

Scenario-based Evaluation of Prediction Models for Automated Vehicles

Manuel Muñoz Sánchez¹ Jos Elfring² Emilia Silvas³ and René van de Molengraft¹

Abstract—To operate safely, an automated vehicle (AV) must anticipate how the environment around it will evolve. For that purpose, it is important to know which prediction models are most appropriate for every situation. Currently, assessment of prediction models is often performed over a set of trajectories without distinction of the type of movement they capture, resulting in the inability to determine the suitability of each model for different situations. In this work we illustrate how standardized evaluation methods result in wrong conclusions regarding a model's predictive capabilities, preventing a clear assessment of prediction models and potentially leading to dangerous on-road situations. We argue that following evaluation practices in safety assessment for AVs, assessment of prediction models should be performed in a scenario-based fashion. To encourage scenario-based assessment of prediction models and illustrate the dangers of improper assessment, we categorize trajectories of the Waymo Open Motion dataset according to the type of movement they capture. Next, three different models are thoroughly evaluated for different trajectory types and prediction horizons. Results show that common evaluation methods are insufficient and the assessment should be performed depending on the application in which the model will operate.

I. INTRODUCTION

Automated vehicles (AVs) have become popular in recent years since they have the potential to increase road safety, efficiency and comfort [1]–[3]. To operate safely, an AV must accurately anticipate the future motion of other road users (RUs) in its surroundings [4]. To build trajectory prediction models, deep learning (DL) techniques are gaining attention [5], since they can effectively learn complex interactions between different RUs [6], [7] and the road infrastructure [8], [9] from past observations to produce more accurate predictions. Traditionally, training these models effectively was a problematic task since the amount of data required was not easily available. However, this issue has been alleviated in recent years with the release of several large public datasets [10]–[14]. A common practice to assess a model's predictive accuracy is to consider a fraction of the dataset reserved for this purpose (commonly referred to as test data), and to compare the model's predictions with the real trajectories. The output of prediction models may vary, hence different metrics exist to quantify the disparity between the real and predicted trajectories [4]. For example, some models

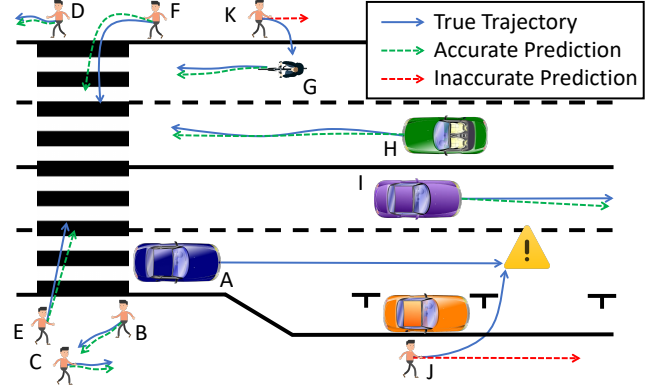


Fig. 1. Example where a model that is accurate on average fails to predict a pedestrian trajectory, leading to a dangerous situation.

produce a single prediction, while others produce a set of feasible trajectories and associated confidence for each.

Despite the existence of various evaluation metrics for prediction models, several challenges remain unaddressed in current evaluation practices, such as the inability of these metrics to capture a model's robustness or generalization capabilities [5]. Perhaps the most severe shortcoming is that all trajectories are considered equal for error computation despite capturing significantly different behaviors, which can lead to dangerous situations due to misjudgement of a model's suitability for specific situations. For instance, consider the situation shown in Fig. 1, where an AV (A) predicts the future trajectory of surrounding RUs (B-K) in a crowded urban scenario. Current evaluation practices would deem this model suitable for RU trajectory prediction in crowded urban scenarios, since its predictions are highly accurate on average. It accurately predicts pedestrians on the sidewalk (B-D), crossing at designated crossings (E,F), and lane-following cyclists and vehicles (G-I). However, in this example only a few of these RUs are relevant to the AV (I, J). Additionally, failure cases like the pedestrians crossing at non-designated crossings (J, K) can go unnoticed since all trajectories are considered equally for error computation.

The importance of a thorough evaluation for different types of trajectories has been recognized previously [15]. However, current efforts to improve evaluation of prediction models focus mainly on interactions between pedestrians (e.g. collision-avoidance [15]), and disregard interactions of RUs with the road infrastructure (e.g. pedestrian stops at a red traffic light). Additionally, the evaluation procedure should provide a transparent assessment of a model's suitability for the intended application. For instance, for AVs, an inaccurate prediction for a pedestrian walking in front of the vehicle should be considered more important or severe than one of a pedestrian that is walking behind the vehicle or far

This work was supported by SAFE-UP under EU's Horizon 2020 research and innovation programme, grant agreement 861570.

¹Manuel Muñoz Sánchez, Emilia Silvas, Jos Elfring and René van de Molengraft are with the Department of Mechanical Engineering, Eindhoven University of Technology, Eindhoven, The Netherlands.

²Jos Elfring is also with the Product Unit Autonomous Driving, TomTom, Amsterdam, The Netherlands.

