

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

Samuel S. Sohn
Rutgers University
samuel.sohn@rutgers.edu

Mihee Lee
Rutgers University
ml1323@rutgers.edu

Seonghyeon Moon
Rutgers University
sm2062@cs.rutgers.edu

Gang Qiao
Rutgers University
gq19@cs.rutgers.edu

Muhammad Usman
York University
usman@cse.yorku.ca

Sejong Yoon
The College of New Jersey
yoons@tcnj.edu

Vladimir Pavlovic
Rutgers University
vladimir@cs.rutgers.edu

Mubbasis Kapadia
Rutgers University
mubbasis.kapadia@rutgers.edu

ABSTRACT

In recent years, human trajectory prediction (HTP) has garnered attention in computer vision literature. Although this task has much in common with the longstanding task of crowd simulation, there is little from crowd simulation that has been borrowed, especially in terms of evaluation protocols. The key difference between the two tasks is that HTP is concerned with forecasting multiple steps at a time and capturing the multimodality of real human trajectories. A majority of HTP models are trained on the same few datasets, which feature small, transient interactions between real people and little to no interaction between people and the environment. Unsurprisingly, when tested on crowd egress scenarios, these models produce erroneous trajectories that accelerate too quickly and collide too frequently, but the metrics used in HTP literature cannot convey these particular issues. To address these challenges, we propose (1) the A2X dataset, which has simulated crowd egress and complex navigation scenarios that compensate for the lack of agent-to-environment interaction in existing real datasets, and (2) evaluation metrics that convey model performance with more reliability and nuance. A subset of these metrics are novel *multiverse metrics*, which are better-suited for multimodal models than existing metrics. The dataset is available at: <https://mubbasis.github.io/HTP-benchmark/>.

CCS CONCEPTS

• **Computing methodologies** → *Motion path planning*; **Model verification and validation**.

KEYWORDS

human trajectory prediction, datasets, evaluation metrics

ACM Reference Format:

Samuel S. Sohn, Mihee Lee, Seonghyeon Moon, Gang Qiao, Muhammad Usman, Sejong Yoon, Vladimir Pavlovic, and Mubbasis Kapadia. 2021. A2X: An Agent and Environment Interaction Benchmark for Multimodal Human Trajectory Prediction. In *MIG'21: Motion, Interaction and Games*. ACM, New York, NY, USA, 9 pages. <https://doi.org/10.1145/1122445.1122456>

1 INTRODUCTION

The study of human navigation has long been of interest to various research communities such as computer graphics [Helbing and Molnar 1995], computer vision [Alahi et al. 2016], cognitive science [Wiener et al. 2009], and robotics [Ferrer et al. 2013]. Advancements in these areas have seen widespread practical application in pandemic response, architectural design, urban planning, transportation engineering, crowd management, socially compliant robot navigation, and entertainment. Accordingly, the influence of human navigation research has reached countless individuals and will continue to do so in the foreseeable future.

Most applications rely on simulation models [Pelechano et al. 2016], which are sufficiently accurate to human behavior and generalizable to unforeseen circumstances. However, the past five years of predictive modeling in computer vision has achieved significantly better accuracy [Rudenko et al. 2020], giving it a strong potential to overtake the longstanding models from computer graphics. This is largely due to the transition from using unimodal, discriminative models [Alahi et al. 2016] that predict a single future trajectory to using multimodal, generative models [Gupta et al. 2018; Mangalam et al. 2020b; Salzmann et al. 2020] that predict a distribution of future trajectories, which captures the inherent uncertainty in human decision-making [Dubey et al. 2019; Scharine and McBeath 2002]. Despite the evolution of models, however, the accuracy metrics that were introduced with the first unimodal models are still in use today. In order to adapt these fundamentally unimodal metrics to multimodal models, the metrics are computed between each predicted trajectory and the ground truth trajectory, and the minimum error for each metric is reported. This results in a gross overestimation of accuracy that we later show is not consistent with the expected accuracy, which may misguide future research efforts. Furthermore, the minimum value is not actionable, because while it is evident that a state-of-the-art (SOTA) multimodal model can find

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

MIG'21, November 10–12, 2021, Lausanne, Switzerland

© 2021 Association for Computing Machinery.

ACM ISBN 978-1-4503-XXXX-X/18/06...\$15.00

<https://doi.org/10.1145/1122445.1122456>

an accurate trajectory, it cannot determine *which* trajectory is most accurate for unseen data. We measure this uncertainty through a decidability metric.

Generalizability cannot be maximized by solely improving upon accuracy metrics. An inaccurate model can be robust by producing realistic trajectories, and an accurate model can fail to be practicable by being undecidable. Models can exist on the continuum between these two extremes, making it critical to consider realism and decidability metrics as well.

