

# Group-based Motion Prediction for Navigation in Crowded Environments

Anonymous Author(s)

Affiliation

Address

email

1       **Abstract:** We focus on the problem of planning the motion of a robot in a dy-  
2       namic multiagent environment such as a pedestrian scene. Enabling the robot  
3       to navigate safely and in a socially compliant fashion in such scenes requires a  
4       representation that accounts for the unfolding multiagent dynamics. Existing ap-  
5       proaches to this problem tend to employ microscopic models of motion prediction  
6       that reason about the individual behavior of other agents. While such models may  
7       achieve high tracking accuracy in trajectory prediction benchmarks, they often  
8       lack an understanding of the group structures unfolding in crowded scenes. In-  
9       spired by the Gestalt theory from psychology, we build a Model Predictive Control  
10      framework (G-MPC) that leverages group-based prediction for robot motion plan-  
11      ning. We conduct an extensive simulation study involving a series of challenging  
12      navigation tasks in scenes extracted from two real-world pedestrian datasets. We  
13      illustrate that G-MPC enables a robot to achieve statistically significantly higher  
14      safety and lower number of group intrusions than a series of baselines featuring  
15      individual pedestrian motion prediction models. Finally, we show that G-MPC  
16      can handle noisy lidar-scan estimates without significant performance losses.

## 17   1 Introduction

18   Over the past three decades, there has been a vivid interest in the area of robot navigation in pedes-  
19   trian environments [1, 2, 3, 4, 5]. Planning robot motion in such environments can be challenging  
20   due to the lack of rules regulating traffic, the close proximity of agents and the complex emerging  
21   multiagent interactions. Further, accounting for human safety and comfort as well as robot efficiency  
22   add to the complexity of the problem.

23   To address such specifications, a common [6, 4, 3, 7, 8] paradigm involves the integration of a  
24   behavior prediction model into a planning mechanism. Recent models tend to predict the individual  
25   interactions among agents to enable the robot to determine collision-free candidate paths [3, 4,  
26   9]. While this paradigm is well-motivated, it tends to ignore the structure of interaction in such  
27   environments. Often, the motion of pedestrians is coupled as a result of social grouping. Further,  
28   the motion of multiple agents can often be *effectively* grouped as a result of similarity in motion  
29   characteristics. Lacking a mechanism for understanding the emergence of this structure, the robot  
30   motion generation mechanism may yield unsafe or uncomfortable paths for human bystanders, often  
31   violating the space of social groups.

32   Motivated by such observations, we draw inspiration from human navigation to propose the use of  
33   group-based prediction for planning in crowd navigation domains. We argue that humans do not  
34   employ detailed individual trajectory prediction mechanisms. In fact, our motion prediction capa-  
35   bilities are short-term and do not scale with the number of agents. We do however employ effective  
36   grouping techniques that enable us to discover safe and efficient paths among motions of crowd  
37   networks. This anecdotal observation is aligned with gestalt theory from psychology [10] which  
38   suggests that organisms tend to perceive and process *formations of entities*, rather than individual  
39   components. Such techniques have recently led to advances in computer vision [11] and computa-  
40   tional photography [12]. Similarly, we envision that a robot could reason the formation of effective  
41   groups in a crowded environment and react to their motion as an effective way to navigate safely.

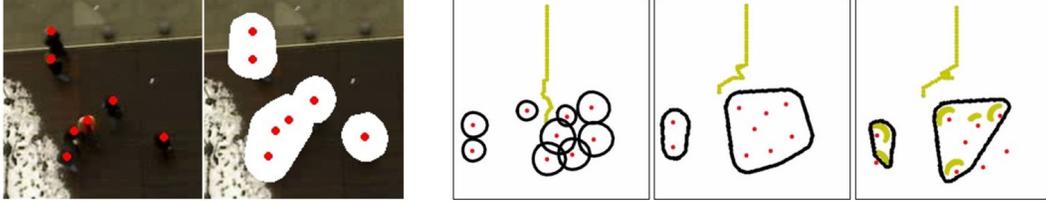


Figure 1: Based on a representation of social grouping [13], we build a group behavior prediction model to empower a robot to perform safe and socially compliant navigation in crowded spaces. The images to the left demonstrate an example of our representation overlaid on top of a scene from a real-world dataset [14]. The images to the right demonstrates that a model predictive controller equipped with our prediction model is able to navigate around the group socially (middle) as opposed to the baseline that cuts through the group (left). Our formulation is also able to handle imperfect state estimates (right) where the yellow arcs are scan points from a simulated 2D lidar laser scan.

42 In this paper, we propose a group-based representation coupled with an autoencoder prediction  
 43 model based on the group-space approximation model of Wang and Steinfeld [13]. This model  
 44 groups a crowd into sets of agents with similar motion characteristics and draws geometric en-  
 45 closures around them, given observation of their states. The prediction module then predicts future  
 46 states of these enclosures. We conduct an extensive empirical evaluation over 5 different human  
 47 datasets [14, 15], each with a flow following and a crossing scenario. We further conduct a same  
 48 set of evaluations with agents powered by ORCA [16] that share the start and end locations in the  
 49 datasets. Last but not least, we conducted evaluation given inputs in the form of simulated laser  
 50 scans, from which pedestrians are only partially observable or even completely occluded. We com-  
 51 pare the performance of our group-based formulation against three individual reasoning baselines:  
 52 a) a reactive baseline with no prediction; b) a constant velocity prediction baseline; c) one based on  
 53 individual S-GAN trajectory predictions [17]. We present statistically significant evidence suggest-  
 54 ing that agents powered by our formulation produce safer and more socially compliant behavior and  
 55 are potentially able to handle imperfect state estimates.

## 56 2 Related Work

57 Over the recent years, a considerable amount of research has been placed to the problem of robot  
 58 navigation in crowded pedestrian environments [4, 3, 18, 8, 19]. Such environments often comprise  
 59 groups of pedestrians, navigating as coherent entities. This has motivated recent work on group  
 60 detection and group motion modeling.

