
A2X: An Agent and Environment Interaction Benchmark for Multimodal Human Trajectory Prediction

Abstract

1 Recent trends in human trajectory prediction are the development of generative
2 models which generate distributions of trajectories. However existing metrics
3 are suited only for single (unimodal) trajectory instances. Furthermore, existing
4 datasets are largely limited to small-scale interactions between people, with little
5 to no agent-to-agent environment interaction. To address these challenges, we
6 propose a dataset that compensates for the lack of agent-to-environment interaction
7 in existing datasets with a new simulated dataset and metrics to convey model
8 performance with more reliability and nuance. A subset of these metrics are
9 novel *multiverse metrics*, which are better-suited for multimodal models than
10 existing metrics but are still applicable to unimodal models. Our results showcase
11 the benefits of the augmented dataset and metrics. The dataset is available at:
12 <https://mubbasir.github.io/HTP-benchmark/>.

13 1 Introduction

14 The study of human navigation has long been of interest to various research communities such as
15 computer graphics [10], computer vision [1], cognitive science [33], and robotics [5]. Advancements
16 in these areas have seen widespread practical application in pandemic response, architectural design,
17 urban planning, transportation engineering, crowd management, socially compliant robot navigation,
18 and entertainment. Accordingly, the influence of human navigation research has reached countless
19 individuals and will continue to do so in the foreseeable future.

20 Most applications rely on simulation models [20], which are sufficiently accurate to human behavior
21 and generalizable to unforeseen circumstances. However, the past five years of predictive modeling
22 in computer vision has achieved significantly better accuracy [23], giving it a strong potential to
23 overtake the longstanding models from computer graphics. This is largely due to the transition from
24 using unimodal, discriminative models [1] that predict a single future trajectory to using multimodal,
25 generative models [7, 24, 18] that predict a distribution of future trajectories, which captures the
26 inherent uncertainty in human decision-making [25, 4]. Despite the evolution of models, however, the
27 accuracy metrics that were introduced with the first unimodal models are still in use today. In order to
28 adapt these fundamentally unimodal metrics to multimodal models, the metrics are computed between
29 each predicted trajectory and the ground truth trajectory, and the minimum error for each metric is
30 reported. This results in a gross overestimation of accuracy that we later show is not consistent with
31 the expected accuracy, which may misguide future research efforts. Furthermore, the minimum value
32 is not actionable, because while it is evident that a state-of-the-art (SOTA) multimodal model can
33 find *an* accurate trajectory, it cannot determine *which* trajectory that is on unseen data. We measure
34 this uncertainty through a decidability metric.

35 Generalizability cannot be maximized by solely improving upon accuracy metrics. An inaccurate
36 model can be robust by producing realistic trajectories, and an accurate model can fail to be practicable

37 by being undecidable. Models can exist on the continuum between these two extremes, making it
 38 critical to consider realism and decidability metrics as well.

39 Furthermore, there is a stark class imbalance in existing datasets. While datasets are abundant in
 40 instances where humans are interacting with each other in open spaces [16, 22, 2, 34, 3, 14], they are
 41 significantly lacking in both environment information and instances where humans are interacting
 42 with their environment. Ultimately, this hinders generalization at a global level and has led to some
 43 models being developed without considering environments at all [1, 7].

44 In this work, we provide an augmented human trajectory prediction dataset that compensates for
 45 the lack of agent-to-environment interaction in existing datasets with a new simulated dataset. To
 46 understand model performance on this new dataset with more reliability and nuance, we propose a
 47 comprehensive set of accuracy, realism, and decidability metrics. A subset of these metrics are novel
 48 *multiverse metrics*, which are better-suited for multimodal models than existing metrics but are still
 49 applicable to unimodal models. The evaluation using these metrics decisively evidences that the new
 50 dataset facilitates better robustness and generalization, that current metrics can be misleading, and
 51 that there are still remaining challenges to modeling human trajectories. We finally showcase that
 52 realism metrics can also be used to decide which prediction to take from an undecidable multimodal
 53 model through the process of *Multimodal Model Collapse*. Henceforth, we refer to humans as agents,
 54 since our conceptual framework is broadly applicable, e.g. to robotic and vehicular agents.

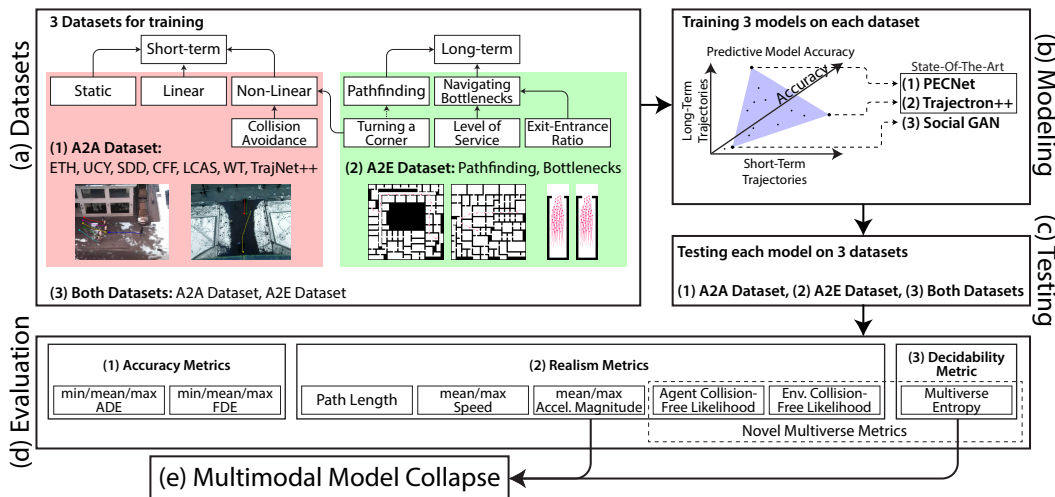


Figure 1: The above framework image shows (a) the differences between the trajectories of existing datasets (A2A) and the novel dataset (A2E), (b-c) the models trained and tested on combinations of A2A and A2E, (d) the proposed set of metrics for evaluating the accuracy, realism, and decidability of models, and (e) a greedy method for selecting the prediction most realistic movement.

55 2 Background and Preliminaries

56 **Models for Human Trajectory Prediction.** Earlier methods such as Social LSTM [1] and Social
 57 Attention [31] proposed a deterministic model which predict a future trajectory given observed
 58 trajectories. However, forecasting trajectories inherently introduce the uncertainty in the future, hence
 59 the utility of those uni-modal models which predict only one future trajectory is limited. Recent
 60 studies [7, 18, 24, 36, 12, 17] assume the multi-modalities in the future human behavior and predict
 61 its distribution to embody the uncertainty. In this paper, we focus on three SOTA methodologies to
 62 showcase our benchmark dataset: SocialGAN [7], PEcNet [18], and Trajectron++ [24].