Furthermore, there is a stark class imbalance in existing datasets. While datasets are abundant in instances where humans are interacting with each other in open spaces [Alahi et al. 2014; Chavdarova et al. 2018; Kothari et al. 2021; Lerner et al. 2007; Robicquet et al. 2016; Yan et al. 2017], they are significantly lacking in both environment information and instances where humans are interacting with their environment. Ultimately, this hinders generalization at a global level and has led to some models being developed without considering environments at all [Alahi et al. 2016; Gupta et al. 2018].

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

2 BACKGROUND

2.1 Models for Human Trajectory Prediction.

While crowd simulation has been well-studied in computer graphics literature [Kapadia et al. 2015; Thalmann and Musse 2012], we focus our attention on the use of machine learning techniques for the growing field of human trajectory prediction. Earlier methods such as Social LSTM [Alahi et al. 2016] and Social Attention [Vemula et al. 2018] proposed a deterministic model which predict a future trajectory given observed trajectories. However, forecasting trajectories inherently introduces uncertainty in the future, hence the utility of those unimodal models which predict only one future trajectory is limited. Recent studies [Gupta et al. 2018; Ivanovic and Pavone 2019; Mangalam et al. 2020a,b; Salzmann et al. 2020; Zhao et al. 2019] assume the multi-modalities in the future human behavior and predict its distribution to embody the uncertainty. In this paper, we focus on three SOTA methodologies to showcase our benchmark dataset: SocialGAN [Gupta et al. 2018], PECNet [Mangalam et al. 2020b], and Trajectron++ [Salzmann et al. 2020].

SocialGAN [Gupta et al. 2018] adopts GAN [Goodfellow et al. 2014] framework to forecast possible future trajectories and it can

avoid collisions among pedestrians by introducing a pooling mechanism that captures between-human interaction. PECNet [Mangalam et al. 2020b] solves the trajectory prediction problem by first modeling the future goal position distribution using a Variational Autoencoder (VAE) [Kingma and Welling 2014], and then predict the future positions by interpolating the observed positions and the estimated goal position. Trajectron++ [Salzmann et al. 2020] proposes a graph structured recurrent model based on conditional VAE [Sohn et al. 2015] to predict the future trajectories. Further details can be found in the Supplementary Materials.

We investigate these three models as the representatives of the various SOTA works. We choose them because PECNet [Mangalam et al. 2020b] shows an outstanding performance on the long-term trajectory while the short-term trajectory is most well predicted in Trajectron++ [Salzmann et al. 2020]. We expect SocialGAN [Gupta et al. 2018], as one of the earliest and most frequently referred models, to be a bound around existing models with respect to PECNet and Trajectron++. Fig. 1.b shows the coverage comparison of SOTA models in terms of the short- and long-term human trajectory prediction accuracy. We differentiate between predictive models of short-term and long-term trajectories on the basis of goal conditioning. A model that is not goal-conditioned will inherently increase in error as the predicted path length increases, sometimes at an exponential rate [Salzmann et al. 2020], whereas goal-conditioned models are expected to predict long paths without the same trade-off between path length and error.

2.2 Datasets for Human Trajectory Prediction.

The computer vision and graphics community have collected several human pedestrian trajectory datasets. ETH [Pellegrini et al. 2009] and UCY [Lerner et al. 2007] are commonly used datasets that contain five outdoor scenes with jointly more than 1,600 pedestrian trajectories. Stanford Drone Dataset (SDD) [Robicquet et al. 2016] consists of eight outdoor scenes tracking 19,000 targets including pedestrians, bicyclists, skateboarders, cars, and buses collected from a drone. Stanford Crowd Dataset (CFF) [Alahi et al. 2014] consists of pedestrian trajectories collected within a train station building of size $25\text{m} \times 100\text{m}$ for 12×2 hours captured by a distributed camera network. L-CAS 3D Point Cloud People Dataset (LCAS) [Yan et al. 2017] consists of 28,002 scan frames collected within a university building by a 3D LiDAR sensor mounted on a robot that is either stationary or moving. WILDTRACK (WT) [Chavdarova et al. 2018] is a collection of annotated dense pedestrian groups captured by seven static HD cameras in a public square for about 60 minutes. The Supplementary Materials provide more details of these datasets. Some datasets, such as TrajNet++ [Kothari et al. 2021], augment upon existing datasets. TrajNet++ combines ETH/UCY, CFF, LCAS, and Wildtrack datasets, as well as a synthetic dataset generated by ORCA [Van Den Berg et al. 2011].

