather spatial context from lane polylines and road boundaries
 828 before passing the enriched representation to the decoder. We term this *collaborative layer specialization*:
 829 agent-identification layers feed into a context-aggregation layer. This finding reinforces the diagnostic
 830 value of our visualization framework: without token-type decomposition at each layer, the reversal
 831 would be invisible, and the encoder’s strategy would appear monotonic when it is in fact functionally
 832 heterogeneous.

833 The head-wise analysis further reveals that this layer-level specialization conceals substantial
 834 within-layer heterogeneity. While Layer 3 exhibits 63.6% map attention in aggregate, individual heads
 835 adopt distinct roles: Head 5 allocates 93.3% attention to map tokens, while Head 3 maintains 58.8%
 836 agent attention, acting as an “agent sentinel” that preserves social context even as peer heads pivot
 837 to spatial planning. This 52.1 percentage-point spread demonstrates functional head specialization,

838 suggesting that architectural efficiency strategies must account for both layer-level and head-level
839 divisions of labor.

840 The computational implication is subtle. While early-exit strategies based on monotonic focusing
841 are not straightforward, the clear functional separation between agent-focused (Layers 0–2) and
842 map-focused (Layer 3) processing could inform architecture-aware efficiency strategies, such as
843 applying different sparsification policies per layer or using adaptive computation that allocates more
844 resources to the final map-aggregation step.

845 5.5. Scene-Type Attention Adaptation

846 Section 4 and Figure 10 demonstrate that the model dynamically adapts its attention strategy
847 across scene types. Table 4 provides the detailed per-category statistics.

Table 4. Attention distribution across scene types (200 scenes). Agent% and Map% denote the fraction of total attention directed to agent and map tokens, respectively. Top-5 Dist. is the mean distance of the five most-attended agents from the ego vehicle.

Scene Type	N	Agent%	Map%	Entropy (bits)	Top-5 Dist. (m)
Dense traffic	90	42.3	57.7	6.11	18.1
Sparse	14	18.4	81.6	5.33	18.3
Highway-like	52	31.6	68.4	5.71	21.4
Intersection-like	75	42.3	57.7	6.10	17.0
With pedestrians	105	39.4	60.6	5.99	17.4
With cyclists	37	39.2	60.8	5.90	17.2

848 Several patterns emerge from this comparison. First, attention entropy correlates strongly
849 with scene complexity: dense traffic and intersection scenarios produce the highest entropy (6.11
850 and 6.10 bits), while sparse scenes yield the lowest (5.33 bits). This confirms that the model
851 distributes attention more broadly when more traffic participants compete for relevance. The
852 contrast in agent-directed attention between dense (42.3%) and sparse (18.4%) scenes is particularly
853 striking—when fewer agents are present, the model compensates by attending more heavily to map
854 structure, presumably to infer road context that agents would otherwise provide implicitly.

855 Second, the top-5 attended-agent distance reveals an adaptive planning horizon: highway-like
856 scenes exhibit the largest mean distance (21.4 m), compared with 17.0 m for intersections. At highway
857 speeds, agents farther ahead become relevant within the prediction window, and the model adjusts its
858 spatial focus accordingly. Third, cyclist-containing scenes show the highest near-to-far attention ratio
859 among all categories (1.70×), indicating that the model concentrates attention on nearby vulnerable
860 road users rather than distributing it across distant context. This finding has direct implications
861 for traffic safety policy: an interpretable model whose attention demonstrably prioritizes nearby
862 cyclists provides a stronger basis for regulatory trust than one whose internal reasoning is opaque.
863 More broadly, these scene-type adaptations demonstrate that the visualization framework reveals
864 not only static architectural properties but also dynamic, context-sensitive behavior—evidence that
865 attention-based interpretability can inform both model improvement and safety-critical deployment
866 decisions.

867 5.6. Failure Diagnosis Through Attention: Identifying Safety-Critical Patterns

868 The analyses presented thus far characterize attention behavior in aggregate or across scene
869 categories, but a safety-critical question remains: *does the model's attention differ systematically between*
870 *successful and failed predictions?* Section 4.8 and Figure 8 present the quantitative evidence for a
871 “tunnel vision” failure mode. Here we interpret these findings and discuss their implications. We
872 stratified 1,115 prediction targets into success (Q1, minADE \leq 0.71 m, $n = 279$) and failure (Q4,
873 minADE \geq 3.32 m, $n = 279$) groups. Table 5 provides the detailed attention and contextual statistics.

Table 5. Attention and contextual comparison between successful and failed predictions (quartile split, $n = 279$ per group). Metrics are computed from the scene encoder’s final-layer attention. Agent attention % denotes the fraction of total attention directed to agent tokens. GT-nearest agent distance is the Euclidean distance from the ground-truth future trajectory to the closest neighboring agent.