61 Groups are often perceived as sub-modular entities that collectively define the behavior of the crowd.  
 62 Šochman and Hogg [20] suggests that 50-70% of pedestrians walk in groups. Many works exist in  
 63 group detection. One popular area in such domain is static group detection, often leveraging F-  
 64 ormation theories [21]. However, dynamic groups often dominate pedestrian rich environments  
 65 and they exhibit different spatial behavior [22]. Among dynamic group detection, the most common  
 66 approach treats grouping as a probabilistic process where groups are a reflection of close proba-  
 67 bilistic association of pedestrian trajectories [23, 24, 25, 26, 27]. Others use graph models to build  
 68 inter-pedestrian relationships with strong graphical connections indicating groups [28, 29]. The so-  
 69 cial force model [30] also inspires Mazzon et al. [31], Šochman and Hogg [20] to develop features  
 70 that indicate groups. Clustering is another common group of technique to group pedestrians with  
 71 similar features into groups [32, 33, 34, 35]. In this paper, we do not intend to explore the state-of-  
 72 the-art grouping practice. For our formulation, it is sufficient to employ a simple clustering based  
 73 grouping method proposed by Chatterjee and Steinfeld [35].

74 Applications on groups often focus on a specific behavior aspect. In terms of interacting with pedes-  
 75 trians, a major focus in this area is how a robot should behave as part of the group formation [36].  
 76 On dyad groups involving a single human and a robot, some researchers examined socially appro-  
 77 priate following behavior [37, 38, 39, 40] and conversely, guiding behavior [41, 42, 43]. In works  
 78 that do not include robots as part of the pedestrian groups, some research teams studied how a robot  
 79 should guide a group of pedestrians [44, 45, 46]. From navigation perspective, Yang and Peters [22]  
 80 leverage groups as obstacles, but their group space is a collection of individual personal spaces with  
 81 occasional O-space modeling from F-formation theories. Without the engineered occasional occur-

82 rence of O-space, their representation reduces to one of our baselines. Katyal et al. [47] introduce an  
 83 additional cost term that leverages robot’s distance to the closest group in a reinforcement learning  
 84 framework. They model groups using convex hulls directly generated from pedestrian coordinates  
 85 instead of taking personal spaces into consideration. This less principled approach often leads to the  
 86 robot approaching dangerously close to pedestrians. In our work, we additionally explore the capa-  
 87 bilities of groups in handling imperfect sensor inputs. While our focus is on analysing the benefits of  
 88 groups, our group based formulation can be easily incorporated into Katyal et al. [47]’s framework.

### 89 3 Problem Statement

90 Consider a robot navigating in a workspace  $\mathcal{W} \subseteq \mathbb{R}^2$  amongst  $n$  other dynamic agents. Denote by  
 91  $s \in \mathcal{W}$  the state of the robot and by  $s^i \in \mathcal{W}$  the state of agent  $i \in \mathcal{N} = \{1, \dots, n\}$ . The robot is  
 92 navigating from a state  $s_0$  towards a destination  $s_T$  by executing a policy  $\pi : \mathcal{W}^{n+1} \times \mathcal{U} \rightarrow \mathcal{U}$  that  
 93 maps the assumed fully observable world state  $\mathbf{S} = s \cup_{i=1:n} s^i$  to a control action  $u \in \mathcal{U}$ , drawn  
 94 from a space of controls  $\mathcal{U} \subseteq \mathbb{R}^2$ . We assume that the robot is not aware of agents’ destinations  $s_T^i$   
 95 or policies  $\pi_i : \mathcal{W}^{n+1} \times \mathcal{U}^i \rightarrow \mathcal{U}^i, i \in \mathcal{N}$ . In this paper, our goal is to design a policy  $\pi$  that enables  
 96 the robot to navigate from  $s_0$  to  $s_T$  safely and socially.

### 97 4 Group-based Prediction

98 We introduce a framework for group-based representations based on [13] and a model for group-  
 99 based prediction that is amenable for use in decentralized multiagent navigation.

#### 100 4.1 Group Representation

101 Define as  $\theta^i \in [0, 2\pi)$  the orientation of agent  $i \in \mathcal{N}$  which is assumed to be aligned with the  
 102 direction of its velocity  $u^i$ , extracted via finite differencing of its position over a timestep  $dt$  and  
 103 denote by  $v^i = \|u^i\| \in \mathbb{R}^+$  its speed. We define an augmented state for agent  $i$  as  $q^i = (s^i, \theta^i, v^i)$ .

104 We treat a social group as a set of agents who are in close proximity and share similar motion  
 105 characteristics. Assume that a set of  $J$  groups,  $\mathcal{J} = \{1, \dots, J\}$  navigate in a scene. Define by  
 106  $g^i \in \mathcal{J}$  a variable indicating the group membership of agent  $i$ . We then define a group  $j \in \mathcal{J}$  as a  
 107 set  $G^j = \{i \in \mathcal{N} \mid g^i = j\}$  and collect the set of all groups in a scene into a set  $\mathbf{G} = \{G^j \mid j \in \mathcal{J}\}$ .

108 **Extracting Group Membership.** We define the combined augmented state of all agents as  $\mathbf{q} =$   
 109  $\cup_{i=1:n} q^i$ . To obtain group memberships for a set of agents  $\mathcal{N}$ , we apply the Density-Based Spatial  
 110 Clustering of Applications with Noise algorithm (DBSCAN) [48] on agent states:

$$\mathbf{G} \leftarrow \text{DBSCAN}(\mathbf{q} \mid \epsilon_s, \epsilon_\theta, \epsilon_v) \quad (1)$$

111 Where  $\epsilon_s, \epsilon_\theta, \epsilon_v$  are respectively threshold values on agent distances, orientation and speeds for the  
 112 clustering method.

113 **Extracting the Social Group Space.** For each group  $G^j, j \in \mathcal{J}$ , we define a *social group space* as  
 114 a geometric enclosure  $\mathcal{G}^j$  around agents of the group. For each agent  $i \in G^j$ , we define a personal  
 115 space  $\mathcal{P}^i$  as a two-dimensional asymmetric Gaussian based on the model introduced by Kirby [49].  
 116 Refer to Appendix A for detailed descriptions.

117 Given the personal spaces  $\mathcal{P}^i, i \in G^j$ , of all agents in a group  $j$ , we extract the social group space  
 118 of the whole group as a convex hull:

$$\mathcal{G}^j = \text{Convexhull}(\{\mathcal{P}^i \mid i \in G^j\}). \quad (2)$$

119 The shape described by  $\mathcal{G}^j$  represents an obstacle space representation of a group containing agents  
 120 in close proximity with similar motion characteristics. For convenience, let us collect the spaces of  
 121 all groups in a scene into a set  $\mathcal{G} = \{\mathcal{G}^j \mid j \in \mathcal{J}\}$ .

