
Latent Variable Models for Multi-Agent Trajectories

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Analyzing the spatiotemporal behavior of multiple agents is of great interest to
2 many communities. Existing probabilistic models in this regime employ either
3 unsupervised generative settings, in which the latent space is described fully by
4 discrete or continuous representations, or, are alternatively formalized in a (fully)
5 supervised framework where weakly preserved labels add explicit information to a
6 continuous latent representation learned from the data. To overcome the resulting
7 limitations, we propose a novel objective function for processing multi-agent trajec-
8 tories based on semi-supervised variational autoencoders, where equivariance and
9 interaction of agents are modeled via customized graph networks. Our formulation
10 disentangles discrete and continuous effects and allows discrete behavioral indica-
11 tors in arbitrary quantity and annotation type to guide the generation process. This
12 lifts applicability to relevant prediction problems beyond the generation of collec-
13 tive movements and provides an effective solution to incorporate expensive domain
14 knowledge into interactive multi-agent systems. Empirically, our model outper-
15 forms various state-of-the-art baselines in generating future agent movements on
16 interactive real-world datasets. We also show that our approach effectively learns to
17 leverage unsupervised multi-agent sequences to improve classification of long-term
18 locations as well as manually annotated situations on sports tracking data.

19 1 Introduction

20 Analyzing the spatiotemporal behavior of multiple agents bears a great deal of value in many domains
21 such as autonomous driving [49, 13], public transportation [30, 58], migration [39, 22] and sports
22 analytics [53]. Particularly, the detection of collective patterns across time and space is of interest to
23 many communities but non-trivial inter-dependencies between agents render estimating the inherent
24 multi-modal distributions often difficult.

25 Existing generative approaches to modeling sequential multi-agent data thus capitalize on variants
26 from graph networks [55, 12] and (sequential) variational autoencoders [26, 7]; being purely unsuper-
27 vised, their functionality is grounded in uncovering the implicit modular structure in agent-wise latent
28 variables directly from the observed trajectories via approximate inference [28, 52, 59, 16, 14, 33].
29 Although a generative view on the problem is very appealing, approaches in this regime either lack
30 interpretable and controllable latent factors (e.g., [59]), and are thus limited in practical applicability,
31 or put a strong focus on relationship discovery (e.g., [28]), which limits their generative capacity.

32 An alternative is offered by supervised approaches that aim to counteract these issues by incorporating
33 domain-knowledge in form of discrete labels into the generation process [64]. Ideally, these labels
34 encode implicit expert knowledge that would be learned by the model and substantially expand the
35 scope of information captured in latent space. Instead of arguing in favor of expensive and tedious
36 manual annotations of spatiotemporal data, however, recent approaches propose to incorporate
37 heuristic surrogates as makeshift labels into the problem [65, 64, 40]. While reporting impressive
38 predictive accuracies, they require the existence of label sequences for all training instances and

39 cannot include inexpensive and possibly abundant unlabeled data due to their purely supervised
40 nature.

41 To close this gap, we propose semi-supervised variational autoencoders for spatiotemporal multi-agent
42 problems to process (weakly) labeled as well as unlabeled data. Our contribution generalizes the
43 class of semi-supervised variational autoencoders [27] to spatiotemporal domains by using ideas from
44 variational recurrent neural networks [7]. At a high level, the derived objective function subsumes
45 previous formalizations into a unified framework via a disentangled latent space that can incorporate
46 supervision signals in arbitrary quantities for structuring the discrete latent subspace. Additional
47 adaptive graph network layers provide order in- and equivariance of agents and render the proposed
48 approach appropriate for multi-agent scenarios. The resulting semi-supervised approach is applicable
49 to a wide range of problems including the generation of collective movements as well as classifying
50 situations of interest. Empirically, our contribution significantly advances the state-of-the art in
51 generation and classification tasks on interactive real-world data.

52 2 Related Work

53 **Semi-supervised generative models** Semi-supervised variational autoencoders have been origi-
54 nally studied by [27] who propose to incorporate additional labels in latent space. Related approaches
55 introduce auxiliary variables that leave the original model unchanged but increase the flexibility of
56 the variational posterior [34]. [24] propose to explicitly capture label characteristics in latent space
57 instead of the label values themselves and [41] study generalizations that allow for learning more
58 complex dependency structures. However, all the above approaches are introduced only for static
59 domains. Our contribution constitutes a generalization of the above approaches to spatiotemporal
60 domains.