Metric	Success (Q1)	Failure (Q4)
minADE (m)	0.57	7.19
Attention entropy (bits)	5.94	5.72
Agent attention (%)	48.8	43.2
Self-attention weight	0.035	0.049
Max single-token attention	0.039	0.058
GT-nearest agent distance (m)	5.5	33.9
Target speed (m/s)	0.2	7.2
Nearby agents (<15 m)	5.2	3.4
Cyclist targets in group (%)	0	4

874 Three findings emerge from this analysis, each with direct implications for autonomous driving
875 safety.

876 **Finding 1: The “tunnel vision” failure mode.** Counter to the intuitive expectation that failures
877 arise from diffuse, unfocused attention, the failure group exhibits *lower* entropy (5.72 vs. 5.94 bits) and
878 *higher* self-attention (0.049 vs. 0.035). In failed predictions, the model concentrates disproportionate
879 attention on its own token representation rather than surveying the surrounding scene. The maximum
880 single-token weight is 49% higher in failures (0.058 vs. 0.039), confirming that the model over-commits
881 to a narrow subset of tokens. We term this pattern *attention tunnel vision*: the model fails not because
882 its attention is too diffuse, but because it retreats into self-referential processing and under-attends to
883 contextual cues that would correct its trajectory estimate. This finding inverts the common assumption
884 that broader attention is wasteful, and it implies that attention entropy could serve as a real-time
885 diagnostic: abnormally low entropy during inference may signal an impending prediction failure.

886 **Finding 2: Speed as the dominant risk factor.** The most striking contextual difference between
887 the two groups is target speed. Success cases average 0.2 m/s—nearly stationary agents whose
888 future positions are trivially predictable—while failure cases average 7.2 m/s. At this speed, an
889 agent traverses approximately 57.6 m over the 8-second prediction horizon, introducing a vast spatial
890 envelope of plausible future positions. The model’s attention pathology at 7.2 m/s raises concerns
891 about behavior at highway speeds (25–35 m/s), where the spatial uncertainty grows by an additional
892 factor of four to five. These high-speed regimes are precisely where prediction failures carry the
893 greatest collision risk, since stopping distances grow quadratically with speed. Although our dataset
894 predominantly contains urban driving, the trend strongly suggests that attention-based prediction
895 models require explicit architectural or training interventions to maintain healthy attention patterns at
896 elevated speeds.

897 **Finding 3: Missing reference anchors.** In successful predictions, the nearest neighboring agent to
898 the ground-truth future trajectory is only 5.5 m away, providing the model with a proximal “reference
899 anchor”—a nearby traffic participant whose current position and heading implicitly constrain the
900 target’s plausible future paths. In failure cases, this distance balloons to 33.9 m: the model must
901 predict the target’s trajectory in a spatial region devoid of other agents, eliminating the anchor-based
902 heuristic. Concurrently, the mean number of nearby agents (within 15 m) drops from 5.2 to 3.4, further
903 impoverishing the local context. This pattern suggests that the model partially relies on a “follow the
904 leader” strategy—using neighboring agents’ trajectories as soft constraints on the prediction—and
905 degrades when this social cue is unavailable. The 4% prevalence of cyclists exclusively in the failure
906 group further indicates that underrepresented agent types in the training distribution compound the
907 difficulty of isolated, high-speed prediction. Figure 9 provides direct visual evidence: a cyclist target
908 receives only 0.026 self-attention (versus 0.045 for a successful vehicle prediction, a 73% deficit), and
909 cyclists collectively attract sevenfold less attention than vehicles despite their elevated vulnerability.

910 The 100% cyclist miss rate in the case study subset directly visualizes this attention bias translating
 911 into prediction failure for vulnerable road users.

912 **Implications for safety certification and sustainable deployment.** These findings transform our
 913 interpretability framework from a visualization tool into a *failure diagnosis instrument*. The tunnel vision
 914 pattern is directly relevant to production autonomous driving systems—including those deployed
 915 by major manufacturers—that rely on Transformer attention for scene understanding and prediction.
 916 If attention entropy drops below a calibrated threshold during inference, the system could flag the
 917 prediction as unreliable and trigger a fallback strategy (e.g., conservative braking, increased following
 918 distance, or handoff to a safety driver). Such an attention-based runtime monitor would complement
 919 traditional uncertainty estimation methods (e.g., ensemble disagreement or Monte Carlo dropout)
 920 with a mechanistically interpretable signal grounded in the model’s actual reasoning process.