#### 122 4.2 Group Space Prediction Oracle

123 Based on the group-space representation of Sec. 4.1, we describe a prediction oracle that outputs an  
 124 estimate of the future spaces occupied by a set of groups  $\mathcal{G}_{t:t_f}$  up to a time  $t_f = t + f$ , where  $f$

Table 1: Autoencoder Performance

	Metric	ETH	HOTEL	ZARA1	ZARA2	UNIV
Baseline	mIoU (%)	83.52	90.37	88.04	89.30	85.32
	fIoU (%)	76.32	85.38	82.14	83.88	77.24
Autoencoder	mIoU (%)	86.66	92.10	89.97	90.94	87.52
	fIoU (%)	78.64	86.83	83.77	85.09	78.55

125 is a future horizon given a past sequence of group spaces  $\mathcal{G}_{t_h:t}$  from time  $t_h = t - h$  where  $h$  is a  
 126 window of past observations:

$$\mathcal{G}_{t:t_f} \leftarrow \mathcal{O}(\mathcal{G}_{t_h:t}) = \cup_{j=1:J} \mathcal{O}_j(\mathcal{G}_{t_h:t}^j), \quad (3)$$

127 where  $\mathcal{O}_j$  is a model generating a group space prediction for group  $G^j$ . Refer to Appendix B for  
 128 detailed description of partial input handling.

129 We implement the oracle  $\mathcal{O}_j$  of eq. (3) using a simple autoencoder. The encoder follows the 3D  
 130 convolutional architecture in [50] whereas the decoder mirrors the model layout of the encoder. The  
 131 autoencoder takes as input a sequence<sup>1</sup>  $\mathcal{G}_{t_h:t}$  and outputs a sequence  $\mathcal{G}_{t+1:t_f}$  which we pass through  
 132 a sigmoid layer. We supervise the autoencoder’s output using the binary cross entropy loss.

133 We verified the effectiveness of our autoencoder on the 5 scenes of our experiments by conducting  
 134 a cross-validation comparison against a baseline. The baseline predicts the future shapes by linearly  
 135 translating the last social group shape using its geometric center velocity. We use Intersection over  
 136 Union (IoU) as our metric. Between the ground truths and the predictions, this metric divides the  
 137 number of overlapped pixels by the number of pixels occupied by either one of them. As shown in  
 138 Table 1, our autoencoder outperforms the baseline.

## 139 5 Model Predictive Control with Group-based Prediction

140 We describe G-MPC, a model predictive control (MPC) framework for navigation in multiagent  
 141 environments that leverages the group-based prediction oracle of Sec. 4.

142 We describe our group-prediction informed MPC, or G-MPC. At planning time  $t$ , given a (possibly  
 143 partial) augmented world state history  $\mathcal{Q}_{t_h:t}$ , we first extract a sequence of group spaces  $\mathcal{G}_{t_h:t}$  based  
 144 on the method of Sec. 4.1. Given these, the robot computes an optimal control trajectory  $\mathbf{u}^* = u_{1:K}^*$   
 145 of length  $K$  by solving the following optimization problem:

$$(\mathbf{s}^*, \mathbf{u}^*) = \arg \min_{u_{1:K}} \sum_{k=1:K} \gamma^k J(s_{k+1}, \mathcal{G}_{k+1}, s_T) \quad (4)$$

$$s.t. \mathcal{G}_{2-h:1} \leftarrow \mathcal{G}_{t_h:t} \quad (5)$$

$$s_1 \leftarrow s_t \quad (6)$$

$$\mathcal{G}_{k+1:k_f} = \mathcal{O}(\mathcal{G}_{k_h:k}) \quad (7)$$

$$u_k \in \mathcal{U} \quad (8)$$

$$s_{k+1} = s_k + u_k \cdot dt, \quad (9)$$

146 where  $\gamma$  is the discount factor and  $J$  represents a cost function, eq. (5) initializes the group space  
 147 history ( $k = 2 - h$  is the timestep displaced a horizon  $h$  in the past from the first MPC-internal  
 148 timestep  $k = 1$ ), eq. (6) initializes the robot state to the current robot state  $s_t$ , eq. (7) is an update  
 149 rule recursively generating a predicted future group sequence up to timestep  $k_f = k + f$  given  
 150 history from time  $k_h = k - h$  up to time  $k$ ,  $\mathcal{O}$  represents a group-space prediction oracle based on  
 151 Sec. 4, and eq. (9) is the robot state transition assuming a fixed time parametrization of step size  $dt$ .

152 We employ a weighted sum of costs  $J_g$  and  $J_d$ , penalizing respectively distance to the robot’s goal  
 153 and proximity to groups:

$$J(s_k, \mathcal{G}_k, s_T) = \lambda J_g(s_k, s_T) + (1 - \lambda) J_d(s_k, \mathcal{G}_k), \quad (10)$$

<sup>1</sup>The oracle input sequence is first converted into image-space coordinates using the homography matrix of the scene. We also preprocess inputs to have normalized scale and group positions. The autoencoder output is converted back into Cartesian coordinates using the inverse homography transform.

Table 2: Number of trials per task and scene.

Task	ETH	HOTEL	ZARA1	ZARA2	UNIV
Flow	58	43	25	127	106
Cross	58	44	28	129	114

154 where  $\lambda$  is a weight representing the balance between the two costs and

$$J_g(s_k) = \begin{cases} 0, & \text{if } s_k \in \mathcal{G}_k \\ \lambda \|s_{k-1} - s_T\|, & \text{else,} \end{cases} \quad (11)$$

155 penalizes a rollout according to the distance of the last collision-free waypoint to the robot’s goal.  
156 Further, we define  $J_d$  as:

$$J_d(s_k, \mathcal{G}_k) = \exp(-\mathcal{D}(s_k, \mathcal{G}_k)), \quad (12)$$

157 where

$$\mathcal{D}(s_k, \mathcal{G}_k) = \begin{cases} \min_{j \in \mathcal{J}} D(s_k - \mathcal{G}_k^j), & \text{if } s_k \notin \mathcal{G}_k^j \\ -\min_{j \in \mathcal{J}} D(s_k - \mathcal{G}_k^j), & \text{else,} \end{cases} \quad (13)$$

158 where  $D(s_k - \mathcal{G}_k^j)$  returns the minimum distance between the robot state and the space occupied by  
159 group  $j$  at time  $k$ . Using  $D$ , function  $\mathcal{D}$  computes the minimum distance to any group for a given  
160 time. In most cases, the robot lies outside of groups, i.e.,  $s_k \notin \mathcal{G}_k^j$ —therefore, the cost  $J_d$  tries to  
161 maximize the distance  $\mathcal{D}$ . Sometimes, the robot might end up entering the group space  $\mathcal{G}$ —in those  
162 cases,  $J_d$  tries to minimize  $\mathcal{D}$ , to steer the robot towards the direction of quickest escape from the  
163 group. In case that the robot is inside a group to begin with, we shrink the group sizes in Sec. 4.1  
164 until the robot is outside the groups again.