³Emilia Silvas is also with the Department of Integrated Vehicle Safety, TNO, Helmond, The Netherlands.

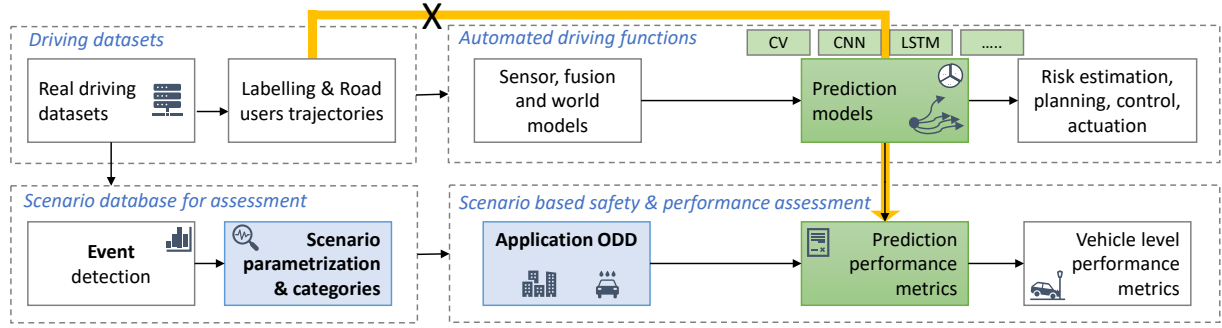


Fig. 2. Overview of a scenario-based assessment pipeline. The orange arrow indicates the standard approach to evaluate prediction models. Green blocks concern prediction-specific activities. Blue blocks are currently lacking but required steps for thorough evaluation of prediction models. White blocks are relevant for vehicle level assessment, but out of the scope of this work.

from the road. As a second example, consider the purpose of vehicle fuel and energy optimization, where accurate long-term predictions are required for optimal path planning. On the contrary, for the development of emergency advanced driver-assistance system features (e.g. emergency braking or emergency steering) accurate short-term predictions become more relevant. Thus, it is important to assess the suitability of models with respect to the functional applications they will be used in, in other words, their operational design domain (ODD), and therein test for various scenarios and their overall impact with vehicle level performance metrics, not only prediction metrics (Fig. 2).

Current evaluation practices reporting averaged errors over all predicted trajectories are beneficial for ease of comparison between different models. If a model achieves a lower error over all the trajectories in the dataset, one can confidently say such a model is more accurate, at least on average. However, it remains unclear under which circumstances this model is preferred over others. To show that an improved assessment of prediction models is needed, and working towards that goal, the contribution of our work is as follows:

- 1) We illustrate the extent to which common evaluation methods, which only report average errors over all trajectories, result in misleading conclusions of a model's predictive capabilities, and argue that a scenario-based assessment is a more suitable approach.
- 2) We facilitate scenario-based evaluation of prediction models providing an open source framework¹, which will allow for a transparent evaluation of a model's capabilities for different situations, leading to an optimal choice of prediction model depending on the application.

The remainder of this article is structured as follows. Section II introduces common trajectory prediction metrics and datasets, and presents related work on scenario-based evaluation. Section III introduces the prediction models compared and outlines how the comparison will be done using standard evaluation practices. Section IV presents an analysis of the results, and Section V concludes the work and highlights future improvements.

II. PRELIMINARIES

This section summarizes the most commonly used performance metrics, recent datasets used to develop AV applications, and related work on scenario-based evaluation.

A. Common Trajectory Prediction Metrics

A plethora of performance indicators exist to evaluate trajectory prediction models [4], with average displacement error (ADE) and final displacement error (FDE) being the most popular [5]. ADE measures the difference between the predicted and ground truth trajectories, averaged over all prediction horizons. FDE measures this difference at a specific horizon. To allow comparison of deterministic models that produce a single trajectory with probabilistic models that produce multiple feasible trajectories, variants of these metrics are used which report the errors of the trajectory that achieved the best accuracy. These variants are commonly referred to as minADE and minFDE. Although these metrics have several limitations [4] and new metrics have been introduced recently to address some of these limitations [10], we use them in this work since they remain the most common performance indicators at the moment.

To formally define these metrics, let $\hat{\mathbf{S}}$ denote a set of trajectory predictions for a set of road users \mathbf{N} at future prediction horizons \mathbf{T} . The minADE of the predictions for prediction horizon t is given by

$$\text{minADE}(\hat{\mathbf{S}}, t) = \sum_{n \in \mathbf{N}} \min_{\hat{s} \in \hat{\mathbf{S}}^n} \sum_{\substack{t' \in \mathbf{T} \\ t' \leq t}} \frac{\|\hat{s}_t - s_t^n\|_2}{|\mathbf{N}| \times |\mathbf{T}|}, \quad (1)$$

where $\hat{\mathbf{S}}^n$ denotes the set of predictions for a road user n , \hat{s}_t denotes the predicted position at time t , and s_t^n denotes the true position of road user n at time t . Additionally, $|\cdot|$ denotes the size of a set and $\|\cdot\|_2$ denotes the L2-norm of a vector. Similarly, the minFDE at a given prediction horizon t is defined as

$$\text{minFDE}(\hat{\mathbf{S}}, t) = \sum_{n \in \mathbf{N}} \min_{\hat{s} \in \hat{\mathbf{S}}^n} \frac{\|\hat{s}_t - s_t^n\|_2}{|\mathbf{N}|}. \quad (2)$$

B. Recent Datasets & Waymo's Motion Prediction Challenge

Several large datasets are publicly available for development and evaluation of prediction models, with some of the most recent and often used ones as summarized in Table I.