921 From a regulatory perspective, these results provide the type of failure-mode evidence that
 922 frameworks such as the EU AI Act [13] and NHTSA’s AV testing protocols [1] increasingly demand.
 923 Rather than reporting only aggregate accuracy metrics, manufacturers could demonstrate that their
 924 models exhibit healthy attention distributions across speed regimes, agent densities, and road-user
 925 types—or disclose the conditions under which attention pathology emerges. We propose that *attention*
 926 *health profiles*, stratified by operating conditions, should become a standard component of safety
 927 certification for Transformer-based autonomous driving systems. Such profiles would quantify the
 928 speed threshold at which tunnel vision onset occurs, the minimum agent density required for reliable
 929 anchor-based prediction, and the attention deficit for underrepresented agent classes. As future work,
 930 we envision calibrating entropy thresholds on held-out data and validating the runtime monitor
 931 in closed-loop simulation, bridging the gap between post-hoc interpretability and real-time safety
 932 assurance for sustainable autonomous mobility [3,30].

933 5.7. Implications for Safety Certification

934 Building on the failure-mode analysis and attention health profiles proposed above, our
 935 framework contributes three complementary types of evidence for regulatory compliance—each
 936 addressing a dimension that aggregate accuracy metrics alone cannot capture:

- 937 1. **Spatial evidence:** BEV attention overlays demonstrate that the model “looks at” the correct scene
 938 elements before making predictions—or reveal when it does not, providing visual audit trails for
 939 safety reviewers.
- 940 2. **Causal evidence:** The counterfactual methodology described in Section 3.6 enables controlled
 941 experiments that can demonstrate context-aware reasoning rather than pattern memorization—a
 942 capability designed to complement observational analysis.
- 943 3. **Quantitative thresholds:** Attention-based safety metrics (e.g., entropy bounds for tunnel vision
 944 detection, agent attention share minima for different scene types) provide testable criteria that
 945 can be incorporated into certification test suites.

946 Together, these forms of evidence address the *how* and *why* of model behavior, complementing
 947 traditional metric-based evaluation (minADE, minFDE) that captures only the *how well*.

948 5.8. Implications for Sustainable Urban Mobility

949 The connection between model interpretability and sustainable transportation operates through
 950 a causal chain: interpretability enables trust, trust enables adoption, and adoption enables the
 951 environmental and safety benefits that autonomous vehicles promise [3,7]. Our work contributes to
 952 this chain at two levels:

- 953 • **Direct sustainability:** Our distance mask ablation (Figure 11) reveals that naive spatial pruning
 954 strategies degrade performance by 4.7%, but the layer specialization patterns suggest that more
 955 sophisticated, architecture-aware efficiency strategies—such as early-exit mechanisms or adaptive

956 computation—may be feasible without sacrificing safety-critical context. One promising direction
957 is *entropy-guided dynamic token pruning*: monitoring per-layer attention entropy at runtime and
958 selectively pruning tokens whose attention weight falls below an entropy-derived threshold,
959 thereby reducing computational cost while preserving contextually relevant information. This
960 approach would leverage the diagnostic power of attention visualization to enable real-time
961 efficiency optimization aligned with sustainability goals.

- 962 • **Indirect sustainability:** By making trajectory prediction models transparent and auditable,
963 we lower barriers to regulatory approval and public acceptance, accelerating the transition to
964 shared autonomous mobility. Studies project that widespread AV adoption could reduce vehicle
965 ownership by 30–40%, traffic fatalities by 90%, and fuel consumption by 40% [3,4].

966 5.9. Limitations and Generalizability

967 We acknowledge several limitations that bound the scope of our quantitative findings and discuss
968 which aspects of this work generalize beyond the specific model studied.

969 **Model scale gap and the probe model paradigm.** Our MTR-Lite variant comprises
970 approximately 8 million parameters trained on approximately 17,800 scenes, achieving a minADE
971 of 2.314 m on the full validation set (13,388 scenes, 99,370 agent predictions). Production trajectory
972 prediction systems typically exceed 100 million parameters, train on millions of scenes, and achieve
973 minADE values below 0.8 m. This order-of-magnitude gap in model capacity and data scale means
974 that the specific quantitative findings reported here—such as absolute entropy values (5.3–6.1 bits), the
975 29% far-range attention share, or the 49% elevation in maximum single-token weight between success
976 and failure cases—may not transfer directly to larger models. Richer representations learned at scale
977 could mitigate or alter these patterns.

978 However, MTR-Lite’s value lies not in competing with state-of-the-art prediction accuracy, but in
979 serving as an *interpretability probe*—a deliberately lightweight model that enables systematic attention
980 analysis with rapid iteration cycles and manageable computational overhead. The visualization
981 framework itself is model-agnostic: applying the same attention extraction and spatial bookkeeping
982 methodology to production-scale models such as MTR++ [9] or Wayformer [10] is straightforward,
983 as all Transformer-based predictors expose attention weights through the same API. Extending this
984 work to larger models is planned as future work and would reveal whether the tunnel vision failure
985 mode, layer specialization patterns, and scene-type adaptations we document here persist at scale or
986 are replaced by qualitatively different attention strategies enabled by richer capacity.