63 SocialGAN [7] adopts GAN [6] framework to forecast possible future trajectories and it can avoid
 64 collisions among pedestrians by introducing a pooling mechanism that captures between-human
 65 interaction. PEcNet [18] solves the trajectory prediction problem by first modeling the future goal
 66 position distribution using a Variational Autoencoder (VAE) [13], and then predict the future positions
 67 by interpolating the observed positions and the estimated goal position. Trajectron++ [24] proposes

68 a graph structured recurrent model based on conditional VAE [28] to predict the future trajectories.
69 Further details can be found in the Supplementary Materials.

70 We investigate these three models as the representatives of the various SOTA works. We choose
71 them because PECNet [18] shows an outstanding performance on the long-term trajectory while
72 the short-term trajectory is most well predicted in Trajectron++ [24]. We expect SocialGAN [7], as
73 one of the earliest and most frequently referred models, to be a bound around existing models with
74 respect to PECNet and Trajectron++. Fig. 1.b shows the coverage comparison of SOTA models in
75 terms of the short- and long-term human trajectory prediction accuracy. We differentiate between
76 predictive models of short-term and long-term trajectories on the basis of goal conditioning. A model
77 that is not goal-conditioned will inherently increase in error as the predicted path length increases,
78 sometimes at an exponential rate [24], whereas goal-conditioned models are expected to predict long
79 paths without the same trade-off between path length and error.

80 **Datasets for Human Trajectory Prediction.** The computer vision and graphics community have
81 collected several human pedestrian trajectory datasets. ETH [21] and UCY [16] are commonly
82 used datasets that contain five outdoor scenes with jointly more than 1,600 pedestrian trajectories.
83 Stanford Drone Dataset (SDD) [22] consists of eight outdoor scenes tracking 19,000 targets including
84 pedestrians, bicyclists, skateboarders, cars, and buses collected from a drone. Stanford Crowd
85 Dataset (CFF) [2] consists of pedestrian trajectories collected within a train station building of size
86 $25\text{m} \times 100\text{m}$ for 12×2 hours captured by a distributed camera network. L-CAS 3D Point Cloud
87 People Dataset (LCAS) [34] consists of 28,002 scan frames collected within a university building
88 by a 3D LiDAR sensor mounted on a robot that is either stationary or moving. WILDTRACK
89 (WT) [3] is a collection of annotated dense pedestrian groups captured by seven static HD cameras
90 in a public square for about 60 minutes. The Supplementary Materials provide more details of
91 these datasets. Some datasets, such as TrajNet++ [14], augment upon existing datasets. TrajNet++
92 combines ETH/UCY, CFF, LCAS, and Wildtrack datasets, as well as a synthetic dataset generated by
93 ORCA [30].

94 Existing human trajectory datasets have limitations in the sense of embodying interactions. They either
95 do not contain agent-to-environment (A2E) interactions [3], or exhibit limited agent-to-agent (A2A)
96 interactions at small scale in simple environments. We speculate that many self-centered pedestrians
97 are prone to avoid or mitigate, consciously or unconsciously, the influence of the environments and
98 other pedestrians during their navigation. In this work, we are proposing datasets that augment A2E
99 and A2A interactions, which may bring benefits for enhancing learning models by encoding more
100 complex trajectory dynamics.

101 **Benchmarks for Human Trajectory Prediction.** In computer graphics community [27], trajectories
102 are, in general, measured by motion statistics such as the number of collisions, average speed, average
103 acceleration, and total distance traveled. On the other hand, in machine learning community [14,
104 1, 7], the most commonly used evaluation metrics for trajectory forecasting models are Average
105 Displacement Error (ADE) and Final Displacement Error (FDE). ADE is the average L_2 distance
106 between the ground truth and the predicted trajectories across all future steps. FDE is the L_2 distance
107 between the ground truth final destination and the predicted final destination at the end of the future
108 steps. More evaluation metrics in machine learning community are discussed in Supplementary
109 Materials.

110 ADE and FDE are applicable to unimodal methods which predict only one future sequence that can be
111 compared with the ground truth future sequence. However, as aforementioned in this section, many
112 multimodal trajectory forecasting models assuming uncertainty and multimodality in pedestrians’
113 future behaviors predict k future sequences (usually $k = 20$). Most of these models report the
114 minimum ADE / FDE results among all k predictions, which, in our view, is over optimistic. Not
115 only is this a significant underestimation of the error, but it is also an impossible standard in that
116 these models are incapable of choosing the prediction with the minimum error. In Section 4 of this
117 work, we propose new metrics that can tackle this issue.

118 3 Agent-to-Agent and Agent-to-Environment Interaction Dataset

119 We propose a comprehensive trajectory prediction dataset **A2X** that consists of a representative set
120 of trajectories, which will enable better generalization under realistic circumstances that are either
121 complex or unsafe and out-of-distribution (OOD) with respect to current datasets.

122 In order to understand what the shortcomings of current datasets are (Sec. 2), we first taxonomize the
 123 characteristics of human trajectories. The TrajNet++ benchmark [14] proposed an initial taxonomy
 124 that only considers short-term characteristics, e.g., standing still, moving linearly, or avoiding
 125 collisions (Fig. 1.a). While the original taxonomy is sufficient for describing the trajectories in many
 126 real datasets and their agent-to-agent (A2A) interactions, models that learn exclusively from these
 127 types are insufficient for most applications, which consider environments that have non-navigable
 128 regions and time frames longer than 5 seconds, which is the practical limit for most models before
 129 they become exponentially erroneous [24]. We have improved upon this by considering long-term
 130 characteristics (Fig. 1.a), i.e., pathfinding alone and navigating through crowded bottlenecks. These
 131 types of trajectories emerge from agent-to-environment (A2E) interactions, which unfold over a
 132 longer time frame than A2A interactions and are essential for navigation within any environment [29].

133 3.1 Agent-to-Agent Interactions

134 For representing A2A interactions, we make use of each prior dataset described in Section 2:
 135 ETH [16], UCY [16], SDD [22], CFF [2], LCAS [34], WT [3], and TrajNet++ [14]. These datasets
 136 feature transient interactions between agents and little interaction with the environment, which is
 137 made difficult to measure by the frequent unavailability of environment information. Therefore, we
 138 approximate environment information based on the principle of stigmergy [19, 11], which observes
 139 the self-organization of human navigation along trails. For each position that agents have traveled
 140 through in either the training or testing sets of the ground truth, a 1-meter radius around the position
 141 is considered to be navigable. This guarantees that predictions with less than 1 meter of displacement
 142 from the ground truth at all times will never intersect with the environment. In addition, in order to
 143 compensate for the imbalance between A2A and A2E interactions in prior datasets, we propose the
 144 generation of synthetic data in addition to that of TrajNet++. While real datasets are valuable for their
 145 veridicality, there are logistical limitations that prevent the acquisition of real data in OOD scenarios
 146 that are unsafe for human participants or prohibitively expensive from an organizational standpoint.