¹Code available at <https://github.com/manolotis/SBEP>

For this research we chose the Waymo Open Motion Dataset (WOMD) since it has the longest horizon, covers several cities, and contains extra information such as traffic light states.

With the release of WOMD, the Waymo motion prediction challenge² (WMPC) was introduced. In this challenge, the task is to predict the trajectories of a subset of RUs for 8 seconds into the future, given their history for the past 1 second and corresponding map of the area. The trajectories required to be predicted are selected to include interesting behavior and a balance of RUs as specified in [10]. In this work, we refer to those trajectories as *trajectories to predict* (TTP), which make up about 7% of all trajectories present in the dataset.

TABLE I
OVERVIEW OF DATASETS (ADAPTED FROM [10])

	Lyft	NuSc	Argo.	Inter.	WOMD
Reference article	[13]	[12]	[11]	[14]	[10]
Prediction horizon [s]	5	6	3	3	8
Number of segments	170k	1k	324k	-	104k
Segment duration [s]	25	20	5	-	20
Sampling rate [Hz]	10	2	10	10	10
Cities	1	2	2	6	6
Map available	✓	✓	✓	✓	✓
Traffic light states	✓				✓

C. Scenario-based Assessment

The need for a more complete assessment of a model's predictive capabilities has been briefly recognized in previous works. Some authors recognize that the largest prediction errors occur in non-linear regions of the trajectory and report the ADE for these regions separately [16], [17]. Some other works report on average maximum errors along the entire predicted trajectory to give an indication of worst-case predictions [18], [19]. However, these practices have not become the standard and lack the ability to capture relevant situations from the point of view of an AV.

In the area of safety assessment for driver assistance systems and AVs, a scenario-based approach has been adopted. This approach presents several benefits, such as the ability to evaluate the coverage of the assessment, and the possibility of a direct translation between test outcomes and an assessment of the AV's performance with respect to a specific ODD, ultimately facilitating legal and public acceptance of AVs [20].

In the area of trajectory prediction assessment, this scenario-based approach has not been generally adopted yet. TrajNet++ is a recent benchmark with the goal of standardizing trajectory prediction applying similar concepts [15]. However, the focus is on interacting RUs (mainly pedestrians) and disregard other interactions with the road infrastructure (e.g. a traffic light). Additionally, specification of scenarios should consider relevant situations for the AV. For example, failing to predict the trajectory of a pedestrian

that ends up on the road to avoid another pedestrian is irrelevant to the AV if this interaction occurs behind or far from the vehicle. However, if it occurs immediately in front of the vehicle, it would be highly relevant. The authors of [21] proposed a framework for scenario-based testing of prediction models in a simulated environment. The framework supports modeling and generation of scenarios involving interactive RUs for a thorough evaluation of prediction models. This approach overcomes several limitations of current evaluation methods. However, it is important to also evaluate how prediction models perform with real driving data, and the impact on the entire AV architecture with vehicle-level performance metrics.

III. METHODOLOGY & EXPERIMENTS

To perform a thorough analysis of the performance of a prediction model for different types of trajectories, we first need to systematically detect these trajectory types in existing datasets. As a first step towards a scenario-based evaluation framework for prediction models, we consider individual RU trajectories present in WOMD and categorize them assigning one or more of the tags summarized in Table II. These tags are chosen to explore differences in performance between trajectories of different shapes (T1 and T2), different behaviors (T3-T5), different availability of observations (T6-T9), and the same or different trajectories as specified in the WMPC (T10 and T11). Future iterations of this work will include more complex scenarios (e.g. pedestrian at non-designated crossing ahead of the AV). Examples of trajectories from some of the selected tags are shown in Fig. 3.

TABLE II
TAGS USED TO LABEL TRAJECTORIES

Tag	Description
T1	<i>Straight</i> - RU closely follows a straight path
T2	<i>Non-straight</i> - RU deviates from a straight path
T3	<i>Starting</i> - RU is still during observation and moves in the future
T4	<i>Stopping</i> - RU moves during observation and stops in the future
T5	<i>Still</i> - RU is still during observation and in the future
T6	<i>Late</i> - RU is detected late (≤ 0.3 sec before prediction)
T7	<i>Very Late</i> - RU is detected very late (0.1 sec before prediction)
T8	<i>Full</i> - RU is detected during the entire observation period
T9	<i>Reappearance</i> - The same RU disappears during observation and reappears in the future
T10	<i>TTP</i> - Trajectories To Predict - Required trajectories to predict for Waymo's Motion Prediction Challenge
T11	<i>NTTP</i> - Trajectories that were not required for Waymo's Motion Prediction Challenge

To illustrate the importance of a thorough assessment, the performance of three different models is compared, first in a manner that adheres to common current evaluation practices, and then considering additional aspects and revealing observations that are crucial for understanding the limitations of each of the methods but that could not be concluded from the initial evaluation.