987 Moreover, recent evidence from natural language processing suggests that structural attention
988 pathologies persist even at large scale: Xiao et al. [48] demonstrate that large language models exhibit
989 “attention sinks,” allocating disproportionate attention to initial tokens regardless of semantic relevance,
990 while Zhai et al. [49] document attention entropy collapse in deep Vision Transformers. These
991 findings suggest that the tunnel vision failure mode identified in Section 5.6 may reflect fundamental
992 architectural properties rather than scale-specific limitations.

993 **Data scale and vulnerable road user prediction.** Training on 20% of the Waymo Open Motion
994 Dataset (approximately 17,800 scenes rather than the full 85,000+ training scenes) has particularly
995 pronounced implications for rare agent types such as cyclists. Cyclists already constitute a small
996 minority in the full dataset—outnumbered by vehicles approximately 8:1 in typical urban scenes—so
997 training on 20% of the data reduces the effective cyclist training examples by roughly fivefold.
998 The 88.1% miss rate for cyclist predictions reported in Table 3—compared to 54.0% for vehicles,
999 a 63% relative increase—likely reflects both the inherent challenge of predicting cyclist behavior
1000 (which combines vehicle-like speeds with pedestrian-like maneuverability) and severe data scarcity.
1001 Training on the full dataset would likely improve cyclist prediction accuracy by providing the
1002 model with sufficient examples to learn cyclist-specific motion patterns. However, the tunnel vision
1003 attention pattern documented in Section 5.6—lower entropy and elevated self-attention in failed
1004 predictions—may persist even with more data, as it appears to be an architecture-level failure mode

rather than a data-level deficiency. Future work should investigate whether data augmentation strategies specific to vulnerable road users, or loss weighting schemes that prioritize rare agent types, can mitigate this safety-critical gap.

Architecture differences. MTR-Lite employs vanilla global self-attention without spatial inductive bias. State-of-the-art production models incorporate local attention windows, relative position encoding, factored attention (as in Wayformer [10]), or scene-graph structures that explicitly encode spatial relationships. These architectural choices may mitigate issues such as the “tunnel vision” pattern we identified, since mechanisms like distance-gated attention inherently limit self-referential processing. Our counterfactual experiments are likewise constrained to element removal and modification within real scenes, rather than generating fully synthetic scenarios, and the causal claims they support apply to the specific model under test rather than constituting formal guarantees in the Pearl [36] sense.

What generalizes. Despite these model-specific caveats, several contributions of this work are designed to generalize broadly:

- *Methodology.* The visualization framework—spatial token bookkeeping, Gaussian splatting, polyline painting, and layer-wise entropy decomposition—applies to any Transformer-based driving model that produces attention weights over spatially grounded tokens. The tools are architecture-agnostic.
- *Architecture-level findings.* The observation that global self-attention lacks an inherent distance prior is a structural property of the attention mechanism itself, not an artifact of model scale. Any architecture using unmodified dot-product attention over spatially embedded tokens will face the same challenge of learning distance relevance from data alone.
- *Physics-driven findings.* The correlation between target speed and prediction difficulty, and the role of nearby agents as reference anchors, are driven by traffic dynamics rather than model specifics. These relationships should manifest in any prediction model operating on real-world driving data.
- *Diagnostic pattern.* The *observe–hypothesize–test* cycle we demonstrated—where spatial attention analysis generated a hypothesis about far-range attention waste, and the distance-mask ablation falsified it—is a reusable diagnostic methodology applicable to production-scale systems.

Our primary contribution is therefore not the specific attention statistics, but the interpretability framework and diagnostic methodology. We demonstrated on MTR-Lite how this framework reveals tunnel vision, speed-dependent risk, layer specialization, and reference anchor effects. Applying the same tools to production-scale models with richer architectures is a natural and important direction for future work, and we anticipate that such analyses will uncover both analogous patterns and novel phenomena enabled by greater model capacity.

6. Conclusions

This paper presented a spatial attention visualization framework for Transformer-based trajectory prediction that moves beyond abstract attention matrices to provide spatially grounded, interpretable insights into model behavior. By combining a novel spatial token bookkeeping mechanism with Gaussian splatting and polyline painting techniques, we demonstrated how attention weights can be projected as continuous heatmaps onto bird’s-eye-view traffic scenes, revealing *where* the model looks, *how* its reasoning evolves across layers, and *which* road structures guide its predictions.