147 3.2 Agent-to-Environment Interactions

148 Two such scenarios are used to sample trajectories exhibiting A2E interactions: (1) pathfinding alone
 149 in a large, complex environment, which has prohibitive logistical cost and (2) navigating through
 150 bottlenecks of varied width with a dense crowd, which can be unsafe. Though simulation models
 151 are normally less accurate than predictive models in predicting human trajectories [1], the prevalent
 152 Social Force model [10] currently outperforms predictive models in terms of robustness, has been
 153 used in several application domains [5, 32, 35], and has ecological validity in these A2E scenarios,
 154 which have not had sufficient real data for training predictive models until A2X.

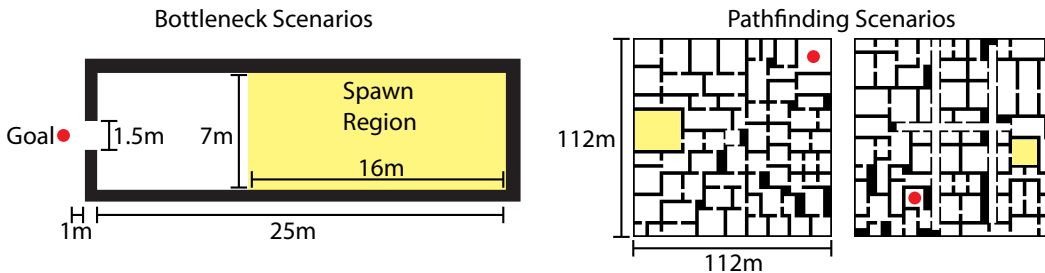


Figure 2: The above images show the exact dimensions of environments from the bottleneck and pathfinding scenarios in A2E.

155 We leverage the Social Force model to simulate 236 scenarios of a single agent navigating between
 156 random points in complex $112 \times 112 \text{ m}^2$ environments from [29] (Fig. 2). This produces long-term
 157 isolated interactions between single agents and the environment. We then use the same model
 158 to simulate well-studied bottleneck scenarios [26, 9] in a $25 \times 7 \text{ m}^2$ room that vary in terms of
 159 (a) the density of agents (Level of Service) from $\{0.2, 0.4, 0.6, 0.8, 1.0\}$ agents/ m^2 and (b) the
 160 ratio between the width of the bottleneck and the width of the room (Exit-Entrance Ratio) from
 161 $\{0.2, 0.3, 0.4, 0.6, 0.7\}$ (Fig. 2). A total of 398 scenarios have been generated across all combinations
 162 of Level of Service and Exit-Entrance Ratio. This produces long-term interactions between agents as
 163 a result of the constricting environment. Exact environment information has been provided for both

164 types of scenarios. We later show that current models trained on existing A2A datasets are unable to
 165 generalize to these critical scenarios, but with the addition of training data on these scenarios, the
 166 accuracy of predictions significantly improves.

167 4 Accuracy, Realism, and Decidability of Human Trajectory Prediction

168 We propose a total of 15 accuracy, realism, and decidability metrics (Fig. 1.d). These metrics are either
 169 borrowed from computer vision and computer graphics literature [21, 1, 27, 8] or newly developed
 170 *multiverse metrics*, which assess the A2A and A2E interactions of both multimodal models with
 171 $k > 1$ and unimodal models with $k = 1$.

172 4.1 Preliminaries

173 In accordance with both unimodal and multimodal predictive models, we
 174 utilize the following notation for their predictions. A prediction scenario is
 175 defined by a set of n agents present in an environment \mathbf{E} at the same time.
 176 Each agent a has t_p frames of past position data as input and t_f frames of
 177 future position data for ground truth $\mathbf{Y}_{a,0} \in \mathbb{R}^{t_f \times 2}$ and for each prediction
 178 $\hat{\mathbf{Y}}_{a,j} \in \mathbb{R}^{t_f \times 2}$, where $0 \leq j < k$. All position data is in meters and has a
 179 frame rate of $1/\Delta t$ hertz based on the dataset. The position at the t -th frame
 180 is $\mathbf{Y}_{a,0,t} \in \mathbb{R}^2$ for the ground truth and $\hat{\mathbf{Y}}_{a,j,t} \in \mathbb{R}^2$ for prediction j , where
 181 $0 \leq t < t_f$. We then compute the velocities corresponding to the ground truth
 182 $\mathbf{V}_{a,0} \in \mathbb{R}^{(t_f-1) \times 2}$ and each prediction $\hat{\mathbf{V}}_{a,j} \in \mathbb{R}^{(t_f-1) \times 2}$.

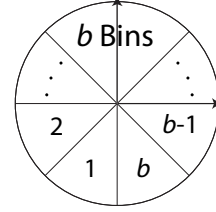


Figure 3: This image shows how $b = 8$ bins would be arranged in 2D space.

183 Many of the following metrics make use of aggregate functions. For any d -dimensional vector
 184 $\mathbf{v} \in \mathbb{R}^d$, we denote the minimum value by $\Omega(\mathbf{v})$, the mean value by $\Theta(\mathbf{v})$, and the maximum value
 185 by $O(\mathbf{v})$. For a matrix of d -many 2D vectors $\mathbf{V} \in \mathbb{R}^{d \times 2}$, function $\Xi(\mathbf{V}, b)$ transforms the 2D vectors
 186 into a probability distribution $\mathbf{p} \in \mathbb{R}^b$ over a vector of b -many equiangular bins, which radiate from
 187 the origin (Fig. 3). Finally, we denote the L_2 norm by $\|\cdot\|$.

188 4.2 Accuracy Metrics: Comparison to Ground Truth

189 Accuracy metrics from computer vision literature are responsible for comparing the ground truth with
 190 the predictions based on the displacement error.

191 **Average Displacement Error (ADE).** ADE is computed for each prediction j as \mathbf{a}_j , the average
 192 distance between a position in the ground truth and a position in the prediction across t_f frames
 193 (Eq. 1) [21]. It is then aggregated across the k predictions in three ways: minimum, mean, and
 194 maximum, which offers a more reliable expectation of a model's accuracy than the minimum alone.