165 To solve eq. (4), we search over a finite set  $\mathcal{U}$  of control trajectories of horizon  $K$ . With the assump-  
166 tion that the robot is holonomic and is not under any kinematic constraints, we use a set of  $R$  control  
167 rollouts  $\mathcal{U} = \{u^1, \dots, u^R\}$  with three levels of tangential speeds and a set of turning speed, i.e.,

$$u_{1:K}^r = (v \cos \psi, v \sin \psi, \omega), \psi = \frac{2\pi r}{R}, v \in \left\{ \frac{1}{3}v_{max}, \frac{2}{3}v_{max}, v_{max} \right\}, \omega \in \left\{ 0, \pm \frac{\pi}{2} \right\} \quad (14)$$

168 To ensure compatibility between our group-based prediction model and our MPC formulation, we  
169 set the control rollout time horizon to be the prediction model’s prediction horizon, or  $K = f$ .

## 170 6 Evaluation

171 We evaluate our framework through a simulation study in which the robot performs a navigation  
172 task (a transition between two points) within a crowds of dynamic agents in a set of scenes.

### 173 6.1 Experimental Setup

174 We consider a set of realistic pedestrian scenes, drawn from the ETH [14] (ETH and HOTEL scenes)  
175 and UCY [15] (ZARA1, ZARA2 and UNIVERSITY scenes) datasets, which often serve as bench-  
176 marking testbeds in the motion prediction and social navigation literature [51, 17, 52, 53]. In each  
177 scene, we define two navigation tasks (see Fig. 2): *Flow*: in which the robot navigates along the  
178 crowd flow and *Cross* in which the robot intersects vertically with the traffic flow. For each task, we  
179 generate a set of trials by segmenting the scene recording into blocks involving challenging interac-  
180 tions. We define a challenging interaction to be a segment involving at least 5 pedestrians inside the  
181 test region drawn in black in Fig. 2. This process provided us with a distribution of trials as shown in  
182 table Table 2. Across all trials, we keep the robot’s maximum at 1.75m/s and use a fixed timestep  
183 size  $dt = 0.1$ .

184 We consider two experimental conditions: an *Offline* and an *Online* one. In the *Offline* one, the  
185 robot navigates among a crowd moving according to a recording of a human crowd. Under this  
186 condition, pedestrians act as dynamic obstacles that do not react to the robot, a situation which  
187 could arise in cases where robots are of shorter size and could thus be easily missed by navigating

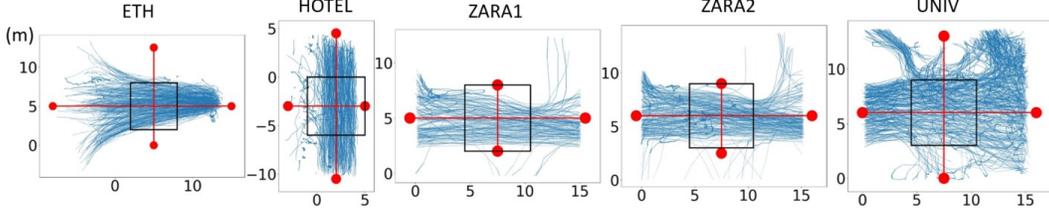


Figure 2: Trajectories of all pedestrians in the datasets. The red dots represent the task start and end locations. The red lines represent the task paths. The black box represents the test region to check for non-trivial tasks.

188 pedestrians. In the *Online* one, the robot navigates among a crowd<sup>2</sup> moving by running ORCA [16],  
 189 a policy that is frequently used as a simulation engine for benchmarking in the social navigation  
 190 literature[53, 8, 54].

191 To investigate the value of G-MPC, we develop three variants of it. **group-auto** is a G-MPC in  
 192 which the autoencoder has a history  $h = 8$  and a horizon  $f = 8$ . **group-nopred** is a variant that  
 193 features no prediction at all –it just reacts to observed groups at every timesteps and it is equivalent  
 194 to the framework of Yang and Peters [22]. Finally, **laser-group-auto** is identical to **group-auto**  
 195 but instead of using ground-truth pose information, it takes as input noisy lidar scan readings. We  
 196 simulate this by modeling pedestrians as  $1m$ -diameter circles and lidar scans as rays projecting from  
 197 the robot. We refer to the spec sheet of a SICK LMS511 2D lidar for simulation parameters. We  
 198 further inject noise into the readings according to the spec sheet. Under this simulation, pedestrians  
 199 may only be partially observable or even completely occluded from the robot.

200 We compare the performance of these policies against a set of MPC variants using mechanisms  
 201 for individual motion prediction. **ped-nopred** is a vanilla MPC that reacts to the current states of  
 202 other agents without making predictions about their future states. **ped-linear** is a vanilla MPC that  
 203 estimates future states of agents by propagating agents’ current velocities forward. This baseline  
 204 is motivated by recent work showing that constant-velocity models yield competitive performance  
 205 in pedestrian motion prediction tasks [55]. Finally, **ped-sgan** is an MPC that uses S-GAN [17] to  
 206 extract a sequence of future state predictions for agents based on inputs of their past states. We  
 207 selected S-GAN because it is a recent highly performing model.

208 We measure the performance of the policies with respect to four different metrics: a) *Success rate*,  
 209 defined as the ratio of successful trials over total number of trials; b) *Comfort*, defined as the ratio  
 210 of trials in which the robot does not enter any social group space over the total number of trials; c)  
 211 *Minimum distance to pedestrians*, defined as the smallest distance between the robot and any agent  
 212 per trial; d) *Path length*, defined as the total distance traversed by the robot in a trial.

213 To track the performance of G-MPC, we design a set of hypotheses targeting aspects of safety and  
 214 group space violation which we investigate under both experimental conditions, i.e., offline and  
 215 online:

216 **H1:** To explore the benefits of group based representations alone, we hypothesize that **group-nopred**  
 217 is safer than **ped-nopred** while achieving similar success rates but worse efficiency.