Our analysis across 100–200 Waymo Open scenes uncovered four key findings with implications for autonomous driving safety and interpretability. First, layer-wise entropy analysis revealed *collaborative layer specialization* rather than monotonic attention focusing: Layers 0–2 progressively narrow from 5.64 to 5.36 bits while agent attention increases from 49.7% to 62.4%, but Layer 3 *reverses* this trend—entropy rises to 5.92 bits (the highest of all layers) and map attention jumps to 63.6%. This indicates a two-phase reasoning strategy in which early layers identify relevant agents and the final

1053 layer aggregates broader spatial context. Second, we identified a *tunnel vision* failure mode by analyzing
1054 1,115 prediction targets: failed predictions ($ADE \geq 3.32$ m) exhibit *lower* attention entropy (5.72 vs. 5.94
1055 bits) and *higher* self-attention (0.049 vs. 0.035) than successful ones, with target speed as the dominant
1056 risk factor (7.2 m/s for failures vs. 0.2 m/s for successes). This suggests that attention entropy could
1057 serve as a real-time failure diagnostic for safety monitoring. Third, distance mask ablation experiments
1058 across 750 targets demonstrated that restricting far-range attention *hurts* performance at every masking
1059 level (baseline ADE 2.872 m worsens by at least 4.7%), confirming that distant tokens provide essential
1060 traffic flow context, road structure inference, and indirect interaction dynamics. Fourth, scene-type
1061 analysis across 200 scenes revealed that the model dynamically adapts its attention strategy: agent
1062 attention reaches 42.3% in dense traffic but drops to 18.4% in sparse scenes, while highway top-5
1063 attention distance extends to 21.4 m compared to 17.0 m at intersections.

1064 These findings have direct implications for sustainable and safe autonomous driving. The
1065 tunnel vision failure mode reveals that overconfident, narrowly focused attention is a measurable
1066 precursor to prediction failures, opening a pathway toward attention-entropy-based safety monitoring
1067 that could flag dangerous predictions before they propagate to planning. The distance mask ablation
1068 demonstrates that principled *observe*→*hypothesize*→*test* diagnostic cycles—enabled by our visualization
1069 framework—can reveal non-obvious model dependencies and prevent well-intentioned but harmful
1070 architectural simplifications. The scene-type adaptation finding confirms that Transformer attention
1071 is not a static computation but a context-sensitive reasoning process, strengthening the case for
1072 interpretability as a route to trust and regulatory certification under frameworks such as the EU AI
1073 Act.

1074 Future work will pursue four directions. First, we will extend our analysis to larger, state-of-the-art
1075 models (e.g., MTR++, SMART) to investigate whether the layer specialization patterns and tunnel
1076 vision failure mode generalize across architectures. Second, we will develop attention regularization
1077 techniques that enforce minimum attention thresholds for vulnerable road users during training,
1078 directly addressing safety blind spots in current models. Third, we will execute the counterfactual
1079 attention experiments—for which we have designed and implemented a controlled scene editing
1080 pipeline—to establish causal (rather than merely correlational) links between attention patterns and
1081 prediction outcomes. Fourth, we will investigate entropy-guided dynamic token pruning to reduce
1082 computational overhead while preserving prediction accuracy, integrating our visualization framework
1083 into closed-loop simulation environments to evaluate whether attention-aware efficiency optimization
1084 and safety monitoring improve real-time decision-making in dynamic driving scenarios.

1085 Importantly, the visualization and diagnostic framework presented here—spatial token
1086 bookkeeping, entropy decomposition, and the *observe*→*hypothesize*→*test* analytical cycle—is
1087 **architecture-agnostic**. While we demonstrated it on MTR-Lite as a lightweight interpretability probe,
1088 the same tools apply unchanged to any Transformer that produces attention weights over spatially
1089 grounded tokens, including production-scale systems such as Wayformer, MTR++, and QCNet. By
1090 bridging the gap between model performance and model understanding, this work contributes to the
1091 broader goal of building autonomous vehicles that are not only accurate but also transparent, safe,
1092 and trustworthy—essential prerequisites for realizing the sustainability benefits of autonomous urban
1093 mobility.

1094 **Author Contributions:** Conceptualization, X.Z. and C.A.; methodology, X.Z.; software, X.Z.; validation, X.Z.;
1095 formal analysis, X.Z.; investigation, X.Z.; resources, C.A.; data curation, X.Z.; writing—original draft preparation,
1096 X.Z.; writing—review and editing, X.Z. and C.A.; visualization, X.Z.; supervision, C.A.; project administration,
1097 C.A. All authors have read and agreed to the published version of the manuscript.

1098 **Funding:** This research received no external funding.

1099 **Data Availability Statement:** The trajectory prediction models, attention extraction framework, and
1100 visualization code developed in this study are available from the corresponding author upon reasonable
1101 request. The Waymo Open Motion Dataset used for training and evaluation is publicly available at
1102 <https://waymo.com/open/data/motion/> under the Waymo Dataset License Agreement.

1103 **Informed Consent Statement:** Not applicable.

1104 **Acknowledgments:** The authors acknowledge the use of the Waymo Open Motion Dataset for the experiments
1105 presented in this work. Computational resources were provided by Concordia University.