Existing human trajectory datasets have limitations in the sense of embodying interactions. They either do not contain agent-to-environment (A2E) interactions [Chavdarova et al. 2018], or exhibit limited agent-to-agent (A2A) interactions at small scale in simple environments. We speculate that many self-centered pedestrians are prone to avoid or mitigate, consciously or unconsciously, the influence of the environments and other pedestrians during their

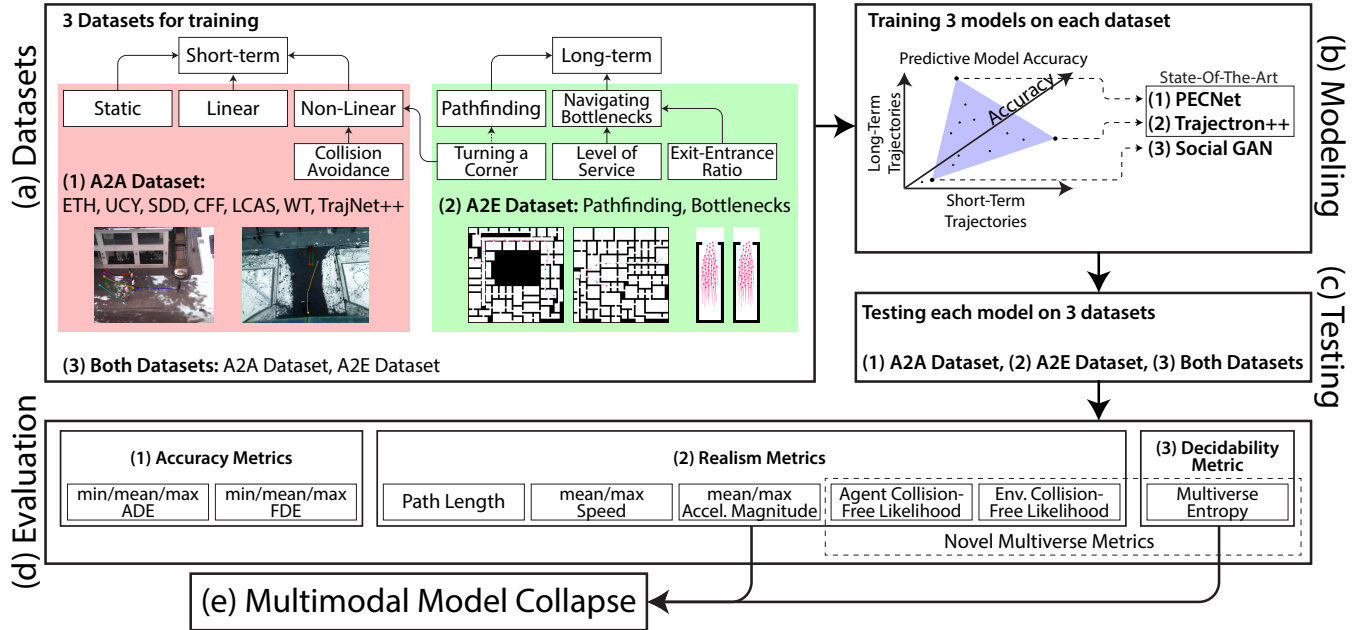


Figure 1: The above framework image shows (a) the differences between the trajectories of existing datasets (A2A) and the novel dataset (A2E), (b-c) the models trained and tested on combinations of A2A and A2E, (d) the proposed set of metrics for evaluating the accuracy, realism, and decidability of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

navigation. In this work, we are proposing datasets that augment A2E and A2A interactions, which may bring benefits for enhancing learning models by encoding more complex trajectory dynamics.

2.3 Evaluation for Human Trajectory Prediction.

In computer graphics community [Singh et al. 2009], trajectories are, in general, measured by motion statistics such as the number of collisions, average speed, average acceleration, and total distance traveled. On the other hand, in machine learning community [Alahi et al. 2016; Gupta et al. 2018; Kothari et al. 2021], the most commonly used evaluation metrics for trajectory forecasting models are Average Displacement Error (ADE) and Final Displacement Error (FDE). ADE is the average L_2 distance between the ground truth and the predicted trajectories across all future steps. FDE is the L_2 distance between the ground truth final destination and the predicted final destination at the end of the future steps. More evaluation metrics in machine learning community are discussed in Supplementary Materials.

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

error. In Section 4 of this work, we propose new metrics that can tackle this issue.