218 **H2:** To explore the full benefit of group based formulation, we hypothesize that **group-auto** is safer  
 219 than **ped-linear** and **ped-sgan** while achieving similar success rates but worse efficiency.

220 **H3:** To explore how our formulation handles imperfect inputs, we hypothesize that **laser-group-**  
 221 **auto** achieves similar safety to **group-auto** while achieving similar success rate and efficiency.

222 **H4:** To check that our formulation is socially compliant, we hypothesize that **group-nopred**, **group-**  
 223 **auto** and **laser-group-auto** violate agents’ group space less often than the baselines.

## 224 6.2 Results

225 **Quantitative Analysis.** Fig. 3 and Fig. 4 contain bar charts representing the performance of G-MPC  
 226 compared with its baselines under Offline and Online settings respectively. Bars indicate means,  
 227 errorbars indicate standard deviations, “F” and “C” are flow and cross scenarios respectively, and

<sup>2</sup>For consistency, the agents in the crowd start and end at the same spots as the agents in the recorded crowd from the Offline condition.

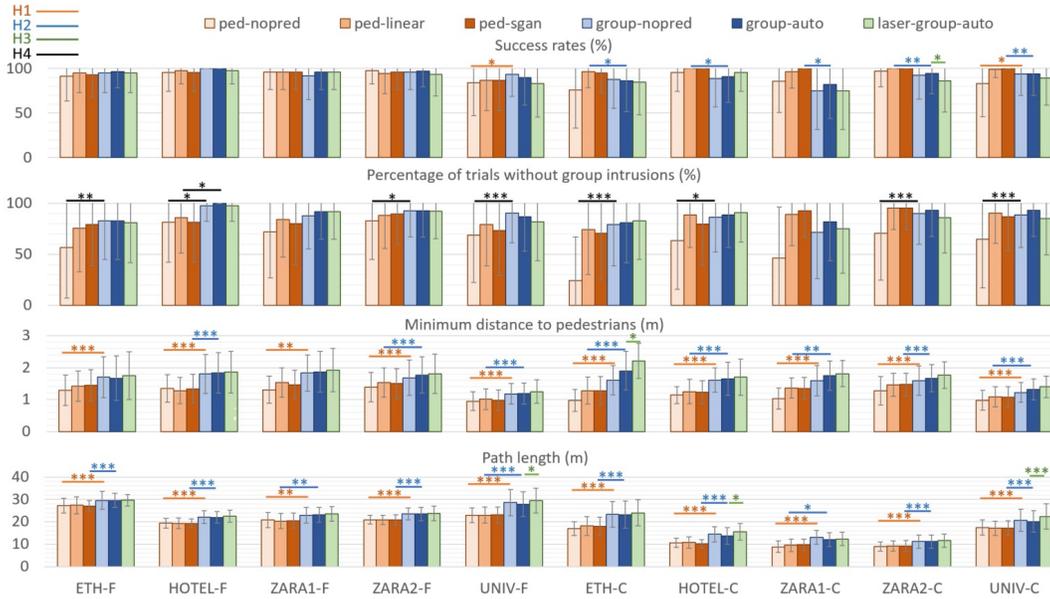


Figure 3: Performance per scene under the *Offline* condition. Horizontal lines indicate statistically significant results corresponding to different hypotheses.

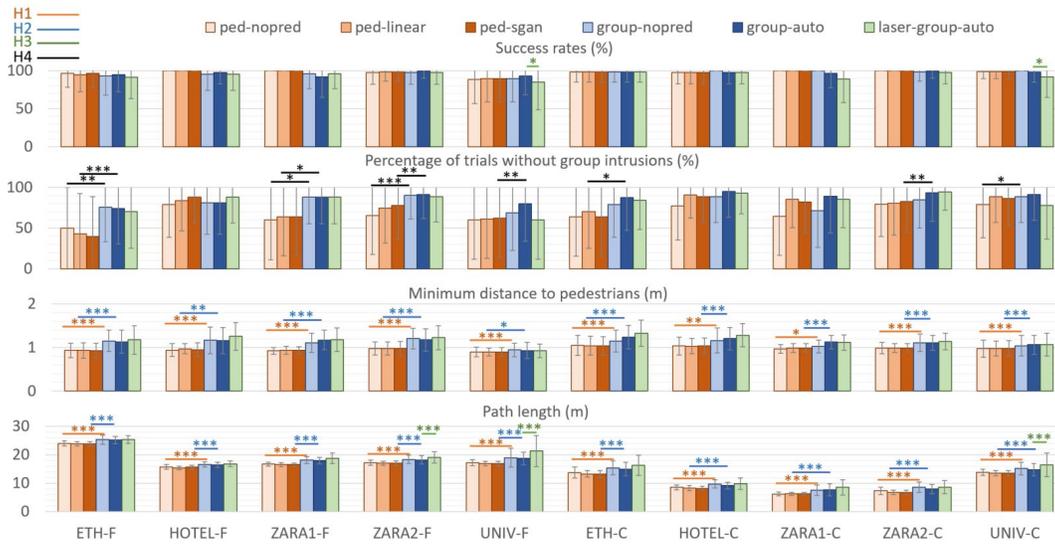


Figure 4: Performance per scene under the *Online* condition (simulated pedestrians powered by ORCA [16]). Horizontal lines indicate statistically significant results corresponding to different hypotheses.

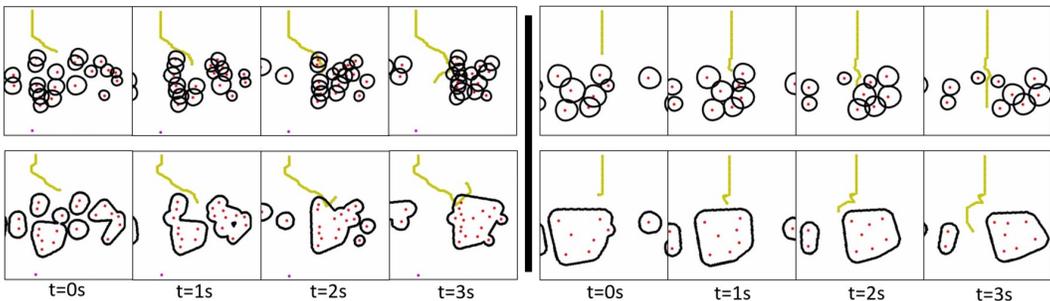


Figure 5: Qualitative performance difference between approaches leveraging pedestrian-based (top) and group-based (bottom) representations. Left: non-reactive agents. Right: reactive agents.