195 **Final Displacement Error (FDE).** FDE is computed for each prediction j as \mathbf{b}_j , the distance
 196 between the final positions of the ground truth and the prediction (Eq. 2) [1]. It is aggregated across
 197 the k predictions in the same ways as ADE for better reliability.

$$\text{ADE}(\mathbf{Y}_a, \hat{\mathbf{Y}}_a) = [\Omega(\mathbf{a}), \Theta(\mathbf{a}), O(\mathbf{a})] \quad (1)$$

$$s.t. \mathbf{a}_j = \frac{1}{t_f} \sum_{t=0}^{t_f-1} \left\| \mathbf{Y}_{a,0,t} - \hat{\mathbf{Y}}_{a,j,t} \right\|, 0 \leq j < k$$

$$\text{FDE}(\mathbf{Y}_a, \hat{\mathbf{Y}}_a) = [\Omega(\mathbf{b}), \Theta(\mathbf{b}), O(\mathbf{b})] \quad (2)$$

$$s.t. \mathbf{b}_j = \left\| \mathbf{Y}_{a,0,t_f-1} - \hat{\mathbf{Y}}_{a,j,t_f-1} \right\|, 0 \leq j < k$$

198 4.3 Realism Metrics: Motion and Interaction Statistics

199 Realism metrics are used to describe the movement and interactions within the ground truth and
 200 the predictions separately. These metrics can then be used to uncover more nuanced differences
 201 between the ground truth and predictions. While they cannot ensure that predictions are accurate,

202 they can ensure that predictions are realistic in their movement and plausible. Every realism metric is
 203 computed in the same way for both the ground truth and predictions, so \mathbf{Y} is interchangeable with $\widehat{\mathbf{Y}}$
 204 and \mathbf{V} with $\widehat{\mathbf{V}}$. For generality, we consider the ground truth as a unimodal model with $k = 1$, but we
 205 refer to it as having k paths instead of predictions.

206 The following motion statistics are used to describe the movement of agent a in either the ground truth
 207 or averaged across the k predictions. They have been used to evaluate crowd simulations in computer
 208 graphics research [27], but have not yet been used to evaluate predictive models in computer vision.

209 **Path Length.** The average path length (m) for an agent a is computed by first finding the length of
 210 each path j and then averaging the values across all k paths (Eq. 3).

211 **Speed.** In order to report the speed (m/s), the magnitudes $\mathbf{S} \in \mathbb{R}^{k \times (t_f - 1)}$ of velocities in \mathbf{V}_a are first
 212 computed for each agent a . Next, two values are reported for speed: the mean speed averaged across
 213 k paths and the maximum speed averaged across k paths. For each path j of agent a , the mean and
 214 maximum speed are computed across $t_f - 1$ frames (Eq. 4).

215 **Acceleration Magnitude.** Similar to speed, we first compute the magnitudes $\mathbf{A} \in \mathbb{R}^{k \times (t_f - 2)}$ of the
 216 difference between every pair of consecutive velocities in \mathbf{V}_a for each agent a . The acceleration
 217 magnitude (m/s²) $\mathbf{A}(\mathbf{V}_a)$ is then reported in the same way as speed: the mean acceleration magnitude
 218 averaged across k paths and the maximum magnitude averaged across k paths (Eq. 5).

$$L(\mathbf{Y}_a) = \left[\frac{1}{k} \sum_{j=0}^{k-1} \sum_{t=0}^{t_f-2} \left\| \mathbf{Y}_{a,j,t+1} - \mathbf{Y}_{a,j,t} \right\| \right] \quad (3)$$

$$S(\mathbf{V}_a) = \left[\frac{1}{k} \sum_{j=0}^{k-1} \Theta(\mathbf{S}_j), \frac{1}{k} \sum_{j=0}^{k-1} O(\mathbf{S}_j) \right] \quad (4)$$

$$s.t. \quad \mathbf{S}_{j,t} = \left\| \mathbf{V}_{a,j,t} \right\|, \quad 0 \leq t < t_f - 1$$

$$\mathbf{A}(\mathbf{V}_a) = \left[\frac{1}{k} \sum_{j=0}^{k-1} \Theta(\mathbf{A}_j), \frac{1}{k} \sum_{j=0}^{k-1} O(\mathbf{A}_j) \right] \quad (5)$$

$$s.t. \quad \mathbf{A}_{j,t} = \left\| (\mathbf{V}_{a,j,t+1} - \mathbf{V}_{a,j,t}) / \Delta t \right\|, \quad 0 \leq t < t_f - 2$$

219 Traditional measures of collision are unsuitable for multimodal models in which an agent a may be
 220 colliding with agent b when it takes the direction of path j , but not when it takes the direction of path
 221 $j + 1$. We therefore propose multiverse metrics such as Agent Collision-Free Likelihood (ACFL)
 222 and Environment Collision-Free Likelihood (ECFL) to measure the A2A and A2E interactions of
 223 multimodal models respectively.

224 **Agent Collision-Free Likelihood (ACFL).** In order to assess the quality of A2A interaction under
 225 the k^n possible futures for n agents, we propose ACFL, which computes the probability that agent a
 226 has a path that is free of collision in all of the $k^{(n-1)}$ possible futures with other agents (Eq. 6). The
 227 indicator function $\mathbf{1}_{\mathbb{R}>0}$ returns 1 when the distance between agents a and b is greater than r meters
 228 at time t , and 0 otherwise. This means that if their centers of mass are within r meters of each other,
 229 they are considered to be colliding. For analysis, r has been set to 0.3 meters (~ 1 foot).

230 **Environment Collision-Free Likelihood (ECFL).** ECFL complements ACFL in that it measures the
 231 quality of A2E interaction under the k possible futures that agent a can interact with the environment
 232 (Eq. 7). Namely, it reports the probability that agent a has a path that is free of collision with the
 233 environment. The environment is represented by a binary matrix \mathbf{E} , in which each cell corresponds
 234 to a square space and is equal to 1 if that space is navigable and 0 otherwise. $\mathbf{E}[0, 0]$ is aligned with
 235 the origin of the position data \mathbf{Y} , but \mathbf{E} has a scale of $1/s$ meters per unit as opposed to 1 meter per
 236 unit like \mathbf{Y} . This means that the position $[x, y] = \mathbf{Y}_{a,j,t}$ of agent a taking path j at time t maps to
 237 $\mathbf{E}[\lfloor s \cdot y \rfloor, \lfloor s \cdot x \rfloor]$. For analysis, s has been set to 2 based on the dataset. When agent a 's center of
 238 mass is intersecting a non-navigable region of the environment like a wall, the agent is considered to
 239 be colliding with the environment.