A. Evaluated Models

Three different prediction models are compared, in increasing level of complexity. The inputs and outputs of each model are summarized in Table III.

²<https://waymo.com/open/challenges/2021/motion-prediction/>

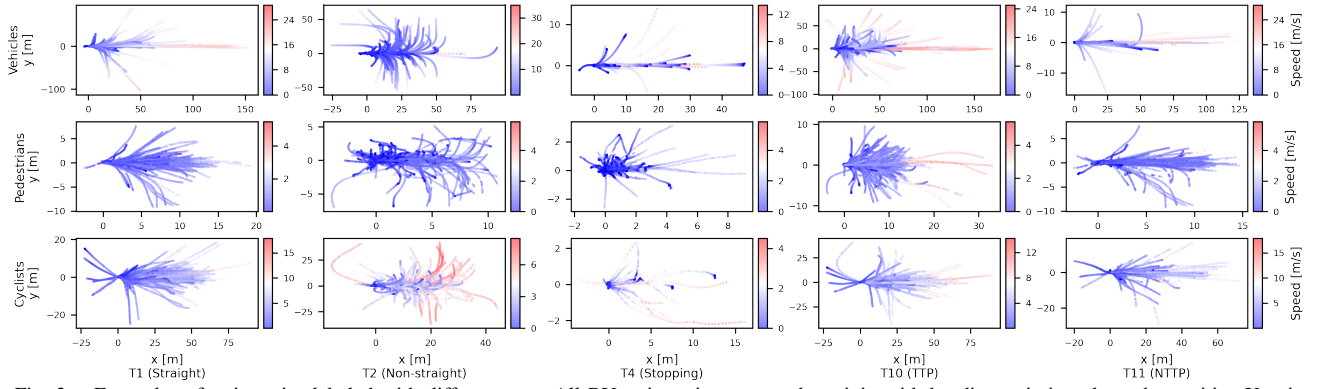


Fig. 3. Examples of trajectories labeled with different tags. All RU trajectories start at the origin with heading pointing along the positive X axis.

a) A Constant Velocity (CV) model, which is often used as a baseline [22]. CV assumes every RU maintains the same velocity over all prediction horizons, considering only the last observed RU state and disregarding the static (e.g. walls) and dynamic (e.g. other RUs) environment.

b) An LSTM encoder-decoder neural network. Virtually any state of the art DL method has some recurrent neural network component to encode the past of RUs, with LSTM-based architectures being a popular choice [23], [24]. The implemented LSTM considers past observations of RUs and does not exploit any knowledge about road topology or surrounding RUs, meaning that its predictions will not avoid overlap with other RUs or static obstacles similarly to CV.

c) MotionCNN [8], the 3rd place solution of the WMPC. MotionCNN provides an elegant solution to exploit different types of information (i.e. past of RUs, surrounding RUs, and road infrastructure) with an architecture based on convolutional neural networks only, and it produces 6 feasible trajectories. We downloaded a pre-trained model from an open repository provided by the authors³.

models' performance in increasing detail and draw conclusions regarding their capabilities. Subsequent, more detailed iterations reveal different, possibly conflicting conclusions, showing that current evaluation methods can be misleading.

A. Overall Evaluation

Table IV shows the performance of the models for trajectories of vehicles (veh.), pedestrians (ped.) and cyclists (cyc.) aggregated over all prediction horizons, for the required trajectories in the WMPC. From these results we might conclude:

- C1) DL models are superior when compared to CV, especially for predicting vehicles' trajectories, as can be seen by the large difference in minADE and minFDE.
- C2) MotionCNN is clearly the most suitable model to predict trajectories of vehicles and cyclists.
- C3) LSTM is better suited for pedestrian trajectory prediction than the other two models.

Many works report the performance of their method in a similar fashion as presented in Table IV [8], [9], over all prediction horizons and only for a subset of trajectories deemed interesting, raising the question of whether or not performance remains the same for other trajectories. Table V summarizes the models' accuracies after evaluation on all trajectories not considered previously. From this table, additional conclusions can be made:

TABLE III
OVERVIEW OF COMPARED MODELS AND THEIR INPUTS AND OUTPUTS

Model	Input			Output	
	Past RU state	Other RUs	Map	Single Trajectory	Multiple Trajectories
CV				✓	
LSTM	✓			✓	
MotionCNN	✓	✓	✓		✓

B. Evaluation Criteria

To assess the accuracy of the predictions at various prediction horizons, the horizons on which the models are evaluated start with the first future timestep (corresponding to 0.1 seconds into the future), and follow the same sampling rate of 2Hz as in the WMPC. Thus, we compute minADE and minFDE as defined in (1) and (2) for $T = \{0.1 + \frac{t}{2} \mid 0 \leq t < 16 \wedge t \in \mathbb{Z}\}$ seconds.