3 AGENT-TO-AGENT AND AGENT-TO-ENVIRONMENT INTERACTION DATASET

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

In order to understand what the shortcomings of current datasets are (Sec. 2), we first taxonomize the characteristics of human trajectories. The TrajNet++ benchmark [Kothari et al. 2021] proposed an initial taxonomy that only considers short-term characteristics, e.g., standing still, moving linearly, or avoiding collisions (Fig. 1.a). While the original taxonomy is sufficient for describing the trajectories in many real datasets and their agent-to-agent (A2A) interactions, models that learn exclusively from these types are insufficient for most applications, which consider environments with obstacles and time frames longer than 5 seconds, which is the practical limit for most models before they become exponentially erroneous [Salzmann et al. 2020]. We have improved upon this by considering long-term characteristics (Fig. 1.a), i.e., pathfinding alone and navigating through crowded bottlenecks. These types of trajectories emerge from agent-to-environment (A2E) interactions, which unfold over a longer time frame than A2A interactions and are essential for navigation within any environment [Sohn et al. 2020].

3.1 Agent-to-Agent Interactions

For representing A2A interactions, we make use of each prior dataset described in Section 2.2: ETH [Lerner et al. 2007], UCY [Lerner et al. 2007], SDD [Robicquet et al. 2016], CFF [Alahi et al. 2014], LCAS [Yan et al. 2017], WT [Chavdarova et al. 2018], and TrajNet++ [Kothari et al. 2021]. These datasets feature transient interactions between agents and little interaction with the environment, which is made difficult to measure by the frequent unavailability of environment information. Therefore, we approximate environment information based on the principle of stigmergy [Helbing et al. 1997; Parunak 2005], which observes the self-organization of human navigation along trails. For each position that agents have traveled through in either the training or testing sets of the ground truth, a 1-meter radius around the position is considered to be navigable. This guarantees that predictions with less than 1 meter of displacement from the ground truth at all times will never intersect with the environment. Additionally, in order to compensate for the imbalance between A2A and A2E interactions in prior datasets, we propose the generation of synthetic data in addition to that of TrajNet++. While real datasets are valuable for their veridicality, there are logistical limitations that prevent the acquisition of real data in OOD scenarios that are unsafe for human participants or prohibitively expensive from an organizational standpoint.

3.2 Agent-to-Environment Interactions

Two such scenarios are used to sample trajectories exhibiting A2E interactions: (1) pathfinding alone in a large, complex environment, which has prohibitive logistical cost and (2) navigating through bottlenecks of varied width with a dense crowd, which can be unsafe. Though simulation models are normally less accurate than predictive models in predicting human trajectories [Alahi et al. 2016], the prevalent Social Force model [Helbing and Molnar 1995] currently outperforms predictive models in terms of robustness, has been used in several application domains [Ferrer et al. 2013; Wei-Guo et al. 2006; Zeng et al. 2014], and has adequate ecological validity in these A2E scenarios, which lack sufficient real data for training predictive models until A2X. On one hand, simulation models are robust enough for producing plausible behavioral data, so all metrics can be used to evaluate A2E-trained predictive models on A2E test cases. On the other hand, simulation models are not perfectly accurate to real human navigation, so A2A-trained models should not be evaluated on A2E test cases using accuracy metrics.

We leverage the Social Force model to simulate 236 scenarios of a single agent navigating between random points in complex $112 \times 112 \text{ m}^2$ environments from [Sohn et al. 2020] (Fig. 2). This produces long-term isolated interactions between single agents and the environment. We then use the same model to simulate well-studied bottleneck scenarios [Haworth et al. 2015; Seyfried et al. 2010] in a $25 \times 7 \text{ m}^2$ room that vary in terms of (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 ratio between the width of the bottleneck and the width of the room (Exit-Entrance Ratio) from $\{0.2, 0.3, 0.4, 0.6, 0.7\}$ (Fig. 2). A total of 398 scenarios have been generated across all combinations of Level of Service and Exit-Entrance Ratio. This produces long-term interactions between agents as a result of the constricting environment. Exact environment information has been provided

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

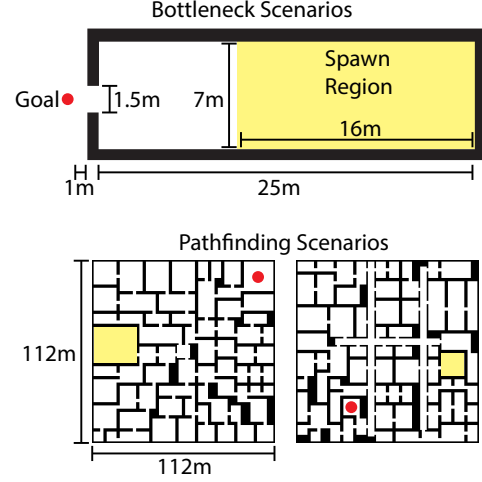


Figure 2: The above images show the exact dimensions of environments from the bottleneck and pathfinding scenarios in A2E.