$$\text{ACFL}(\mathbf{Y}, a) = \left[\frac{1}{k} \sum_{j=0}^{k-1} \prod_{b=0}^{n-1} \prod_{i=0}^{k-1} \prod_{t=0}^{t_f-1} \mathbf{1}_{\mathbb{R}^{>0}} \left(\left\| \mathbf{Y}_{a,j,t} - \mathbf{Y}_{b,i,t} \right\| - r \right) \right] \text{ s.t. } a \neq b \quad (6)$$

$$\text{ECFL}(\mathbf{Y}_a, \mathbf{E}) = \left[\frac{1}{k} \sum_{j=1}^k \prod_{t=0}^{t_f-1} \mathbf{E} \left[[s \cdot \mathbf{Y}_{a,j,t,1}], [s \cdot \mathbf{Y}_{a,j,t,0}] \right] \right] \quad (7)$$

$$\text{MVE}(\mathbf{Y}_a) = - \sum_{\mathbf{p} \in \mathbf{p}} p \cdot \log_2(p) \text{ s.t. } \mathbf{p} = \Xi(\mathbf{D}, 20), \quad (8)$$

$$\mathbf{D}_j = \frac{1}{t_f - 1} \left(\sum_{t=1}^{t_f-1} \mathbf{Y}_{a,j,t} \right) - \mathbf{Y}_{a,j,0}, \quad 0 \leq j < k$$

240 4.4 Decidability Metric: Certainty in Movement Direction

241 Decidability is a measure of a model’s uncertainty in the movement direction of agents, and it is not
 242 strictly opposite between unimodal and multimodal models. If a multimodal model has low enough
 243 uncertainty in an agent’s direction of movement, we consider it to be decidable.

244 **Multiverse Entropy (MVE).** We compute MVE to measure the decidability for agent a . We first
 245 transform each path j into an average direction vector $\mathbf{D}_j \in \mathbb{R}^2$ as the vector from the initial position
 246 $\mathbf{Y}_{a,j,0}$ to the average position of the $t_f - 1$ subsequent points (Eq. 8). The average direction vectors
 247 \mathbf{D} are then transformed into a probability distribution $\mathbf{p} \in \mathbb{R}^b$ over a vector of b -many equiangular
 248 bins (Fig. 3). Finally, the entropy of \mathbf{p} is reported as MVE. High values of ACFL and ECFL are
 249 contingent on low MVE (high decidability), because high certainty in the direction that an agent
 250 will travel along will cause fewer potential collisions with other agents (ACFL) and the environment
 251 (ECFL). For experimental purposes, b has been set to k , so that MVE is maximized when every
 252 prediction is in a different direction.

253 4.5 Comparing Realism Metrics

254 In order to compare realism metrics between the ground truth and predictions for an agent a , we
 255 first compute a feature vector for the ground truth $\mathbf{F}_a = \langle \mathbf{L}(\mathbf{Y}_{a,0}), \mathbf{S}(\mathbf{V}_a), \mathbf{A}(\mathbf{V}_a), \text{ACFL}(\mathbf{Y}, a),$
 256 $\text{ECFL}(\mathbf{Y}_a, \mathbf{E}) \rangle$, where $\langle \cdot, \cdot \rangle$ denotes vector concatenation. The same vector concatenation is used
 257 to compute the feature vector $\widehat{\mathbf{F}}_{a,j} \in \mathbb{R}^7$ for each prediction j . Equation 9 returns the percent
 258 differences $\widehat{\mathbf{C}}_a \in \mathbb{R}^k$ between the feature vectors of each prediction j and the ground truth of agent a .

$$\widehat{\mathbf{C}}_{a,j} = \frac{100}{7} \sum_{f=0}^6 \frac{|\widehat{\mathbf{F}}_{a,j,f} - \mathbf{F}_{a,0,f}|}{\mathbf{F}_{a,0,f}} \text{ s.t. } \mathbf{F}_{a,0,f} > 0, \quad 0 \leq j < k \quad (9)$$

259 5 Results

260 In order to understand the limits of not only the SOTA but also the models that paved the way towards
 261 the SOTA, we evaluate three critical multimodal models that are capable of either short-term or
 262 long-term trajectory prediction and provide a large coverage over the performance of prior models
 263 (Fig. 1.b). In particular, we have selected (1) Social GAN (SGAN) [7], one of the earliest models;
 264 (2) Trajectron++ (T++) [24], a SOTA model for short-term trajectory prediction; and (3) PECNet
 265 (PECN) [18], a SOTA model for long-term trajectory prediction.

266 **Training Protocol.** Each of the three models was trained on 3 combinations from the **A2X** Dataset:
 267 A2A interaction, A2E interaction, and both (Fig. 1.b), producing a total of 9 models. Each trained
 268 model was then evaluated on the testing sets of the 3 combinations (Fig. 1.c). The results of the
 269 evaluations on A2A and A2E are reported in Table 1, while the results on both A2A and A2E
 270 combined and corresponding visualizations are reported in the Supplementary Materials. According
 271 to the dataset, the following parameters have been set for the evaluation: $k = 20$, $t_p = 8$, $t_f = 12$,

272 and $\Delta t = 0.4$, meaning that each agent is receiving 3.2 seconds of input data and predicting 4.8
 273 seconds into the future.

274 Each row of Table 1 reports the accuracy, realism, and decidability metrics of a model averaged
 275 across the agents of every testing scenario for a given dataset. The first 5 columns of realism metrics
 276 correspond to the dimensions of \mathbf{F} and $\hat{\mathbf{F}}$, the feature vectors used to compute the percent difference
 277 between the ground truth (GT) and predictions. The mean percent difference $\Theta(\hat{\mathbf{C}}_a)$ of each agent a
 278 is averaged across all agents and reported in the final column of the realism metrics. For all accuracy
 279 metrics, the realism percent difference, and the decidability metric, a lower value is favorable, while
 280 for the remaining realism metrics, a value closer to the ground truth is favorable.