IV. RESULTS

An iterative approach is taken to report the performance of the evaluated models. At each iteration, we analyze the

TABLE IV
MODEL PERFORMANCE OVER ALL PREDICTION HORIZONS (TTP)

Model	Metric					
	minADE			minFDE		
	Veh.	Ped.	Cyc.	Veh.	Ped.	Cyc.
CV	3.893	0.633	1.632	10.468	1.477	3.964
LSTM	2.025	0.545	1.481	5.681	1.295	3.618
MotionCNN	1.482	0.564	1.247	3.855	1.315	3.012

Red and blue indicate highest and lowest error per RU

TABLE V
MODEL PERFORMANCE OVER ALL PREDICTION HORIZONS (NTTP)

Model	Metric					
	minADE			minFDE		
	Veh.	Ped.	Cyc.	Veh.	Ped.	Cyc.
CV	0.640	0.565	1.070	1.834	1.284	2.816
LSTM	0.391	0.320	1.014	1.114	0.728	2.441
MotionCNN	0.806	0.637	1.178	2.297	1.575	3.291

³<https://github.com/kbrodt/waymo-motion-prediction-2021/releases>

- C4)** A simple model such as CV outperforms complex DL methods such as MotionCNN, which contradicts C1 and C2. It seems that trajectories deemed uninteresting (NTTP) originate from RUs moving at constant velocity.
- C5)** MotionCNN achieves the lowest accuracy for all three types of RUs, contradicting C1 and C2. This suggests MotionCNN suffers from overfitting, since it does not generalize well to different trajectories, which is problematic since in real-world applications one does not know beforehand if a trajectory would be considered TTP or not.
- C6)** LSTM presents the highest accuracy for prediction of pedestrian trajectories (C3 still holds), and additionally vehicles and possibly cyclists (contradicting C2).
- C7)** The errors are significantly lower, so these trajectories are less challenging to predict.

Following C7, the reader might wonder how it is possible that MotionCNN performs the best for selected challenging trajectories, yet it performs the worst for easier trajectories. The reason might not be obvious, since we purposely left out some practical details on how the models were trained to emphasize the importance of reporting practical development details of data-driven models. LSTM was trained using all the available trajectories, and MotionCNN using only those required to be predicted in the WMPC. As such, it is not extraordinary to observe the differences in Tables IV and V, as evaluation on only TTP or NTTP trajectories would give an advantage to the model trained with those trajectories. However, this reveals yet another pitfall of standard metrics: they do not provide any indication of a model's performance for trajectories that are uncommon in the data. This generalization issue, which is quite common in data-driven models, has already been recognized in the context of motion prediction [5], but it remains unaddressed. If we were to further train MotionCNN also using the unseen trajectories, its performance would improve for these, but it remains unclear if its performance on the selected trajectories would be negatively affected.

B. Evaluation Per Time Horizon

A model's performance might vary significantly depending on the prediction horizon considered, so it is important take that into account, since different applications have different accuracy requirements for different horizons. Evaluating accuracy at different horizons is also common in trajectory prediction literature [25], [26]. Additionally, further analyzing other derived metrics such as standard deviation or maximum prediction errors, which is not commonly done, might reveal interesting insights. Table VI summarizes such an analysis (only for trajectories labeled TTP), revealing new insights that were not obvious from the previous superficial analysis from Table IV:

- C8)** CV can match and even outperform complex DL methods, *for very short prediction horizons*, which contradicts C1-C3 and C6, and further details C4.
- C9)** CV can no longer match the performance of the other models for horizons longer than 1 second, which adds detail to C4.
- C10)** LSTM presents the highest errors *for very short prediction horizons*, even for pedestrians, which contradicts C3 and C6.
- C11)** LSTM predictions of pedestrian trajectories are the closest to the real trajectory overall, but not necessarily for long prediction horizons (as seen by the lowest minADE, but not minFDE at time 7.6 seconds).
- C12)** MotionCNN is best suited for trajectory prediction of vehicles and cyclists overall, but its worst predictions can be less accurate than those of LSTM, even when considering the best out of the 6 predicted trajectories.

Several additional observations could be made from Table VI, which would either contradict or confirm the initial conclusions with a finer level of detail. Thus, it is important to always keep the intended application in mind and analyze the aspects that are most relevant for this purpose.

TABLE VI
MODEL PERFORMANCE AT DIFFERENT PREDICTION HORIZONS (ONLY TRAJECTORIES LABELED TTP)