4 ACCURACY, REALISM, AND DECIDABILITY OF HUMAN TRAJECTORY PREDICTION

We propose a total of 15 accuracy, realism, and decidability metrics (Fig. 1.d). These metrics are either borrowed from computer vision and computer graphics literature [Alahi et al. 2016; Guy et al. 2012; Pellegrini et al. 2009; Singh et al. 2009] or newly developed *multi-verse metrics*, which assess the A2A and A2E interactions of both multimodal models with $k > 1$ and unimodal models with $k = 1$.

4.1 Preliminaries

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

Many of the following metrics make use of aggregate functions. For any d -dimensional vector $\mathbf{v} \in \mathbb{R}^d$, we denote the minimum value by $\Omega(\mathbf{v})$, the mean value by $\Theta(\mathbf{v})$, and the maximum value 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 into a probability distribution $\mathbf{p} \in \mathbb{R}^b$ over a vector of b -many equiangular bins, which radiate from the origin (Fig. 3) and can optionally be divided along the radial dimension according to a maximum vector magnitude. Finally, we denote the L_2 norm by $\|\cdot\|$.

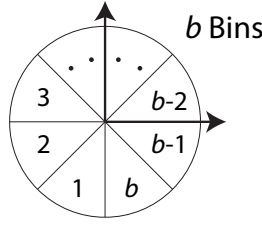


Figure 3: This images shows how $b = 8$ bins would be arranged in 2D space.

4.2 Accuracy Metrics: Comparison to Ground Truth

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

Average Displacement Error (ADE). ADE is computed for each prediction j as \mathbf{a}_j , the average distance between a position in the ground truth and a position in the prediction across t_f frames (Eq. 1) [Pellegrini et al. 2009]. It is then aggregated across the k predictions in three ways: minimum, mean, and maximum, which offers a more reliable expectation of a model’s accuracy than the minimum alone.

Final Displacement Error (FDE). FDE is computed for each prediction j as \mathbf{b}_j , the distance between the final positions of the ground truth and the prediction (Eq. 2) [Alahi et al. 2016]. It is aggregated across the k predictions in the same ways as ADE for better reliability.

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

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

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

$$\text{s.t. } \mathbf{b}_j = \|\mathbf{Y}_{a,0,t_f-1} - \hat{\mathbf{Y}}_{a,j,t_f-1}\|, \quad 0 \leq j < k \quad (2)$$

4.3 Realism Metrics: Motion and Interaction Statistics

Realism metrics are used to describe the movement and interactions within the ground truth and the predictions separately. These metrics can then be used to uncover more nuanced differences between the ground truth and predictions. While they cannot ensure that predictions are accurate, they can ensure that predictions are realistic in their movement and plausible. Every realism metric is computed in the same way for both the ground truth and predictions, so \mathbf{Y} is interchangeable with $\hat{\mathbf{Y}}$ and \mathbf{V} with $\hat{\mathbf{V}}$. For generality, we consider the ground truth as a unimodal model with $k = 1$, but we refer to it as having k paths instead of predictions.

The following motion statistics are used to describe the movement of agent a in either the ground truth or averaged across the k predictions. They have been used to evaluate crowd simulations

in computer graphics research [Singh et al. 2009], but have not yet been used to evaluate predictive models in computer vision.

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

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

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

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

$$\text{s.t. } \mathbf{A}_{j,t} = \|(\mathbf{V}_{a,j,t+1} - \mathbf{V}_{a,j,t}) / \Delta t\|, \quad 0 \leq t < t_f - 2 \quad (5)$$

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

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 computed for each agent a . Next, two values are reported for speed: the mean speed averaged across k paths and the maximum speed averaged across k paths. For each path j of agent a , the mean and maximum speed are computed across $t_f - 1$ frames (Eq. 4).

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

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

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

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

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

$$\mathbf{D}_j = \frac{1}{t_f - 1} \left(\sum_{t=1}^{t_f-1} \mathbf{Y}_{a,j,t} \right) - \mathbf{Y}_{a,j,0}, \quad 0 \leq j < k \quad (8)$$

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

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

4.4 Decidability Metric: Certainty in Movement Direction

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

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

Test	Model	Train	Accuracy Metrics		Realism Metrics						Decidab.
			ADE ↓	FDE ↓	Length	Speed	Accel.	ACFL	ECFL	%Diff. ↓	MVE ↓
			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. Models where minimum accuracy metrics disagree with the averages are red.

4.5 Comparing Realism Metrics

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

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

5 RESULTS

In order to understand the limits of not only the SOTA but also the models that paved the way towards the SOTA, we evaluate three critical multimodal models that are capable of either short-term or long-term trajectory prediction and provide a large coverage over the performance of prior models (Fig. 1.b). In particular, we have selected (1) Social GAN (SGAN) [Gupta et al. 2018], one of the earliest models; (2) Trajectron++ (T++) [Salzmann et al. 2020], a SOTA model for short-term trajectory prediction; and (3) PECNet (PECN) [Mangalam et al. 2020b], a SOTA model for long-term trajectory prediction.