61 **Sequential generative models** Approaches designed for predicting agent trajectories are frequently
62 based on ideas from sequential generative models to encode temporal dependencies and address
63 the stochasticity inherent in future predictions. Sequential extensions of variational autoencoders
64 with dynamic latent variables include [3, 7, 11, 15]. Particularly, the VRNN [7] is related to our
65 contribution as both deploy a VAE instance per time step conditioned on a recurrent state. However,
66 their work lacks modeling of a social dimension as well as means to integrate (potentially expensive)
67 discrete behavioral indicators. Other approaches associate a single global latent variable with each
68 sequence [4, 17, 9, 8].

69 **Graph neural networks** Deep graph-based (GNNs) approaches [50] are a natural methodological
70 choice for applications requiring modeling datasets composed of sets. GNNs operate by learning a
71 chain of hidden representations for each node through an iterative process that relies on aggregating
72 messages along edges in each network layer. Prominent instantiations such as GCN [29], GraphSAGE
73 [18], GAT [55], and others [12, 54, 20, 2, 63, 57] mainly differ in their notion of message passing and
74 neighborhood aggregation [60]. See [5] for a discussion on the function approximation capabilities
75 of this neural network family.

76 **Generative models for sequential multi-agent data** Given the potential benefit across various
77 domains, a great deal of recent publications focus on movements of pedestrians or self-driving cars
78 [37, 49, 13, 38, 61, 48, 1, 21, 62, 32]. [36], however, show that these standard benchmark datasets
79 generally exhibit weak social interaction patterns, so in the remainder, we focus on methods using
80 team sport data. Related approaches show that learning and computing realistic rollouts significantly
81 benefit from the incorporation of makeshift annotations and heuristic surrogates into the problem
82 [65, 64, 40]. [10, 33] propose an architecture based on conditional VAEs [51] to model distributions
83 of future movements of basketball players. [59, 52] combine ideas from graph networks [12] and
84 VRNNs [7] into a unified framework aiming at modeling data from basketball and soccer. Other
85 related approaches [31, 28, 16, 14] aim to explicitly infer discrete latent variables to represent
86 interaction types while executing a trajectory prediction task.

87 Hence, existing generative models for spatiotemporal pattern detection focus on multi-agent trajec-
88 tory prediction. Our contribution constitutes the first semi-supervised generalization of variational
89 autoencoders for spatiotemporal domains that gives rise to addressing tasks beyond representational
90 or generative modeling.

91 3 Methods

92 **3.1 Preliminaries**

93 Variational autoencoders (VAEs) [26, 43] allow to capture intrinsic conditional dependencies and
 94 are often used in structured problems. Essentially, they consist of a θ -parameterized generative
 95 model $p_\theta(\mathbf{x}, \mathbf{z}) = p(\mathbf{z})p_\theta(\mathbf{x}|\mathbf{z})$ and an associated ϕ -parameterized variational distribution $q_\phi(\mathbf{z}|\mathbf{x})$.
 96 Optimal parameters $\{\theta, \phi\}$ are obtained via maximizing the variational lower bound on the marginal
 97 likelihood, given by

$$\log p_\theta(\mathbf{x}) \geq \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})}[\log p_\theta(\mathbf{x}|\mathbf{z})] - \mathcal{KL}[q_\phi(\mathbf{z}|\mathbf{x}) \parallel p(\mathbf{z})].$$

98 It is often helpful to augment discrete labels y into the model to guide the generating process [27].
 99 The resulting model $p(\mathbf{x}, \mathbf{z}, y) = p(\mathbf{x}|\mathbf{z}, y)p(y)p(\mathbf{z})$ then acts on labeled and unlabeled instances.
 100 Consequentially, the variational distribution $q_\phi(y, \mathbf{z}|\mathbf{x}) = q_\phi(y|\mathbf{x})q_\phi(\mathbf{z}|y, \mathbf{x})$ needs to be defined
 101 over both quantities (\mathbf{z} and y). *Semi-supervised VAEs (SSVAE)* are trained by maximizing

$$\sum_{(\mathbf{x}, y)} (\mathcal{L}(\mathbf{x}, y) + \lambda \log q_\phi(y|\mathbf{x})) + \sum_{\mathbf{x}} \mathcal{U}(\mathbf{x}),$$