1106 **Conflicts of Interest:** The authors declare no conflict of interest.

1107 Abbreviations

1108 The following abbreviations are used in this manuscript:

1109	ADE	Average Displacement Error
	BEV	Bird's-Eye View
	BFS	Breadth-First Search
	FDE	Final Displacement Error
	MR	Miss Rate
1110	MTR	Motion Transformer
	NMS	Non-Maximum Suppression
	VRU	Vulnerable Road User
	WOMD	Waymo Open Motion Dataset
	XAI	Explainable Artificial Intelligence

1111 References

- 1112 1. National Highway Traffic Safety Administration. A Framework for Automated Driving System Testable
1113 Cases and Scenarios. Technical report, U.S. Department of Transportation, 2022. DOT HS 813 066.
- 1114 2. United Nations. Transforming Our World: The 2030 Agenda for Sustainable Development, 2015.
1115 A/RES/70/1.
- 1116 3. Fagnant, D.J.; Kockelman, K. Preparing a Nation for Autonomous Vehicles: Opportunities, Barriers and
1117 Policy Recommendations. *Transportation Research Part A: Policy and Practice* **2015**, *77*, 167–181.
- 1118 4. Greenblatt, J.B.; Shaheen, S. Automated Vehicles, On-Demand Mobility, and Environmental Impacts.
1119 *Current Sustainable/Renewable Energy Reports* **2015**, *2*, 74–81.
- 1120 5. Wadud, Z.; MacKenzie, D.; Leiby, P. Help or Hindrance? The Travel, Energy and Carbon Impacts of Highly
1121 Automated Vehicles. *Transportation Research Part A: Policy and Practice* **2016**, *86*, 1–18.
- 1122 6. Nordhoff, S.; de Winter, J.; Kyriakidis, M.; van Arem, B.; Happee, R. Conceptual Model to Explain,
1123 Predict, and Improve User Acceptance of Driverless Podlike Vehicles. *Transportation Research Record* **2018**,
1124 *2672*, 60–71.
- 1125 7. Milakis, D.; Van Arem, B.; Van Wee, B. Policy and Society Related Implications of Automated Driving: A
1126 Review of Literature and Directions for Future Research. *Journal of Intelligent Transportation Systems* **2017**,
1127 *21*, 324–348.
- 1128 8. Shi, S.; Jiang, L.; Dai, D.; Schiele, B. Motion Transformer with Global Intention Localization and Local
1129 Movement Refinement. *Advances in Neural Information Processing Systems (NeurIPS)*, 2022, Vol. 35, pp.
1130 6531–6543.
- 1131 9. Shi, S.; Jiang, L.; Dai, D.; Schiele, B. MTR++: Multi-Agent Motion Prediction with Symmetric Scene
1132 Modeling and Pair-Wise Interaction. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **2024**,
1133 *46*, 3039–3051.
- 1134 10. Nayakanti, N.; Al-Rfou, R.; Zhou, A.; Goel, K.; Refaat, K.S.; Sapp, B. Wayformer: Motion Forecasting
1135 via Simple & Efficient Attention Networks. *IEEE International Conference on Robotics and Automation*
1136 *(ICRA)*. IEEE, 2023, pp. 2980–2987.
- 1137 11. Huang, Z.; Liu, H.; Lv, C. GameFormer: Game-Theoretic Modeling and Learning of Transformer-Based
1138 Interactive Prediction and Planning for Autonomous Driving. *Proceedings of the IEEE/CVF International*
1139 *Conference on Computer Vision (ICCV)*, 2023, pp. 3903–3913.
- 1140 12. Ngiam, J.; Vasudevan, V.; Caine, B.; Zhang, Z.; Chiang, H.T.L.; Ling, J.; Roelofs, R.; Bewley, A.; Liu, C.;
1141 Vaswani, A.; others. Scene Transformer: A Unified Architecture for Predicting Future Trajectories of
1142 Multiple Agents. *International Conference on Learning Representations (ICLR)*, 2022.