5.1 Training Protocol

Each of the three models was trained on 3 combinations from the A2X Dataset: A2A interaction, A2E interaction, and both (Fig. 1.b), producing a total of 9 models. We denote that either a model has been trained on a particular combination using a subscript, e.g., $SGAN_{Both}$. Each trained model was then evaluated on the testing sets of the 3 combinations (Fig. 1.c). The results of the evaluations on A2A and A2E are reported in Table 1, while the results on both A2A and A2E combined and corresponding visualizations are reported in the Supplementary Materials. According to the dataset, the following parameters have been set for the evaluation: $k = 20$, $t_p = 8$, $t_f = 12$, and $\Delta t = 0.4$, meaning that each agent is receiving 3.2 seconds of input data and predicting 4.8 seconds into the future.

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

5.2 Analysis

5.2.1 Training on both types of interaction consistently has near-best accuracy. As expected, we find that in terms of all accuracy metrics, models trained on a single type of interaction perform very poorly on test scenarios that feature the other type of interaction.

By training any of the three models (SGAN, PECN, or T++) on both types of interactions, we find that the accuracy is consistently near-best among all three training datasets by a small margin. For testing on A2A, a model trained on both types is closer in accuracy to the same model trained on A2A, and for testing on A2E, it is closer to the same model trained on A2E. In fact, when testing on A2A, training SGAN and T++ on both types achieves the best mean/maximum ADE and mean FDE among all training datasets. This makes training on both types of interactions an excellent compromise for balancing accuracy between real-world cases from A2A and critical synthetic cases from A2E.

5.2.2 Existing evaluation metrics can misjudge model accuracy. When testing on A2A, $SGAN_{A2A}$ and $T++_{A2A}$ are misjudged as being better than $SGAN_{Both}$ and $T++_{Both}$ according to minimum ADE and minimum FDE (highlighted in red). Reliance on these overly optimistic existing metrics will lead to choosing models that are less accurate than others on average.

5.2.3 Realism metrics influence model choice based on the use case. We cannot rely only on the accuracy of models to determine which is best, since anything short of perfect accuracy carries risk. The realism metrics allow us to better understand a model's performance in the context of its application. For example, we find that the maximum speed and acceleration for $T++_{Both}$ are significantly higher than the ground truth, which for an application in socially compliant robot navigation can discomfort or potentially harm surrounding humans [Kruse et al. 2013]. In contrast, $SGAN_{Both}$ has lower average accuracy by a small margin, but it boasts higher realism by a large margin in terms of maximum speed, maximum acceleration magnitude, and ACFL. We attribute $SGAN_{Both}$'s higher ACFL to the tighter spread of its predictions than $T++_{Both}$ according to MVE. Ultimately, the choice of a model depends on the application, but without the joint consideration of the proposed accuracy and realism metrics, a practitioner may be led to choose an unsuitable model.

5.2.4 A2E is essential for learning collision avoidance. Models trained exclusively on A2E interactions tend to have lower likelihoods of A2A collision (higher ACFL) than models trained on A2A interactions alone or on both types of interactions. This highlights the importance of A2E for improving robustness even in real-world scenarios such as A2A.

5.2.5 ECFL indicates that A2A scenarios have trivial A2E interactions. Models trained on A2E achieve the lowest likelihood of A2E collision (highest ECFL) when testing on A2E, but still have some room to improve. In contrast, we find that ECFL is nearly perfect for A2A scenarios, indicating that A2A scenarios do not challenge models with A2E interactions.

5.2.6 Multimodal models can be decidable. Although PECN is a multimodal model, it has a near-zero MVE, which is significantly lower than SGAN and T++. This indicates that PECN has certainty in the direction that agents will travel along (regardless of whether the direction is correct). PECN also achieves the highest ACFL owing to its low MVE, which is low enough to consider PECN as being decidable and likely helps it in performing long-term trajectory prediction.

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							
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
	SGAN	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
	SGAN	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 worse 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 $k = 1$.

5.3 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 the single most socially compliant 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 . This prediction is the closest in realism to prior ground truth for the same type of scenario (Eq. 9). Table 2 shows the result of applying this technique to all 9 models. Across all models, 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. Furthermore, the social compliance of models is drastically improved through MMC, making them less likely to produce collisions with other agents.

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 evaluation. 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 critical 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 much room for improvement even among the SOTA models.

7 ACKNOWLEDGEMENTS

The research was supported in part by NSF awards: IIS-1703883, IIS-1955404, IIS-1955365, RETTL-2119265, and EAGER-2122119. The authors acknowledge use of the TCNJ ELSA HPC cluster, funded by NSF grant OAC-1828163, for conducting the research reported in this paper.