228 the number of asterisks indicates increasing significance levels:  $\alpha = 0.05, 0.01, 0.001$  according to  
229 two-sided Mann-Whitney U-tests.

230 **H1:** We can see from both Fig. 3 and Fig. 4 that G-MPC achieves statistically significantly larger  
231 minimum distances to pedestrians across all scenarios, often with  $p < 0.001$ . This illustrates that  
232 the group representation is in itself capable of upgrading a simple MPC with no prediction. As  
233 expected, we observe that the price G-MPC pays for that is a larger average path length. We also see  
234 that success rates are comparable. Overall, we conclude that H1 holds.

235 **H2:** When future state predictions are considered, G-MPC obtains statistically significant results in  
236 most scenes supporting its attributes of being safer at the cost of worse efficiency. Thus H2 is partially  
237 confirmed. In offline scenarios, G-MPC has lower success rates in crossing scenarios. Upon  
238 closer inspection, most failure cases are due to timeouts from G-MPC’s conservative behavior. How-  
239 ever, in online scenarios where pedestrians react to the robot, G-MPC achieves high success rates.  
240 In real-world situations, to cross dense traffic, the robot needs to plan its actions with expectations of  
241 reactive pedestrians. Otherwise, the robot will most probably run into *the freezing robot problem* [4].

242 **H3:** Group-based representations have the potential to robustly account for imperfect state-  
243 estimates. Overall, we observe that with simulated imperfect states, G-MPC does not perform  
244 statistically significantly worse in terms of safety, but in dense crowds of the UNIV scenes it has  
245 worse efficiency and worse success rates in online cases. This shows that H3 holds in terms of  
246 safety and, in moderately dense human crowds, holds in terms of efficiency. Future work on better  
247 group representation is needed to achieve better efficiency in high-density human crowds given  
248 imperfect states.

249 **H4:** From Fig. 3 and Fig. 4, we can see that G-MPC often has fewer group-space intrusions than  
250 its baselines. While this relationship is not always statistically significant, we do see a general trend  
251 of the group-based approaches to respect group spaces more often than individual ones. Thus, we  
252 conclude that H4 is partially confirmed.

253 **Qualitative Analysis.** Qualitatively, it is a more common occurrence for regular MPC to perform  
254 aggressive and socially inappropriate maneuvers than G-MPC. As shown in the two examples in  
255 Fig. 5 executed by **ped-sgan** and **group-auto** agents, we can see that in offline conditions, the MPC  
256 agent aggressively cuts in front of the two pedestrians to the left before proceeding headlong into  
257 the cluster of pedestrians, only managing to avoid the deadlock by escaping through the narrow gap  
258 that opens up. While for G-MPC, it tracks the movements of the two pedestrian groups coming from  
259 the left. When the two pedestrian groups merge, the agent turns around and reevaluates its approach  
260 to cross. In the online condition, we observe that the MPC agent cuts through a pedestrian group to  
261 reach the other side, forcing a member of the group to stop and yield as indicated by the pedestrian’s  
262 shrinking personal space, which is proportional to its speed. In the same situation, the G-MPC agent  
263 chooses to circumvent behind the social group.

## 264 7 Conclusion

265 We introduced a methodology of generating group-based representations and predicting their future  
266 states. Through an extensive evaluation over the flow and crossing scenarios drawn from 10 different  
267 real-world scenes from 2 different human datasets with both reactive and non-reactive agents, we  
268 demonstrate that our approach is safer and more socially compliant. Through experimentation with  
269 simulated laser scans, our model displays promising potential to process noisy sensor inputs without  
270 much performance downgrade.

271 Various improvements to our control framework are possible. For example, we could incorporate  
272 state-of-the-art oracles in the form of advanced video prediction models [56]. Further, additional  
273 considerations such as the set of rollouts or the cost functions could possibly increase performance.  
274 Finally, alternative control frameworks such as reinforcement learning approaches could be appli-  
275 cable. However, our goal in this paper was to illustrate the value of group-based representations for  
276 navigation tasks. Future work will involve improving both the prediction and the control components  
277 of our framework.

278 Finally, we plan on validating our findings on a real-world robot to fully test the capability of G-  
279 MPC to handle noisy sensor inputs. We also plan to investigate better group representation to reduce  
280 computation time and improve its effectiveness in high density human crowds.