- 1143 13. European Parliament and Council. Regulation (EU) 2024/1689 of the European Parliament and of the
1144 Council Laying Down Harmonised Rules on Artificial Intelligence (AI Act). Official Journal of the
1145 European Union, 2024.
- 1146 14. Koopman, P.; Wagner, M. Autonomous Vehicle Safety: An Interdisciplinary Challenge. *IEEE Intelligent*
1147 *Transportation Systems Magazine* **2017**, *9*, 90–96.
- 1148 15. Zablocki, E.; Ben-Younes, H.; Pérez, P.; Cord, M. Explainability of Deep Vision-Based Autonomous Driving
1149 Systems: Review and Challenges. *International Journal of Computer Vision* **2022**, *130*, 2425–2452.
- 1150 16. Ribeiro, M.T.; Singh, S.; Guestrin, C. “Why Should I Trust You?”: Explaining the Predictions of Any
1151 Classifier. Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and
1152 Data Mining, 2016, pp. 1135–1144.
- 1153 17. Lundberg, S.M.; Lee, S.I. A Unified Approach to Interpreting Model Predictions. Advances in Neural
1154 Information Processing Systems (NeurIPS), 2017, Vol. 30, pp. 4768–4777.
- 1155 18. Selvaraju, R.R.; Cogswell, M.; Das, A.; Vedantam, R.; Parikh, D.; Batra, D. Grad-CAM: Visual Explanations
1156 from Deep Networks via Gradient-Based Localization. Proceedings of the IEEE International Conference
1157 on Computer Vision (ICCV), 2017, pp. 618–626.
- 1158 19. Vig, J. A Multiscale Visualization of Attention in the Transformer Model. *arXiv preprint arXiv:1906.05714*
1159 **2019**.
- 1160 20. Abnar, S.; Zuidema, W. Quantifying Attention Flow in Transformers. Proceedings of the 58th Annual
1161 Meeting of the Association for Computational Linguistics (ACL), 2020, pp. 4190–4197.
- 1162 21. Chefer, H.; Gur, S.; Wolf, L. Transformer Interpretability Beyond Attention Visualization. Proceedings of
1163 the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), 2021, pp. 782–791.
- 1164 22. Jain, S.; Wallace, B.C. Attention is not Explanation. Proceedings of the 2019 Conference of the North
1165 American Chapter of the Association for Computational Linguistics (NAACL-HLT), 2019, pp. 3543–3556.
- 1166 23. Wiegrefe, S.; Pinter, Y. Attention is not not Explanation. Proceedings of the 2019 Conference on Empirical
1167 Methods in Natural Language Processing (EMNLP), 2019, pp. 11–20.
- 1168 24. Da Silva Martins, S.; Aldea, E.; Le Hégarat-Masclé, S. VISTA: A Vision and Intent-Aware Social Attention
1169 Framework for Multi-Agent Trajectory Prediction. *arXiv preprint arXiv:2511.10203* **2025**.
- 1170 25. Yadav, H.; Schaefer, M.; Zhao, K.; Meisen, T. LMFormer: Lane based Motion Prediction Transformer. *arXiv*
1171 *preprint arXiv:2504.10275* **2025**.
- 1172 26. Zhou, X.; Alecsandru, C. Local Lane Graph Conditioning as a General Inductive Bias for Trajectory
1173 Prediction: A Multi-Architecture Study on the Waymo Open Motion Dataset. *Sustainability* **2026**.
- 1174 27. Ettinger, S.; Cheng, S.; Caine, B.; Liu, C.; Zhao, H.; Pradhan, S.; Chai, Y.; Sapp, B.; Qi, C.R.; Zhou, Y.; others.
1175 Large Scale Interactive Motion Forecasting for Autonomous Driving: The Waymo Open Motion Dataset.
1176 Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV), 2021, pp. 9710–9719.
- 1177 28. Zhou, Z.; Wang, J.; Li, Y.H.; Huang, Y.K. Query-Centric Trajectory Prediction. Proceedings of the IEEE/CVF
1178 Conference on Computer Vision and Pattern Recognition (CVPR), 2023, pp. 17863–17873.
- 1179 29. Atakishiyev, S.; Salameh, M.; Yao, H.; Goebel, R. Explainable Artificial Intelligence for Autonomous
1180 Driving: A Comprehensive Overview and Field Guide for Future Research Directions. *IEEE Access* **2024**,
1181 *12*, 1–30.
- 1182 30. Taiebat, M.; Brown, A.L.; Safford, H.R.; Qu, S.; Xu, M. A Review on Energy, Environmental, and
1183 Sustainability Implications of Connected and Automated Vehicles. *Environmental Science & Technology*
1184 **2018**, *52*, 11449–11465.
- 1185 31. Vaswani, A.; Shazeer, N.; Parmar, N.; Uszkoreit, J.; Jones, L.; Gomez, A.N.; Kaiser, Ł.; Polosukhin, I.
1186 Attention is All You Need. Advances in Neural Information Processing Systems (NeurIPS), 2017, Vol. 30,
1187 pp. 5998–6008.
- 1188 32. Gao, J.; Sun, C.; Zhao, H.; Shen, Y.; Anguelov, D.; Li, C.; Schmid, C. VectorNet: Encoding HD Maps and
1189 Agent Dynamics from Vectorized Representation. Proceedings of the IEEE/CVF Conference on Computer
1190 Vision and Pattern Recognition (CVPR), 2020, pp. 11525–11533.
- 1191 33. Zeng, Z.; Mao, J.; Dai, B.; Anguelov, D. Heterogeneous Polyline Transformer with Relative Pose Encoding
1192 for Map-Aware Motion Prediction. Advances in Neural Information Processing Systems (NeurIPS), 2023,
1193 Vol. 36.
- 1194 34. Wu, W.; Feng, X.; Gao, Z.; Kan, Y. SMART: Scalable Multi-Agent Real-Time Motion Generation via
1195 Next-Token Prediction. Advances in Neural Information Processing Systems (NeurIPS), 2024, Vol. 37.