REFERENCES

- Alexandre Alahi, Kratharth Goel, Vignesh Ramanathan, Alexandre Robicquet, Li Fei-Fei, and Silvio Savarese. 2016. Social LSTM: Human Trajectory Prediction in Crowded Spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 961–971.
- Alexandre Alahi, Vignesh Ramanathan, and Li Fei-Fei. 2014. Socially-Aware Large-Scale Crowd Forecasting. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2203–2210.
- Tatjana Chavdarova, Pierre Baqué, Stéphane Bouquet, Andrii Maksai, Cijo Jose, Timur Bagautdinov, Louis Lettry, Pascal Fua, Luc Van Gool, and François Fleuret. 2018. Wildtrack: A Multi-Camera HD Dataset for Dense Unscripted Pedestrian Detection. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 5030–5039.
- Rohit K Dubey, Samuel S Sohn, Christoph Hoelscher, and Mubbasir Kapadia. 2019. Fusion-Based Wayfinding Prediction Model for Multiple Information Sources. In *2019 22th International Conference on Information Fusion (FUSION)*. IEEE, 1–8.
- Gonzalo Ferrer, Anais Garrell, and Alberto Sanfeliu. 2013. Social-aware robot navigation in urban environments. In *2013 European Conference on Mobile Robots*. IEEE, 331–336.
- Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. 2014. Generative Adversarial Nets. In *Proceedings of the 27th International Conference on Neural Information Processing Systems (NIPS)* (Montreal, Canada). 2672–2680.
- Agrim Gupta, Justin Johnson, Li Fei-Fei, Silvio Savarese, and Alexandre Alahi. 2018. Social GAN: Socially Acceptable Trajectories with Generative Adversarial Networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. 2255–2264.
- Stephen J Guy, Jur Van Den Berg, Wenxi Liu, Rynson Lau, Ming C Lin, and Dinesh Manocha. 2012. A Statistical Similarity Measure for Aggregate Crowd Dynamics. *ACM Transactions on Graphics (TOG)* 31, 6 (2012), 1–11.
- Brandon Haworth, Muhammad Usman, Glen Berseth, Mubbasir Kapadia, and Petros Faloutsos. 2015. Evaluating and optimizing level of service for crowd evacuations. In *Proceedings of the 8th ACM SIGGRAPH Conference on Motion in Games*. 91–96.
- Dirk Helbing and Peter Molnar. 1995. Social Force Model for Pedestrian Dynamics. *Physical review E* 51, 5 (1995), 4282.
- Dirk Helbing, Frank Schweitzer, Joachim Keltsch, and Peter Molnar. 1997. Active walker model for the formation of human and animal trail systems. *Physical review E* 56, 3 (1997), 2527.
- Boris Ivanovic and Marco Pavone. 2019. The Trajectron: Probabilistic Multi-Agent Trajectory Modeling With Dynamic Spatiotemporal Graphs. In *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. 2375–2384. <https://doi.org/10.1109/ICCV.2019.00246>
- Mubbasir Kapadia, Nuria Pelechano, Jan Allbeck, and Norm Badler. 2015. Virtual crowds: Steps toward behavioral realism. *Synthesis lectures on visual computing: computer graphics, animation, computational photography, and imaging* 7, 4 (2015), 1–270.
- Diederik P. Kingma and Max Welling. 2014. Auto-Encoding Variational Bayes. In *2nd International Conference on Learning Representations (ICLR), Banff, AB, Canada, April 14–16, 2014, Conference Track Proceedings*.
- Parth Kothari, Sven Kreiss, and Alexandre Alahi. 2021. Human Trajectory Forecasting in Crowds: A Deep Learning Perspective. *IEEE Transactions on Intelligent Transportation Systems* (2021). doi: 10.1109/ITITS.2021.3069362.
- Thibault Kruse, Amit Kumar Pandey, Rachid Alami, and Alexandra Kirsch. 2013. Human-aware robot navigation: A survey. *Robotics and Autonomous Systems* 61, 12 (2013), 1726–1743.
- Alon Lerner, Yiorgos Chrysanthou, and Dani Lischinski. 2007. Crowds by Example. In *Computer Graphics Forum*, Vol. 26. Wiley Online Library, 655–664.
- Karttikeya Mangalam, Yang An, Harshayu Girase, and Jitendra Malik. 2020a. From Goals, Waypoints Paths To Long Term Human Trajectory Forecasting. *arXiv preprint arXiv:2012.01526* (2020).
- Karttikeya Mangalam, Harshayu Girase, Shreyas Agarwal, Kuan-Hui Lee, Ehsan Adeli, Jitendra Malik, and Adrien Gaidon. 2020b. It is Not the Journey but the Destination: Endpoint Conditioned Trajectory Prediction. *arXiv preprint arXiv:2004.02025* (2020).