Model	RU	Time [s] Metric [m]	0.1			1.1			3.1			5.1			7.6		
			std	max		std	max		std	max		std	max		std	max	
CV	Veh.	minADE	0.034	0.084	6.90	0.393	0.717	72.889	2.187	3.021	205.598	5.26	6.805	336.425	9.878	12.312	465.797
		minFDE	0.034	0.084	6.90	0.842	0.971	72.889	5.663	4.303	205.598	14.142	9.796	336.425	27.323	17.929	465.797
	Ped.	minADE	0.023	0.025	0.825	0.124	0.163	6.123	0.434	0.573	17.715	0.847	1.114	29.263	1.415	1.858	43.429
		minFDE	0.023	0.025	0.825	0.235	0.211	6.123	0.969	0.85	17.715	1.974	1.703	29.263	3.473	2.929	43.429
	Cyc.	minADE	0.054	0.048	0.449	0.285	0.333	3.567	1.061	1.59	53.026	2.177	3.017	58.425	3.798	5.186	67.002
		minFDE	0.054	0.048	0.449	0.539	0.41	3.567	2.431	2.278	51.616	5.271	4.623	58.425	9.659	8.238	67.002
LSTM	Veh.	minADE	0.051	0.109	9.492	0.223	0.425	66.493	1.076	1.617	68.542	2.655	3.952	68.542	5.313	7.957	119.801
		minFDE	0.051	0.109	9.492	0.435	0.564	65.049	2.762	2.481	61.582	7.368	6.487	63.179	15.828	13.615	119.801
	Ped.	minADE	0.059	0.056	2.79	0.114	0.127	2.852	0.36	0.514	8.716	0.72	1.04	14.252	1.236	1.761	25.532
		minFDE	0.059	0.056	2.79	0.183	0.173	2.541	0.808	0.812	8.716	1.719	1.656	14.252	3.124	2.859	25.532
	Cyc.	minADE	0.216	0.184	2.558	0.30	0.272	2.925	0.923	1.496	54.011	1.936	2.887	62.589	3.462	5.004	72.982
		minFDE	0.216	0.184	2.558	0.428	0.347	2.925	2.109	2.23	54.011	4.793	4.555	62.589	9.029	8.047	72.982
MotionCNN	Veh.	minADE	0.038	0.058	4.487	0.219	0.432	66.633	0.888	1.404	68.55	1.973	3.138	94.614	3.64	6.001	150.044
		minFDE	0.038	0.058	4.487	0.419	0.584	65.263	2.141	2.198	63.885	5.089	5.358	94.614	10.08	11.026	150.044
	Ped.	minADE	0.022	0.019	0.477	0.108	0.134	1.723	0.384	0.523	7.25	0.758	1.049	16.647	1.261	1.745	35.719
		minFDE	0.022	0.019	0.477	0.202	0.173	1.723	0.868	0.801	7.25	1.777	1.654	16.647	3.05	2.822	35.719
	Cyc.	minADE	0.047	0.039	0.321	0.215	0.257	2.803	0.814	1.429	53.809	1.67	2.676	62.243	2.887	4.578	72.027
		minFDE	0.047	0.039	0.321	0.404	0.326	2.803	1.874	2.12	53.809	4.062	4.365	62.243	7.213	7.775	72.027

C. Evaluation Per Type of Trajectory

Next we analyze model performance for different types of trajectories according to the tags described earlier. Table VII shows the results and new observations can be made:

C13) Some behavior is missing in the trajectories used for evaluation. There are not pedestrian or cyclist trajectories where the RU stands still during the observation period and either remains still or moves in the next 8 seconds⁴ (denoted by “-” in Table VII). Additionally, there are no trajectories where the RU disappears shortly before making the prediction. Thus, if an AV is presented with these situations, it will not be possible to determine which model is most suitable.

C14) CV is most suitable for predictions where the RU remains still (trivial). MotionCNN performs the worst by far in these cases. If an AV is presented with this situation, it might use a model that is not the most appropriate in this case despite being the most accurate overall.

C15) CV is most suitable for predicting starting behavior. This conclusion could be misleading, as CV naively predicts the RU remains still, and it does not predict starting behavior. If the RU remains still for most of the future 8 seconds, the overall error will be low, but it does not capture CV’s ability to predict when movement will start.

C16) LSTM is particularly well suited for prediction of all types of pedestrian trajectories, which supports some of our previous conclusions (i.e. C3, C6 and partially C11), but contradicts some others (i.e. C8, C10 and C14).

C17) MotionCNN is well suited for prediction of trajectories labeled TTP, straight trajectories of vehicles and cyclists, and non-straight trajectories of vehicles. However, for other types of trajectories, such as late detections or trajectories labeled NTTP, it can even be outperformed by CV.

⁴An RU was considered still if its speed did not exceed 0.01 m/s.

After analyzing model performance considering different aspects like various prediction horizons and trajectory types, it is still not possible to conclude on the potential suitability of a data-driven model, since their performance is heavily affected by the data used to train them. Recall that MotionCNN was trained to predict only those trajectories labeled TTP, while LSTM was trained using all trajectories. Figure 4 shows the percentage of trajectories labeled with each tag, both for the entire dataset and only considering trajectories labeled TTP, which partly explains the difference in performance between LSTM and MotionCNN:

C18) LSTM outperforms MotionCNN for most trajectory types because it has been trained with more instances of each type. MotionCNN was trained on trajectories labeled TTP, which are only about 7% of available data.

C19) Within the sets of trajectories used to train each model, the relative frequency of some behaviors is very different. For instance, MotionCNN is unable to predict starting and still trajectories accurately because this behavior is significantly underrepresented in TTP trajectories (approximately 0.08 and 0.1%), as opposed to their frequency in the entire dataset used to train LSTM (18 and 53%).