- 1196 35. Gao, Z.; Liu, D.; Yu, T.; Hu, H.; Gao, F.; Zhao, R. ISE-GT: Interaction Strength-Enhanced Graph Transformer
1197 for Explainable Multi-Agent Trajectory Prediction. *Transportation Research Part C: Emerging Technologies*
1198 **2025**.
- 1199 36. Pearl, J. Causal Inference in Statistics: An Overview. *Statistics Surveys* **2009**, *3*, 96–146.
- 1200 37. Goyal, Y.; Wu, Z.; Ernst, J.; Batra, D.; Parikh, D.; Lee, S. Counterfactual Visual Explanations. Proceedings
1201 of the 36th International Conference on Machine Learning (ICML), 2019, pp. 2376–2384.
- 1202 38. Tan, S.; Wong, K.; Wang, S.; Manivasagam, S.; Ren, M.; Urtasun, R. SceneGen: Learning to Generate
1203 Realistic Traffic Scenes. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern
1204 Recognition (CVPR), 2021, pp. 892–901.
- 1205 39. Suo, S.; Regalado, S.; Casas, S.; Urtasun, R. TrafficSim: Learning to Simulate Realistic Multi-Agent
1206 Behaviors. Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR),
1207 2021, pp. 10400–10409.
- 1208 40. Zhong, Z.; Rempe, D.; Xu, D.; Chen, Y.; Veer, S.; Che, T.; Ray, B.; Pavone, M. Guided Conditional Diffusion
1209 for Controllable Traffic Simulation. IEEE International Conference on Robotics and Automation (ICRA),
1210 2023, pp. 3560–3566.
- 1211 41. Ding, W.; Xu, C.; Arief, M.; Lin, H.; Li, B.; Zhao, D. A Survey on Safety-Critical Driving Scenario
1212 Generation—A Methodological Perspective. *IEEE Transactions on Intelligent Transportation Systems* **2023**,
1213 *24*, 6971–6988.
- 1214 42. Arrieta, A.B.; Díaz-Rodríguez, N.; Del Ser, J.; Bennetot, A.; Tabik, S.; Barbado, A.; Garcia, S.; Gil-Lopez,
1215 S.; Molina, D.; Benjamins, R.; others. Explainable Artificial Intelligence (XAI): Concepts, Taxonomies,
1216 Opportunities and Challenges toward Responsible AI. *Information Fusion* **2020**, *58*, 82–115.
- 1217 43. Litman, T. Autonomous Vehicle Implementation Predictions: Implications for Transport Planning.
1218 Technical report, Victoria Transport Policy Institute, 2023.
- 1219 44. Qi, C.R.; Su, H.; Mo, K.; Guibas, L.J. PointNet: Deep Learning on Point Sets for 3D Classification and
1220 Segmentation. Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR),
1221 2017, pp. 652–660.
- 1222 45. Ba, J.L.; Kiros, J.R.; Hinton, G.E. Layer Normalization. *arXiv preprint arXiv:1607.06450* **2016**.
- 1223 46. Loshchilov, I.; Hutter, F. Decoupled Weight Decay Regularization. International Conference on Learning
1224 Representations (ICLR), 2019.
- 1225 47. Micikevicius, P.; Narang, S.; Alben, J.; Diamos, G.; Elsen, E.; Garcia, D.; Ginsburg, B.; Houston, M.;
1226 Kuchaiev, O.; Venkatesh, G.; Wu, H. Mixed Precision Training. International Conference on Learning
1227 Representations (ICLR), 2018.
- 1228 48. Xiao, G.; Tian, Y.; Chen, B.; Han, S.; Lewis, M. Efficient Streaming Language Models with Attention Sinks.
1229 *arXiv preprint arXiv:2309.17453* **2023**.
- 1230 49. Zhai, S.; Likhomanenko, T.; Littwin, E.; Busbridge, D.; Ramapuram, J.; Zhang, Y.; Gu, J.; Susskind, J.M.
1231 Stabilizing Transformer Training by Preventing Attention Entropy Collapse. Proceedings of the 40th
1232 International Conference on Machine Learning (ICML), 2023, pp. 40770–40803.

1233 © 2026 by the authors. Submitted to *Sustainability* for possible open access publication
1234 under the terms and conditions of the Creative Commons Attribution (CC BY) license
1235 (<http://creativecommons.org/licenses/by/4.0/>).