- H Van Dyke Parunak. 2005. A survey of environments and mechanisms for human-human stigmergy. In *International workshop on environments for multi-agent systems*. Springer, 163–186.
- Nuria Pelechano, Jan M Allbeck, Mubbasir Kapadia, and Norman I Badler. 2016. *Simulating heterogeneous crowds with interactive behaviors*. CRC Press.
- Stefano Pellegrini, Andreas Ess, Konrad Schindler, and Luc Van Gool. 2009. You'll Never Walk Alone: Modeling Social Behavior for Multi-Target Tracking. In *2009 IEEE 12th International Conference on Computer Vision (CVPR)*. IEEE, 261–268.
- Alexandre Robicquet, Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. 2016. Learning Social Etiquette: Human Trajectory Understanding in Crowded Scenes. In *European Conference on Computer Vision (ECCV)*. Springer, 549–565.
- Andrey Rudenko, Luigi Palmieri, Michael Herman, Kris M Kitani, Dariu M Gavrilu, and Kai O Arras. 2020. Human motion trajectory prediction: A survey. *The International Journal of Robotics Research* 39, 8 (2020), 895–935.
- Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. 2020. Trajectron++: Dynamically-Feasible Trajectory Forecasting with Heterogeneous Data. In *European Conference on Computer Vision (ECCV)*. Springer, 683–700.
- Angelique A Scharine and Michael K McBeath. 2002. Right-Handers and Americans Favor Turning to the Right. *Human Factors* 44, 2 (2002), 248–256.
- Armin Seyfried, Maik Boltes, Jens Kähler, Wolfram Klingsch, Andrea Portz, Tobias Rupprecht, Andreas Schadschneider, Bernhard Steffen, and Andreas Winkens. 2010. Enhanced empirical data for the fundamental diagram and the flow through bottlenecks. *Pedestrian and Evacuation Dynamics 2008* (2010), 145–156.
- Shawn Singh, Mubbasir Kapadia, Petros Faloutsos, and Glenn Reinman. 2009. Steer-bench: A Benchmark Suite for Evaluating Steering Behaviors. *Computer Animation and Virtual Worlds (CAVW)* 20, 5–6 (2009), 533–548.
- Kihyuk Sohn, Honglak Lee, and Xinchen Yan. 2015. Learning Structured Output Representation using Deep Conditional Generative Models. In *Neural Information Processing Systems (NIPS)*.
- Samuel S Sohn, Honglu Zhou, Seonghyeon Moon, Sejong Yoon, Vladimir Pavlovic, and Mubbasir Kapadia. 2020. Laying the Foundations of Deep Long-Term Crowd Flow Prediction. In *European Conference on Computer Vision (ECCV)*. Springer, 711–728.
- Daniel Thalmann and Soraia Raupp Musse. 2012. *Crowd simulation*. Springer Science & Business Media.
- Jur Van Den Berg, Stephen J Guy, Ming Lin, and Dinesh Manocha. 2011. Reciprocal n-Body Collision Avoidance. In *Robotics Research*. Springer, 3–19.
- Anirudh Vemula, Katharina Muelling, and Jean Oh. 2018. Social Attention: Modeling Attention in Human Crowds. In *2018 IEEE International Conference on Robotics and Automation (ICRA)*. 4601–4607. <https://doi.org/10.1109/ICRA.2018.8460504>
- Song Wei-Guo, Yu Yan-Fei, Wang Bing-Hong, and Fan Wei-Cheng. 2006. Evacuation behaviors at exit in CA model with force essentials: A comparison with social force model. *Physica A: Statistical Mechanics and its Applications* 371, 2 (2006), 658–666.
- Jan M Wiener, Simon J Büchner, and Christoph Hölscher. 2009. Taxonomy of human wayfinding tasks: A knowledge-based approach. *Spatial Cognition & Computation* 9, 2 (2009), 152–165.
- Zhi Yan, Tom Duckett, and Nicola Bellotto. 2017. Online Learning for Human Classification in 3D LiDAR-based Tracking. In *Proceedings of the 2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*. Vancouver, Canada.
- Weiliang Zeng, Peng Chen, Hideki Nakamura, and Miho Iryo-Asano. 2014. Application of social force model to pedestrian behavior analysis at signalized crosswalk. *Transportation research part C: emerging technologies* 40 (2014), 143–159.
- Tianyang Zhao, Yifei Xu, Mathew Monfort, Wongun Choi, Chris Baker, Yibiao Zhao, Yizhou Wang, and Ying Nian Wu. 2019. Multi-Agent Tensor Fusion for Contextual Trajectory Prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. 12118–12126.