Even if a model’s performance can be better explained by an in-depth analysis of the training data, it does not mean it is always necessary to do so. If a model will only be used in specific situations (e.g. late detections), then it must be accurate in these situations no matter how inaccurate it might be in others. Similarly, if a model’s purpose is to increase safety or fuel efficiency, then vehicle-level performance metrics after integrating this prediction model in the vehicle should be the main assessment criteria, since a marginal improvement in predictive accuracy might yield little to no improvement for the intended application. Thus, evaluation of these models should be done according to the application in which they operate.

TABLE VII
MODEL PERFORMANCE FOR DIFFERENT TYPES OF TRAJECTORIES

Model	RU	Metric [m]	T1 Straight	T2 NonStraight	T3 Starting	T4 Stopping	T5 Still	T6 Late	T7 VeryLate	T8 Full	T9 Reappear	T10 TTP	T11 NTTP
CV	Veh.	minADE	2.872	2.747	0.821	1.442	0	1.387	1.316	1.089	-	3.893	0.64
		minFDE	8.422	7.285	2.847	3.73	0	4.395	4.175	3.15	-	10.468	1.834
	Ped.	minADE	0.466	0.974	-	0.43	-	1.092	1.102	0.486	-	0.633	0.565
		minFDE	1.127	2.226	-	0.908	-	2.721	2.752	1.152	-	1.477	1.284
	Cyc.	minADE	1.445	2.004	-	1.071	-	1.91	2.037	1.381	-	1.632	1.07
		minFDE	3.67	5.208	-	2.357	-	4.976	5.38	3.54	-	3.964	2.816
LSTM	Veh.	minADE	1.549	1.693	0.828	0.581	0.03	1.364	1.449	0.574	-	2.025	0.391
		minFDE	4.679	4.748	2.853	1.608	0.027	3.692	3.642	1.739	-	5.681	1.114
	Ped.	minADE	0.401	0.597	-	0.182	-	0.715	0.843	0.344	-	0.545	0.32
		minFDE	0.973	1.495	-	0.38	-	1.619	1.835	0.845	-	1.295	0.728
	Cyc.	minADE	1.327	1.93	-	0.967	-	2.308	3.143	1.234	-	1.481	1.014
		minFDE	3.304	5.027	-	2.142	-	4.39	5.129	3.203	-	3.618	2.441
MotionCNN	Veh.	minADE	1.342	1.297	1.058	0.861	0.629	1.868	2.745	0.836	-	1.482	0.806
		minFDE	3.633	3.331	3.109	2.40	1.861	4.483	6.232	2.429	-	3.855	2.297
	Ped.	minADE	0.478	0.981	-	0.764	-	1.748	2.372	0.506	-	0.564	0.637
		minFDE	1.126	2.425	-	1.833	-	4.285	5.708	1.277	-	1.315	1.575
	Cyc.	minADE	1.19	1.994	-	1.261	-	2.408	2.399	1.072	-	1.247	1.178
		minFDE	2.854	5.222	-	3.185	-	5.967	6.104	2.782	-	3.012	3.291

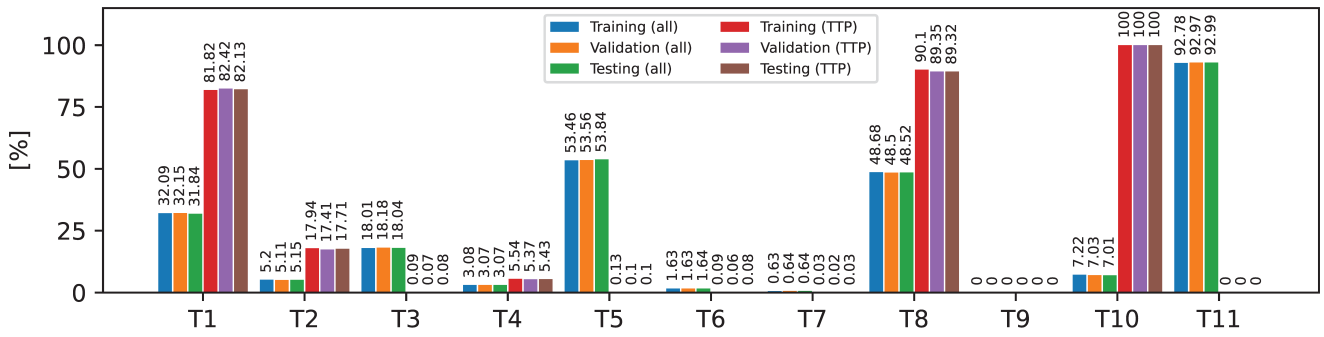


Fig. 4. Percentage of trajectories with a specific tag for the entire dataset and the set of trajectories labeled TTP.

V. CONCLUSION AND FUTURE WORK

To operate safely, an AV must anticipate the future motion of other RUs in its surroundings through trajectory prediction. Assessment of prediction models is commonly performed over a set of trajectories without distinction over the type of movement captured by each trajectory, which does not provide a clear overview of the suitability of each model for different situations. Furthermore, the impact of a marginal increase in predictive accuracy at the vehicle level remains unclear, as other components of an AV are normally not considered for assessment of prediction models.