## References

- 281
- 282 [1] S. Thrun, M. Bennewitz, W. Burgard, A. Cremers, F. Dellaert, D. Fox, D. Hahnel, C. Rosen-  
283 berg, N. Roy, J. Schulte, and D. Schulz. MINERVA: a second-generation museum tour-guide  
284 robot. In *Proceedings of the IEEE International Conference on Robotics and Automation*  
285 *(ICRA)*, volume 3, pages 1999–2005, 1999.
- 286 [2] T. Kruse, A. K. Pandey, R. Alami, and A. Kirsch. Human-aware robot navigation: A survey.  
287 *Robotics and Autonomous Systems*, 61(12):1726–1743, 2013.
- 288 [3] H. Kretzschmar, M. Spies, C. Sprunk, and W. Burgard. Socially compliant mobile robot navi-  
289 gation via inverse reinforcement learning. *The International Journal of Robotics Research*, 35  
290 (11):1289–1307, 2016.
- 291 [4] P. Trautman, J. Ma, R. M. Murray, and A. Krause. Robot navigation in dense human crowds:  
292 Statistical models and experimental studies of human-robot cooperation. *International Journal*  
293 *of Robotics Research*, 34(3):335–356, 2015.
- 294 [5] C. I. Mavrogiannis, A. M. Hutchinson, J. Macdonald, P. Alves-Oliveira, and R. A. Knepper.  
295 Effects of distinct robotic navigation strategies on human behavior in a crowded environment.  
296 In *Proceedings of the 2019 ACM/IEEE International Conference on Human-Robot Interaction*  
297 *(HRI '19)*. ACM, 2019.
- 298 [6] M. Luber, L. Spinello, J. Silva, and K. Arras. Socially-aware robot navigation: A learning  
299 approach. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and*  
300 *Systems (IROS)*, pages 902–907, 2012.
- 301 [7] B. Kim and J. Pineau. Socially adaptive path planning in human environments using inverse  
302 reinforcement learning. *International Journal of Social Robotics*, 8(1):51–66, 2016.
- 303 [8] M. Everett, Y. F. Chen, and J. P. How. Motion planning among dynamic, decision-making  
304 agents with deep reinforcement learning. In *IEEE/RSJ International Conference on Intelligent*  
305 *Robots and Systems (IROS)*, Madrid, Spain, Sept. 2018.
- 306 [9] C. Mavrogiannis, V. Blukis, and R. A. Knepper. Socially competent navigation planning by  
307 deep learning of multi-agent path topologies. In *Proceedings of the IEEE/RSJ International*  
308 *Conference on Intelligent Robots and Systems (IROS)*, pages 6817–6824, 2017.
- 309 [10] K. Koffka. *Principles of Gestalt psychology*. Harcourt, Brace, 1935.
- 310 [11] A. Desolneux, L. Moisan, and J.-M. Morel. *From Gestalt Theory to Image Analysis: A Prob-*  
311 *abilistic Approach*. Springer Publishing Company, Incorporated, 1st edition, 2007. ISBN  
312 0387726357.
- 313 [12] M. Vázquez and A. Steinfeld. An assisted photography method for street scenes. In *2011 IEEE*  
314 *Workshop on Applications of Computer Vision (WACV)*, pages 89–94, 2011.
- 315 [13] A. Wang and A. Steinfeld. Group split and merge prediction with 3D convolutional networks.  
316 *IEEE Robotics and Automation Letters*, 5(2):1923–1930, 2020.
- 317 [14] S. Pellegrini, A. Ess, K. Schindler, and L. van Gool. You’ll never walk alone: Modeling social  
318 behavior for multi-target tracking. In *Proc. IEEE Int. Conf. Comput. Vis.*, pages 261–268, Sept  
319 2009.
- 320 [15] A. Lerner, Y. Chrysanthou, and D. Lischinski. Crowds by example. *Comput. Graph. Forum*,  
321 26(3):655–664, 2007.
- 322 [16] J. van den Berg, S. J. Guy, M. Lin, and D. Manocha. Reciprocal n-body collision avoidance.  
323 In *Robotics Research*, pages 3–19. Springer Berlin Heidelberg, 2011.
- 324 [17] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi. Social GAN: Socially acceptable  
325 trajectories with generative adversarial networks. In *Proceedings of the IEEE Conference on*  
326 *Computer Vision and Pattern Recognition (CVPR)*, pages 2255–2264, 2018.

- 327 [18] C. I. Mavrogiannis, W. B. Thomason, and R. A. Knepper. Social momentum: A frame-  
328 work for legible navigation in dynamic multi-agent environments. In *Proceedings of the 2018*  
329 *ACM/IEEE International Conference on Human-Robot Interaction (HRI '18)*, pages 361–369.  
330 ACM, 2018.
- 331 [19] C. Chen, Y. Liu, S. Kreiss, and A. Alahi. Crowd-robot interaction: Crowd-aware robot navi-  
332 gation with attention-based deep reinforcement learning. In *Proceedings of the IEEE Interna-*  
333 *tional Conference on Robotics and Automation (ICRA)*, pages 6015–6022, 2019.
- 334 [20] J. Šochman and D. C. Hogg. Who knows who - inverting the social force model for finding  
335 groups. In *Proceedings of the IEEE International Conference on Computer Vision (ICCV)*  
336 *Workshops*, pages 830–837, 2011.
- 337 [21] A. Kendon. Conducting interaction : Patterns of behavior in focused encounters. *Studies in*  
338 *International Sociolinguistics*, 7, 1990.
- 339 [22] F. Yang and C. Peters. Social-aware navigation in crowds with static and dynamic groups. In  
340 *2019 11th International Conference on Virtual Worlds and Games for Serious Applications*  
341 *(VS-Games)*, pages 1–4, 2019.
- 342 [23] L. Bazzani, M. Cristani, and V. Murino. Decentralized particle filter for joint individual-group  
343 tracking. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*  
344 *(CVPR)*, pages 1886–1893, 2012.
- 345 [24] M. Chang, N. Krahnstoever, and W. Ge. Probabilistic group-level motion analysis and scenario  
346 recognition. In *Proceedings of the International Conference on Computer Vision (ICCV)*, pages  
347 747–754, 2011.
- 348 [25] G. Gennari and G. D. Hager. Probabilistic data association methods in visual tracking of  
349 groups. In *Proceedings of the IEEE Computer Society Conference on Computer Vision and*  
350 *Pattern Recognition (CVPR)*, volume 2, pages II–II, 2004.
- 351 [26] S. Pellegrini, A. Ess, and L. Van Gool. Improving data association by joint modeling of pedes-  
352 trian trajectories and groupings. In K. Daniilidis, P. Maragos, and N. Paragios, editors, *Com-*  
353 *puter Vision – ECCV 2010*, pages 452–465, Berlin, Heidelberg, 2010. Springer Berlin Heidel-  
354 berg. ISBN 978-3-642-15549-9.
- 355 [27] M. Zanotto, L. Bazzani, M. Cristani, and V. Murino. Online bayesian nonparametrics for  
356 group detection. In *Proceedings of the British Machine Vision Conference (BMVC)*, pages  
357 111.1–111.12, 2012.
- 358 [28] I. Chamveha, Y. Sugano, Y. Sato, and A. Sugimoto. Social group discovery from surveillance  
359 videos: A data-driven approach with attention-based cues. In *Proceedings of the The British*  
360 *Machine Vision Association (BMVC)*, 2013.
- 361 [29] S. D. Khan, G. Vizzari, S. Bandini, and S. Basalamah. Detection of social groups in pedestrian  
362 crowds using computer vision. In S. Battiato, J. Blanc-Talon, G. Gallo, W. Philips, D. Popescu,  
363 and P. Scheunders, editors, *Advanced Concepts for Intelligent Vision Systems*, pages 249–260.  
364 Springer International Publishing, Cham, 2015.
- 365 [30] D. Helbing and P. Molnár. Social force model for pedestrian dynamics. *Physical Review E*, 51  
366 (5):4282–4286, 1995.
- 367 [31] R. Mazzon, F. Poiesi, and A. Cavallaro. Detection and tracking of groups in crowd. In *Proceed-*  
368 *ings of the IEEE International Conference on Advanced Video and Signal Based Surveillance*  
369 *(AVSS)*, pages 202–207, 2013.
- 370 [32] F. Solera, S. Calderara, and R. Cucchiara. Socially constrained structural learning for groups  
371 detection in crowd. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 38(5):  
372 995–1008, 2016.
- 373 [33] W. Ge, R. T. Collins, and R. B. Ruback. Vision-based analysis of small groups in pedestrian  
374 crowds. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 34(5):1003–1016,  
375 2012.