Test	Model	Train	Accuracy Metrics				Realism Metrics					Decidab. MVE ↓	
			ADE ↓			FDE ↓	Length	Speed mean / max	Accel. mean / max	ACFL	ECFL		%Diff. ↓
			min / mean / max	min / mean / max									
Agent-to-Agent Interaction	GT	N/A	0.00 / 0.00 / 0.00	0.00 / 0.00 / 0.00		4.43	1.01 / 1.32	0.29 / 1.04	0.95	1.00	0	0.00	
	SGAN	A2A	0.36 / 0.77 / 1.50	0.62 / 1.61 / 3.33		4.22	0.96 / 1.42	0.09 / 0.56	0.30	0.98	48	0.90	
		A2E	2.21 / 2.48 / 2.81	4.02 / 4.65 / 5.48		3.15	0.72 / 1.38	0.12 / 0.40	0.58	0.97	51	0.70	
		Both	0.37 / 0.74 / 1.35	0.65 / 1.55 / 2.97		4.13	0.94 / 1.32	0.06 / 0.33	0.33	0.98	51	0.84	
	PECN	A2A	0.63 / 0.65 / 0.68	1.12 / 1.28 / 1.45		4.50	1.02 / 2.15	0.48 / 3.41	0.56	0.98	56	0.07	
		A2E	1.25 / 1.28 / 1.31	1.83 / 2.00 / 2.20		4.50	1.02 / 4.16	1.13 / 8.80	0.59	0.98	166	0.10	
		Both	0.73 / 0.76 / 0.79	1.44 / 1.59 / 1.74		4.78	1.08 / 2.61	0.49 / 4.57	0.57	0.98	85	0.10	
	T++	A2A	0.22 / 0.66 / 1.85	0.42 / 1.51 / 4.16		4.38	1.00 / 2.32	0.36 / 3.09	0.22	0.98	47	1.08	
		A2E	0.56 / 1.06 / 1.77	1.13 / 2.29 / 3.90		4.22	0.96 / 1.79	0.29 / 2.18	0.25	0.98	46	1.41	
Both		0.23 / 0.64 / 1.76	0.43 / 1.48 / 4.02		4.35	0.99 / 2.27	0.35 / 2.96	0.22	0.98	47	1.13		
Agent-to-Env. Interaction	GT	N/A	0.00 / 0.00 / 0.00	0.00 / 0.00 / 0.00		5.51	1.25 / 1.40	0.18 / 0.51	1.00	1.00	0	0.00	
	SGAN	A2A	0.28 / 0.66 / 1.33	0.50 / 1.48 / 3.14		5.42	1.23 / 1.70	0.08 / 0.45	0.29	0.90	47	0.82	
		A2E	0.19 / 0.41 / 0.96	0.27 / 0.86 / 2.17		4.19	0.95 / 1.33	0.09 / 0.28	0.35	0.94	48	0.64	
		Both	0.19 / 0.56 / 1.25	0.32 / 1.28 / 3.02		5.03	1.14 / 1.57	0.08 / 0.40	0.32	0.92	49	0.65	
	PECN	A2A	0.47 / 0.49 / 0.51	0.98 / 1.12 / 1.27		5.35	1.22 / 1.72	0.32 / 2.79	0.64	0.92	117	0.03	
		A2E	0.29 / 0.31 / 0.34	0.63 / 0.75 / 0.90		5.64	1.28 / 2.44	0.40 / 3.50	0.60	0.94	148	0.04	
		Both	0.32 / 0.34 / 0.37	0.70 / 0.81 / 0.92		5.64	1.28 / 2.29	0.34 / 3.41	0.60	0.93	157	0.06	
	T++	A2A	0.17 / 0.81 / 2.43	0.34 / 1.86 / 5.54		5.48	1.25 / 3.10	0.53 / 4.41	0.18	0.90	43	1.24	
		A2E	0.10 / 0.29 / 0.64	0.19 / 0.69 / 1.61		5.41	1.23 / 1.63	0.18 / 1.38	0.47	0.95	40	0.73	
Both		0.12 / 0.37 / 1.11	0.23 / 0.87 / 2.55		5.41	1.23 / 2.00	0.27 / 2.04	0.42	0.93	40	0.76		

Table 1: This table showcases the evaluation results of Social GAN (SGAN), PECNet (PECN), and Trajectron++ (T++) after training on either A2A, A2E, or both A2A and A2E and testing on A2A and A2E separately. For every metric in a testing set, the best value has been made bold for each model.

281 **Analysis.** As expected, we find that models trained on a single type of interaction perform poorly on
 282 test scenarios that feature the other type of interaction (Tab. 1). By training any of the three models
 283 on both types of interactions, we find that the accuracy of this model is either nearly the highest or
 284 the highest according to mean ADE/FDE compared to the same model trained on either A2A or A2E.
 285 For instance, $T++_{Both}$ trained on both types of interactions achieves the lowest mean ADE on A2A
 286 across all 9 trained models.

287 However, we cannot rely only on the accuracy of models to determine which is best, since anything
 288 short of perfect accuracy carries risk. The realism metrics allow us to better understand the model’s
 289 performance in the context of its application. For example, we find that the maximum speed and
 290 acceleration for $T++_{Both}$ are significantly higher than the ground truth, which for an application in
 291 socially compliant robot navigation can discomfort or potentially harm surrounding humans [15]. In
 292 contrast, $SGAN_{Both}$ has lower average accuracy by a small margin, but it boasts higher realism by a
 293 large margin in terms of maximum speed, maximum acceleration magnitude, and ACFL. We attribute
 294 $SGAN_{Both}$ ’s higher ACFL to the tighter spread of its predictions than $T++_{Both}$ according to MVE.
 295 Ultimately, the choice of a model depends on the application, but without the joint consideration of
 296 the proposed accuracy and realism metrics, a practitioner may be led to choose an unsuitable model.

297 We have made 5 other notable observations from Table 1. (1) There are instances of models
 298 (highlighted in red) where relying on the optimistic lens of existing evaluations (i.e., minimum ADE
 299 and FDE) would lead to choosing models that are less accurate than others on average. (2) Models
 300 trained exclusively on A2E interactions tend to have lower likelihoods of A2A collision (higher
 301 ACFL) than models trained on A2A interactions alone or on both types of interactions, highlighting
 302 the important of A2E for improving robustness even in OOD scenarios such as A2A. (3) While this
 303 also holds true for the likelihood of A2E collision (ECFL) when testing on A2E, we find that ECFL
 304 is nearly perfect for A2A scenarios, indicating that A2A scenarios do not challenge models with A2E

interactions. (4) PECNet has the highest ACFL by an enormous margin owing to its MVE, which is low enough to consider PECNet as decidable and likely helps it in performing long-term trajectory prediction. Finally, (5) models trained on both types of interactions do not yet generalize to A2E better than models trained on A2E alone as some models have for A2A, meaning that there is still much room for improvement.

Multimodal Model Collapse (MMC). Accuracy metrics cannot be computed on never-before-seen data, because the ground truth is unknown. Consequently, it becomes impossible to find the predicted path with minimum error in accuracy and selecting an arbitrary prediction risks the maximum error. We therefore propose MMC, a baseline greedy method which can make use of the realism metrics to collapse the k predictions of an undecidable multimodal model into a single well-informed prediction. In particular, we rely on the proposed comparison of realism metrics (Sec. 4.5), but instead of computing F_a from ground truth testing data $Y_{a,0}$ for each agent a , we compute it as the average across *all* agents in the ground truth *training* data from the same environment. We then replace the k predictions \hat{Y}_a with the single prediction j that minimizes the percent difference $\hat{C}_{a,j}$ for each agent a , which is the closest in realism to prior ground truth for the same type of scenario (Eq. 9). This, certainly, does not guarantee the optimal selection for a single agent. But it minimizes the overall error in selecting predictions for all agents. Table 2 shows the result of applying this technique to all 9 models. On average, we find that the ADE/FDE of the collapsed prediction is only $\sim 15.76\%$ worse than the mean ADE/FDE of the uncollapsed predictions, and $\sim 31.63\%$ better than the maximum ADE/FDE. Although the accuracy of the most realistic prediction is lower than the average accuracy over 20 predictions, its performance is consistently much better than the worst-case and it ultimately makes the undecidable model applicable to unseen data.