In this work, we have illustrated the extent to which standard evaluation practices result in misleading conclusions of a model's predictive capabilities. Additionally, we have made publicly available a scenario-based framework for evaluation of prediction models, which allows classification of individual trajectories according to the type of movement they capture and facilitates a clear assessment of a model's suitability to predict each trajectory type.

Future work will extend the framework proposed in [20] to model and capture relevant scenarios for prediction algorithms, facilitating assessment for their intended application.

REFERENCES

- [1] M. M. Morando *et al.*, "Studying the Safety Impact of Autonomous Vehicles Using Simulation-Based Surrogate Safety Measures," *Journal of Advanced Transportation*, 2018.
- [2] J. Ploeg, A. F. A. Serrarens, and G. J. Heijenk, "Connect & Drive: design and evaluation of cooperative adaptive cruise control for congestion reduction," *Journal of Modern Transportation*, vol. 19, no. 3, pp. 207–213, 2011.
- [3] D. Milakis, B. Van Arem, and B. Van Wee, "Policy and society related implications of automated driving: A review of literature and directions for future research," *Journal of ITS: Technology, Planning, and Operations*, 2017.
- [4] A. Rasouli, "Deep Learning for Vision-based Prediction: A Survey," 2020.
- [5] A. Rudenko *et al.*, "Human motion trajectory prediction: a survey," *The International Journal of Robotics Research*, vol. 39, no. 8, pp. 895–935, 2020.
- [6] A. Alahi *et al.*, "Social LSTM: Human Trajectory Prediction in Crowded Spaces," in *CVPR*, 2016.
- [7] A. Gupta *et al.*, "Social GAN: Socially Acceptable Trajectories with Generative Adversarial Networks," in *CVPR*, 2018.
- [8] S. Konev, K. Brodt, and A. Sanakoyeu, "MotionCNN: A Strong Baseline for Motion Prediction in Autonomous Driving," 2021.
- [9] J. Gu, C. Sun, and H. Zhao, "Densetnt: End-to-end trajectory prediction from dense goal sets," in *IEEE/CVF ICCV*, 2021.
- [10] S. Ettinger *et al.*, "Large Scale Interactive Motion Forecasting for Autonomous Driving: The Waymo Open Motion Dataset," 2021.
- [11] M.-F. Chang *et al.*, "Argoverse: 3D Tracking and Forecasting With Rich Maps," in *CVPR*, 2019.
- [12] H. Caesar *et al.*, "nuScenes: A Multimodal Dataset for Autonomous Driving," in *CVPR*, 2020.
- [13] R. Kesten *et al.*, *Lyft level 5 perception dataset 2020*, 2019.
- [14] W. Zhan *et al.*, "INTERACTION Dataset: An INTERNATIONAL, Adversarial and Cooperative motion Dataset in Interactive Driving Scenarios with Semantic Maps," 2019.
- [15] P. Kothari, S. Kreiss, and A. Alahi, "Human Trajectory Forecasting in Crowds: A Deep Learning Perspective," *IEEE Transactions on ITS*, 2021.
- [16] A. Alahi *et al.*, "Social LSTM: Human Trajectory Prediction in Crowded Spaces," in *CVPR*, 2016.
- [17] Y. Xu, Z. Piao, and S. Gao, "Encoding Crowd Interaction with Deep Neural Network for Pedestrian Trajectory Prediction," in *CVPR*, 2018.
- [18] N. Lee *et al.*, "DESIRE: Distant Future Prediction in Dynamic Scenes with Interacting Agents," in *CVPR*, 2017.
- [19] P. Felsen, P. Lucey, and S. Ganguly, "Where Will They Go? Predicting Fine-Grained Adversarial Multi-agent Motion Using Conditional Variational Autoencoders," in *Lecture Notes in Computer Science*, vol. 11215 LNCS, 2018, pp. 761–776.
- [20] E. De Gelder *et al.*, "Towards an Ontology for Scenario Definition for the Assessment of Automated Vehicles: An Object-Oriented Framework," *IEEE Trans. on IVs*, 2022.
- [21] F. Indaheng *et al.*, "A Scenario-Based Platform for Testing Autonomous Vehicle Behavior Prediction Models in Simulation," 2021.
- [22] S. Pellegrini *et al.*, "You'll never walk alone: Modeling social behavior for multi-target tracking," in *ICCV*, 2009.
- [23] B. Varadarajan *et al.*, "MultiPath++: Efficient Information Fusion and Trajectory Aggregation for Behavior Prediction," 2021.
- [24] M. Muñoz Sánchez *et al.*, "A Hybrid Framework Combining Vehicle System Knowledge with Machine Learning Methods for Improved Highway Trajectory Prediction," in *SMC*, 2020.
- [25] N. Deo and M. M. Trivedi, "Convolutional Social Pooling for Vehicle Trajectory Prediction," in *CVPR Workshops*, 2018.
- [26] R. Chandra *et al.*, "RobustTP: End-to-end trajectory prediction for heterogeneous road-agents in dense traffic with noisy sensor inputs," in *ACM CSCS*, 2019.