- 376 [34] A. Taylor, D. M. Chan, and L. D. Riek. Robot-centric perception of human groups. *ACM*  
377 *Transactions on Human-Robot Interaction*, 9(3):1–21, 2020.
- 378 [35] I. Chatterjee and A. Steinfeld. Performance of a low-cost, human-inspired perception approach  
379 for dense moving crowd navigation. In *Proceedings of the IEEE International Symposium on*  
380 *Robot and Human Interactive Communication (RO-MAN)*, pages 578–585, Aug 2016.
- 381 [36] N. P. Cuntoor, R. Collins, and A. J. Hoogs. Human-robot teamwork using activity recognition  
382 and human instruction. In *Proceedings of the IEEE/RSJ International Conference on Intelligent*  
383 *Robots and Systems (IROS)*, pages 459–465, 2012.
- 384 [37] R. Gockley, J. Forlizzi, and R. Simmons. Natural person-following behavior for social  
385 robots. In *Proceedings of the ACM/IEEE International Conference on Human-Robot Inter-*  
386 *action (HRI)*, pages 17–24, 2007.
- 387 [38] C. Granata and P. Bidaud. A framework for the design of person following behaviors for social  
388 mobile robots. In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots*  
389 *and Systems (IROS)*, pages 4652–4659, 2012.
- 390 [39] E. Jung, B. Yi, and S. Yuta. Control algorithms for a mobile robot tracking a human in front.  
391 In *Proceedings of the IEEE/RSJ International Conference on Intelligent Robots and Systems*  
392 *(IROS)*, pages 2411–2416, 2012.
- 393 [40] H. Zender, P. Jensfelt, and G. M. Kruijff. Human- and situation-aware people following. In  
394 *Proceedings of the IEEE International Symposium on Robot and Human Interactive Commu-*  
395 *nication (RO-MAN)*, pages 1131–1136, 2007.
- 396 [41] A. Nanavati, X. Z. Tan, J. Connolly, and A. Steinfeld. Follow the robot: Modeling coupled  
397 human-robot dyads during navigation. In *2019 IEEE/RSJ International Conference on Intelli-*  
398 *gent Robots and Systems (IROS)*, pages 3836–3843, 2019.
- 399 [42] D. Feil-Seifer and M. Mataric. People-aware navigation for goal-oriented behavior involving  
400 a human partner. In *Proceedings of the IEEE International Conference on Development and*  
401 *Learning (ICDL)*, volume 2, pages 1–6, 2011.
- 402 [43] A. K. Pandey and R. Alami. A step towards a sociable robot guide which monitors and adapts  
403 to the person’s activities. In *2009 International Conference on Advanced Robotics*, pages 1–8,  
404 2009.
- 405 [44] A. Garrell and A. Sanfeliu. Local optimization of cooperative robot movements for guiding  
406 and regrouping people in a guiding mission. In *Proceedings of the IEEE/RSJ International*  
407 *Conference on Intelligent Robots and Systems (IROS)*, pages 3294–3299, 2010.
- 408 [45] M. Shiomi, T. Kanda, S. Koizumi, H. Ishiguro, and N. Hagita. Group attention control for com-  
409 munication robots with wizard of oz approach. In *Proceedings of the ACM/IEEE International*  
410 *Conference on Human-Robot Interaction (HRI)*, pages 121–128, 2007.
- 411 [46] E. A. Martinez-Garcia, Ohya Akihisa, and Shin’ichi Yuta. Crowding and guiding groups of hu-  
412 mans by teams of mobile robots. In *Proceedings of the IEEE Workshop on Advanced Robotics*  
413 *and its Social Impacts (ARSO)*, pages 91–96, 2005.
- 414 [47] K. Katyal, Y. Gao, J. Markowitz, I.-J. Wang, and C.-M. Huang. Group-Aware Robot Naviga-  
415 tion in Crowded Environments. *arXiv e-prints*, art. arXiv:2012.12291, Dec. 2020.
- 416 [48] M. Ester, H.-P. Kriegel, J. Sander, and X. Xu. A density-based algorithm for discovering  
417 clusters a density-based algorithm for discovering clusters in large spatial databases with noise.  
418 In *Proc. Int. Conf. Knowl. Discovery and Data Mining*, pages 226–231, 1996.
- 419 [49] R. Kirby. *Social Robot Navigation*. PhD thesis, Carnegie Mellon University, Pittsburgh, PA,  
420 May 2010.
- 421 [50] D. Tran, L. Bourdev, R. Fergus, L. Torresani, and M. Paluri. Learning spatiotemporal features  
422 with 3d convolutional networks. In *Proc. IEEE Int. Conf. Comput. Vis.*, December 2015.

- 423 [51] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese. Social LSTM:  
424 Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE Conference on*  
425 *Computer Vision and Pattern Recognition (CVPR)*, pages 961–971, 2016.
- 426 [52] P. Zhang, W. Ouyang, P. Zhang, J. Xue, and N. Zheng. Sr-lstm: State refinement for lstm to-  
427 wards pedestrian trajectory prediction. In *Proc. IEEE Comput. Soc. Conf. Comput. Vis. Pattern*  
428 *Recognit.*, pages 12085–12094, June 2019.
- 429 [53] C. Cao, P. Trautman, and S. Iba. Dynamic channel: A planning framework for crowd naviga-  
430 tion. In *2019 International Conference on Robotics and Automation (ICRA)*, pages 5551–5557,  
431 2019.
- 432 [54] C. Mavrogiannis, F. Baldini, A. Wang, D. Zhao, P. Trautman, A. Steinfeld, and J. Oh. Core  
433 Challenges of Social Robot Navigation: A Survey. *arXiv e-prints*, art. arXiv:2103.05668, Mar.  
434 2021.
- 435 [55] C. Schöller, V. Aravantinos, F. Lay, and A. Knoll. What the constant velocity model can teach  
436 us about pedestrian motion prediction. *IEEE Robotics and Automation Letters*, 5(2):1696–  
437 1703, 2020.
- 438 [56] V. L. Guen and N. Thome. Disentangling physical dynamics from unknown factors for unsu-  
439 pervised video prediction. In *Proceedings of the IEEE/CVF Conference on Computer Vision*  
440 *and Pattern Recognition (CVPR)*, June 2020.