Test	Model	Train	Accuracy Metrics				Realism Metrics					Decidab.		
			ADE ↓		FDE ↓		Length	Speed mean / max	Accel. mean / max	ACFL	ECFL		%Diff. ↓	MVE ↓
			min = mean = max	min = mean = max	min = mean = max	min = mean = max								
Agent-to-Agent Interaction	GT	N/A	0.00	0.00	4.43	1.01 / 1.32	0.29 / 1.04	0.95	1.00	0	0.00			
	SCAN	A2A	0.91	1.99	4.28	0.97 / 1.20	0.16 / 0.41	0.69	0.99	37	0.00			
		A2E	2.57	4.97	3.75	0.85 / 1.32	0.20 / 0.37	0.79	0.97	40	0.00			
		Both	0.86	1.86	4.25	0.97 / 1.15	0.11 / 0.23	0.70	0.99	41	0.00			
	PECN	A2A	0.65	1.27	4.44	1.01 / 1.56	0.33 / 1.79	0.66	0.98	56	0.00			
		A2E	1.28	2.03	4.33	0.98 / 3.23	1.02 / 6.37	0.68	0.98	166	0.00			
		Both	0.76	1.55	4.70	1.07 / 2.12	0.44 / 3.18	0.64	0.98	85	0.00			
	T++	A2A	0.81	1.83	4.51	1.03 / 1.31	0.44 / 0.98	0.66	0.99	26	0.00			
		A2E	1.05	2.27	4.53	1.03 / 1.32	0.42 / 0.97	0.63	0.98	30	0.00			
		Both	0.81	1.84	4.51	1.03 / 1.31	0.44 / 1.00	0.65	0.99	26	0.00			
	Agent-to-Env. Interaction	GT	N/A	0.00	0.00	5.51	1.25 / 1.40	0.18 / 0.51	1.00	1.00	0	0.00		
		SCAN	A2A	0.76	1.84	5.00	1.14 / 1.44	0.15 / 0.33	0.63	0.96	38	0.00		
A2E			0.69	1.60	4.73	1.08 / 1.30	0.13 / 0.23	0.68	0.98	40	0.00			
Both			0.73	1.77	4.55	1.03 / 1.36	0.16 / 0.27	0.66	0.97	40	0.00			
PECN		A2A	0.49	1.11	5.39	1.22 / 1.45	0.25 / 1.10	0.69	0.93	117	0.00			
		A2E	0.30	0.71	5.54	1.26 / 1.71	0.31 / 1.41	0.62	0.93	148	0.00			
		Both	0.34	0.78	5.60	1.27 / 1.97	0.32 / 1.41	0.64	0.94	157	0.00			
T++		A2A	0.90	2.06	4.99	1.13 / 1.48	0.57 / 1.27	0.46	0.97	31	0.00			
		A2E	0.34	0.86	5.36	1.22 / 1.44	0.29 / 0.85	0.61	0.98	24	0.00			
		Both	0.52	1.20	5.34	1.21 / 1.48	0.41 / 0.99	0.57	0.97	28	0.00			

Table 2: This table reports the results of MMC on each of the 9 trained models. On average, MMC produces predictions that are consistently better than the worst case prediction prior to MMC. Only one value is reported for ADE and FDE, because the minimum, mean, and maximum are equal when $k = 1$. The MVE is always 0 when $k = 1$.

6 Conclusion

With the growing attention toward human trajectory prediction, it has become more important than ever to unify future research efforts in the right direction in terms of datasets and benchmark. In this work, we have brought to light the shortcomings of existing datasets, which hinder generalization, and existing evaluation metrics, which misrepresent model performance. By augmenting existing datasets with scenarios that feature substantial interactions between pedestrian agents and the environment, we have evidenced that models can generalize better. By proposing a comprehensive set of novel and existing evaluation metrics, we have not only proven the unreliability of existing evaluation metrics, but also highlighted the subtle factors that are essential for choosing the best trajectory prediction model for a particular application. Together, these contributions show that there is still room for much improvement even among the SOTA models.

References

- 338
- 339 [1] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese. Social LSTM: Human
340 Trajectory Prediction in Crowded Spaces. In *Proceedings of the IEEE Conference on Computer Vision and
341 Pattern Recognition (CVPR)*, pages 961–971, 2016.
- 342 [2] A. Alahi, V. Ramanathan, and L. Fei-Fei. Socially-Aware Large-Scale Crowd Forecasting. In *Proceedings
343 of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2203–2210, 2014.
- 344 [3] T. Chavdarova, P. Baqué, S. Bouquet, A. Maksai, C. Jose, T. Bagautdinov, L. Lettry, P. Fua, L. Van Gool,
345 and F. Fleuret. Wildtrack: A Multi-Camera HD Dataset for Dense Unscripted Pedestrian Detection.
346 In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages
347 5030–5039, 2018.
- 348 [4] R. K. Dubey, S. S. Sohn, C. Hoelscher, and M. Kapadia. Fusion-Based Wayfinding Prediction Model for
349 Multiple Information Sources. In *2019 22th International Conference on Information Fusion (FUSION)*,
350 pages 1–8. IEEE, 2019.
- 351 [5] G. Ferrer, A. Garrell, and A. Sanfeliu. Social-aware robot navigation in urban environments. In *2013
352 European Conference on Mobile Robots*, pages 331–336. IEEE, 2013.
- 353 [6] I. J. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and
354 Y. Bengio. Generative Adversarial Nets. In *Proceedings of the 27th International Conference on Neural
355 Information Processing Systems (NIPS)*, page 2672–2680, 2014.
- 356 [7] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi. Social GAN: Socially Acceptable Trajectories
357 with Generative Adversarial Networks. In *Proceedings of the IEEE Conference on Computer Vision and
358 Pattern Recognition (CVPR)*, pages 2255–2264, 2018.
- 359 [8] S. J. Guy, J. Van Den Berg, W. Liu, R. Lau, M. C. Lin, and D. Manocha. A Statistical Similarity Measure
360 for Aggregate Crowd Dynamics. *ACM Transactions on Graphics (TOG)*, 31(6):1–11, 2012.
- 361 [9] B. Haworth, M. Usman, G. Berseth, M. Kapadia, and P. Faloutsos. Evaluating and optimizing level of
362 service for crowd evacuations. In *Proceedings of the 8th ACM SIGGRAPH Conference on Motion in
363 Games*, pages 91–96, 2015.
- 364 [10] D. Helbing and P. Molnar. Social Force Model for Pedestrian Dynamics. *Physical review E*, 51(5):4282,
365 1995.
- 366 [11] D. Helbing, F. Schweitzer, J. Keltsch, and P. Molnar. Active walker model for the formation of human and
367 animal trail systems. *Physical review E*, 56(3):2527, 1997.
- 368 [12] B. Ivanovic and M. Pavone. The Trajectron: Probabilistic Multi-Agent Trajectory Modeling With Dynamic
369 Spatiotemporal Graphs. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*, pages
370 2375–2384, 2019.
- 371 [13] D. P. Kingma and M. Welling. Auto-Encoding Variational Bayes. In *2nd International Conference on
372 Learning Representations (ICLR), Banff, AB, Canada, April 14-16, 2014, Conference Track Proceedings*,
373 2014.
- 374 [14] P. Kothari, S. Kreiss, and A. Alahi. Human Trajectory Forecasting in Crowds: A Deep Learning Perspective.
375 *IEEE Transactions on Intelligent Transportation Systems*, 2021. doi: 10.1109/TITS.2021.3069362.
- 376 [15] T. Kruse, A. K. Pandey, R. Alami, and A. Kirsch. Human-aware robot navigation: A survey. *Robotics and
377 Autonomous Systems*, 61(12):1726–1743, 2013.
- 378 [16] A. Lerner, Y. Chrysanthou, and D. Lischinski. Crowds by Example. In *Computer Graphics Forum*,
379 volume 26, pages 655–664. Wiley Online Library, 2007.
- 380 [17] K. Mangalam, Y. An, H. Girase, and J. Malik. From Goals, Waypoints Paths To Long Term Human
381 Trajectory Forecasting. *arXiv preprint arXiv:2012.01526*, 2020.
- 382 [18] K. Mangalam, H. Girase, S. Agarwal, K.-H. Lee, E. Adeli, J. Malik, and A. Gaidon. It is Not the Journey
383 but the Destination: Endpoint Conditioned Trajectory Prediction. *arXiv preprint arXiv:2004.02025*, 2020.
- 384 [19] H. V. D. Parunak. A survey of environments and mechanisms for human-human stigmergy. In *International
385 workshop on environments for multi-agent systems*, pages 163–186. Springer, 2005.
- 386 [20] N. Pelechano, J. M. Allbeck, M. Kapadia, and N. I. Badler. *Simulating heterogeneous crowds with
387 interactive behaviors*. CRC Press, 2016.