102 where $\mathcal{L}(\mathbf{x}, y) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x}, y)}[\log p_\theta(\mathbf{x}|\mathbf{z}, y) + \log p(y) + \log p(\mathbf{z}) - \log q_\phi(\mathbf{z}|\mathbf{x}, y)]$ and $\mathcal{U}(\mathbf{x}) =$
 103 $\sum_y q_\phi(y|\mathbf{x})(\mathcal{L}(\mathbf{x}, y) + \mathcal{H}(q_\phi(y|\mathbf{x})))$ denote (variational) lower bounds on labeled and unlabeled
 104 instances, respectively, and λ is a balancing term. In fact, the original motivation for this framework
 105 was inspired by semi-supervised classification tasks. However, their contribution is also regularly
 106 used for learning meaningful representations and generating new data, thus rendering it a suitable
 107 methodological foundation for performing tasks with all sorts of data. The dependency structures of
 108 q_ϕ and p_θ for SSVAEs are displayed visually as part of Figure 1 (right side).

109 **3.2 Problem Formulation**

110 We are interested in modeling spatiotemporal multi-agent scenarios and focus on positions $\mathbf{x}_t^{(a)} \in \mathbb{R}^2$
 111 of an agent $a \in \mathcal{A}$ over time $1 \leq t \leq T$, denoted by $\mathbf{x}_{\leq T}^{(a)}$. We assume that a non-empty subset of
 112 agents in \mathcal{A} is present in any observed segment and collect their trajectories (in random order) to form
 113 collective movements $\mathbf{x}_{\leq T} := \{\mathbf{x}_{\leq T}^{(a)} \mid a \in \mathcal{A}\}$. We will make use of velocities $d\mathbf{x}/dt$ and linearized
 114 motion $\Delta\mathbf{x}_t = \mathbf{x}_{t'} - \mathbf{x}_t$, for $t' > t$, of agents in the remainder.

115 In addition, we introduce discrete labels $y_t^{(a)} \in \mathcal{Y}$ associated with agent $a \in \mathcal{A}$ at time t . The
 116 collection of sequences $\mathbf{y}_{\leq T} := \{\mathbf{y}_{\leq T}^{(a)} \mid a \in \mathcal{A}\}$ may be best thought of as discriminative causes
 117 of variation in the corresponding trajectories $\mathbf{x}_{\leq T}$. In practice, $\mathbf{y}_{\leq T}$ may contain anything that
 118 aids in generating multi-agent rollouts, including behavioral indicators, long-term goals, or current
 119 tasks. Finally, observations $\mathbf{x}_{\leq T}$ accompanied by labels $\mathbf{y}_{\leq T}$ assemble the labeled share of the
 120 data, denoted by $\mathcal{D}_L = \{(\mathbf{x}_{\leq T}, \mathbf{y}_{\leq T})_n\}_{1 \leq n \leq N}$, while unlabeled observations are collected in
 121 $\mathcal{D}_U = \{(\mathbf{x}_{\leq T})_m\}_{1 \leq m \leq M}$. We refer to all data by $\mathcal{D} = \mathcal{D}_L \cup \mathcal{D}_U$.

122 We aim at estimating the underlying distribution that generated the observations in \mathcal{D} . A straight
 123 forward solution is to maximize the log likelihood given by

$$\log p(\mathcal{D}) = \sum_{\mathcal{D}_L} \log p_\theta(\mathbf{x}_{\leq T}, \mathbf{y}_{\leq T}) + \sum_{\mathcal{D}_U} \log p_\theta(\mathbf{x}_{\leq T}) \quad (1)$$

124 However, human movement is inherently non-deterministic and multi-modal. Thus, it is common to
 125 model the stochasticity inherent in future agent behaviors via latent representations of VAEs. We
 126 propose to lower-bound Eq. 1 as follows.

127 **3.3 SSVAEs for Multi-Agent Trajectories**

128 We associate disentangled latent realizations $\{z_t, y_t\}$ with each time step t of the segment while
 129 injecting the dependency structures introduced in Section 3.1 for the generative and inference parts.
 130 To model temporal dependencies, we propose to additionally condition the model components on the
 131 past observations $\mathbf{x}_{<t}$ and latent variables $\mathbf{z}_{<t}$ and $\mathbf{y}_{<t}$ (cmp. [7]). The joint distribution factorizes
 132 into

$$p_\theta(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}, \mathbf{y}_{\leq T}) = \prod_t \left[p_\theta(\mathbf{x}_t | \mathbf{x}_{<t}, \mathbf{z}_{\leq t}, \mathbf{y}_{\leq t}) p_\theta(\mathbf{z}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t}, \mathbf{y}_{<t}) p_\theta(\mathbf{y}_t | \mathbf{x}_{<t}, \mathbf{z}_{<t}, \mathbf{y}_{<t}) \right],$$

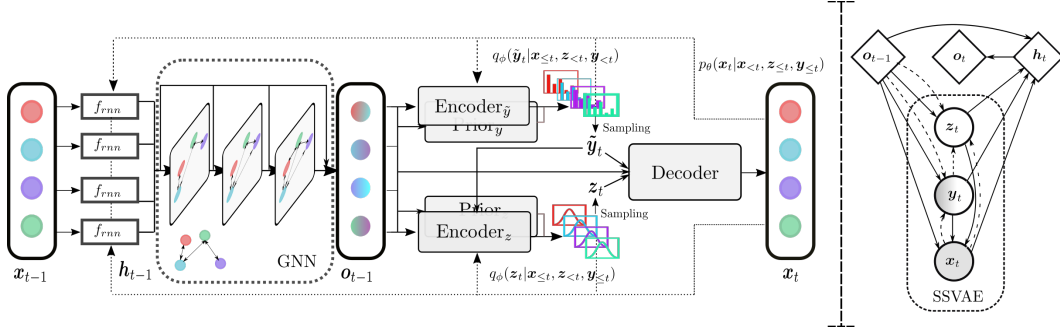


Figure 1: *Right*: Illustrative graphical model at t . Dashed lines indicate the encoding procedure, and solid lines the generative model. We use a semi-supervised VAE (SSVAE, Section 3.1) per time step and additionally introduce variables h_t and o_t . The representations in set o_t jointly encapsulate both temporal and spatial patterns. The proposed model can accommodate all contingent proportions of supervision for the discrete latent subspace y_t . *Left*: The full computational logic when processing unobserved multi-agent segments. Given temporal context encoded using an RNN module, the model leverages an customized GNN architecture based on [55, 57] (Section 3.5) to infer posteriors (and priors) over the discrete (y_t) and continuous (z_t) latent factors. The realized latent values, in conjunction with the GNN output, serve as input for the decoding module, which generates distribution over agent movements.

133 with generating distribution $p_\theta(x_t|x_{<t}, z_{<t}, y_{<t})$ and prior distributions $p_\theta(z_t|x_{<t}, z_{<t}, y_{<t})$ and
 134 $p_\theta(y_t|x_{<t}, z_{<t}, y_{<t})$ for latents z_t and y_t . Note that we differentiate between observed labels y and
 135 latent variables \tilde{y} to derive the variational distributions.

136 Also note that in our sequential setting, the priors are included in the optimization rather than being
 137 represented by fixed distributions to account for latent dynamics and effective sampling at inference
 138 time. Hence, the surrogate posteriors q_ϕ also need to condition on previous observations and random
 139 variables to allow for accurate approximations of the true posterior and provide a tractable lower
 140 bound on the log-likelihood in Eqn (1). Using the factorization introduced in Section 3.1 yields

$$q_\phi(z_{\leq T}|x_{\leq T}, y_{\leq T}) = \prod_t q_\phi(z_t|x_{\leq t}, z_{<t}, y_{\leq t}) \text{ and } q_\phi(\tilde{y}_{\leq T}|x_{\leq T}) = \prod_t q_\phi(\tilde{y}_t|x_{\leq t}, z_{<t}, y_{<t}). \quad (2)$$

141 Labeled instances render reasoning over latents $\tilde{y}_{\leq T}$ unnecessary since the true $y_{\leq T}$ is already
 142 known. Hence, the right-hand side in Eqn (2) can be ignored when sampling from \mathcal{D}_L . We arrive at
 143 the following results and refer to Appendix A for the proofs.