- 388 [21] S. Pellegrini, A. Ess, K. Schindler, and L. Van Gool. You'll Never Walk Alone: Modeling Social Behavior
389 for Multi-Target Tracking. In *2009 IEEE 12th International Conference on Computer Vision (CVPR)*,
390 pages 261–268. IEEE, 2009.
- 391 [22] A. Robicquet, A. Sadeghian, A. Alahi, and S. Savarese. Learning Social Etiquette: Human Trajectory
392 Understanding in Crowded Scenes. In *European Conference on Computer Vision (ECCV)*, pages 549–565.
393 Springer, 2016.
- 394 [23] A. Rudenko, L. Palmieri, M. Herman, K. M. Kitani, D. M. Gavrilu, and K. O. Arras. Human motion
395 trajectory prediction: A survey. *The International Journal of Robotics Research*, 39(8):895–935, 2020.
- 396 [24] T. Salzmann, B. Ivanovic, P. Chakravarty, and M. Pavone. Trajectron++: Dynamically-feasible trajectory
397 forecasting with heterogeneous data. In *European Conference on Computer Vision (ECCV)*, pages 683–700.
398 Springer, 2020.
- 399 [25] A. A. Scharine and M. K. McBeath. Right-Handers and Americans Favor Turning to the Right. *Human*
400 *Factors*, 44(2):248–256, 2002.
- 401 [26] A. Seyfried, M. Boltes, J. Kähler, W. Klingsch, A. Portz, T. Rupperecht, A. Schadschneider, B. Steffen,
402 and A. Winkens. Enhanced empirical data for the fundamental diagram and the flow through bottlenecks.
403 *Pedestrian and Evacuation Dynamics 2008*, pages 145–156, 2010.
- 404 [27] S. Singh, M. Kapadia, P. Faloutsos, and G. Reinman. Steerbench: A Benchmark Suite for Evaluating
405 Steering Behaviors. *Computer Animation and Virtual Worlds (CAVW)*, 20(5-6):533–548, 2009.
- 406 [28] K. Sohn, H. Lee, and X. Yan. Learning Structured Output Representation using Deep Conditional
407 Generative Models. In *Neural Information Processing Systems (NIPS)*, 2015.
- 408 [29] S. S. Sohn, H. Zhou, S. Moon, S. Yoon, V. Pavlovic, and M. Kapadia. Laying the Foundations of Deep
409 Long-Term Crowd Flow Prediction. In *European Conference on Computer Vision (ECCV)*, pages 711–728.
410 Springer, 2020.
- 411 [30] J. Van Den Berg, S. J. Guy, M. Lin, and D. Manocha. Reciprocal n-Body Collision Avoidance. In *Robotics*
412 *Research*, pages 3–19. Springer, 2011.
- 413 [31] A. Vemula, K. Muelling, and J. Oh. Social Attention: Modeling Attention in Human Crowds. In *2018*
414 *IEEE International Conference on Robotics and Automation (ICRA)*, pages 4601–4607, 2018.
- 415 [32] S. Wei-Guo, Y. Yan-Fei, W. Bing-Hong, and F. Wei-Cheng. Evacuation behaviors at exit in ca model
416 with force essentials: A comparison with social force model. *Physica A: Statistical Mechanics and its*
417 *Applications*, 371(2):658–666, 2006.
- 418 [33] J. M. Wiener, S. J. Büchner, and C. Hölscher. Taxonomy of human wayfinding tasks: A knowledge-based
419 approach. *Spatial Cognition & Computation*, 9(2):152–165, 2009.
- 420 [34] Z. Yan, T. Duckett, and N. Bellotto. Online Learning for Human Classification in 3D LiDAR-based
421 Tracking. In *In Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and*
422 *Systems (IROS)*, Vancouver, Canada, September 2017.
- 423 [35] W. Zeng, P. Chen, H. Nakamura, and M. Iryo-Asano. Application of social force model to pedestrian
424 behavior analysis at signalized crosswalk. *Transportation research part C: emerging technologies*, 40:143–
425 159, 2014.
- 426 [36] T. Zhao, Y. Xu, M. Monfort, W. Choi, C. Baker, Y. Zhao, Y. Wang, and Y. N. Wu. Multi-Agent Tensor
427 Fusion for Contextual Trajectory Prediction. In *Proceedings of the IEEE/CVF Conference on Computer*
428 *Vision and Pattern Recognition (CVPR)*, pages 12118–12126, June 2019.