144 **Theorem 1.** A lower bound on $\log p_\theta(x_{\leq T}, y_{\leq T})$ in Eqn (1) is given by

$$\log p_\theta(x_{\leq T}, y_{\leq T}) \geq \sum_t \log p_\theta(y_t|x_{<t}, z_{<t}, y_{<t}) + \mathbb{E}_{q_\phi(z_t|x_{\leq t}, z_{<t}, y_{\leq t})} \left[\log p_\theta(x_t|x_{<t}, z_{\leq t}, y_{\leq t}) \right] \\ - \mathcal{KL}[q_\phi(z_t|x_{\leq t}, z_{<t}, y_{\leq t}) \parallel p_\theta(z_t|x_{<t}, z_{<t}, y_{<t})] \equiv \sum_{t=1}^T -\mathcal{L}(x_t, y_t). \quad (3)$$

145 Correspondingly, for unlabeled instances drawn from \mathcal{D}_U , sequential label information need to be
 146 estimated as shown in the following theorem.

147 **Theorem 2.** Let $\mathcal{H}(\beta)$ be the entropy of quantity β . A lower bound on $\log p_\theta(x_{\leq T})$ in Eqn (1) is
 148 given by

$$\log p_\theta(x_{\leq T}) \geq \sum_t \left(\mathcal{H}(q_\phi(\tilde{y}_t|x_{\leq t}, z_{<t}, y_{<t})) - \mathbb{E}_{q_\phi(\tilde{y}_t|x_{\leq t}, z_{<t}, y_{<t})} [\mathcal{L}(x_t, \tilde{y}_t)] \right) \equiv \sum_t -\mathcal{U}(x_t) \quad (4)$$

149 The full evidence lower bound (ELBO) is obtained by combining the lower bounds for all data \mathcal{D} ,

$$p(\mathcal{D}) \geq \sum_{\mathcal{D}_L} \sum_t -\mathcal{L}(x_t, y_t) + \sum_{\mathcal{D}_U} \sum_{t=1}^T -\mathcal{U}(x_t). \quad (5)$$

150 This formalization focuses on either fully labeled or unlabeled observations, however, a more general
 151 formulation can be obtained to also allow for partially labeled observations.

152 3.4 Semi-supervised Sequential Learning

153 Intrinsic to our contribution is an encoding module that argues over the label space, $q_\phi(\mathbf{y}_{<T}|\mathbf{x}_{\leq T})$,
 154 and hence can be used as a classification network annotating new observations. However, the ELBO
 155 in Eqn (5) is oblivious to classification and simply ignores relevant gradient updates for $\mathcal{L}(\mathbf{x}_t, \mathbf{y}_t)$
 156 by discarding Eqn (2). This is clearly inappropriate for semi-supervised learning where the overall
 157 goal is to learn an unknown mapping $\mathbf{x} \mapsto \mathbf{y}$. Following [27], we circumvent this by heuristically
 158 incorporating Eqn (2) into the label-dependent objective. Hence, the full training criterion to learn
 159 $\{\theta, \phi\}$ is given by

$$\mathcal{J}_{MAT}(\theta, \phi; \mathcal{D}) = \lambda_0 \sum_{\mathcal{D}_U} \sum_t \mathcal{U}(\mathbf{x}_t) + \sum_{\mathcal{D}_L} \sum_t \left(\mathcal{L}(\mathbf{x}_t, \mathbf{y}_t) - \lambda_1 \log q_\phi(\mathbf{y}_t|\mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) \right), \quad (6)$$

160 where λ_1 balances generative and discriminative learning and λ_0 the contribution of labeled and
 161 unlabeled data, respectively. Thus, the semi-supervised variant contains fully (un)supervised modeling
 162 tasks as special cases via adjusting λ_1, λ_2 accordingly. The resulting supervised target differs from
 163 its unsupervised counterpart only in observing \mathbf{y}_t instead of factorizing over \mathcal{Y} , which results in
 164 computing the log-likelihood (auxiliary loss) instead of the entropy $\mathcal{H}(q_\phi(\mathbf{y}_t|\mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}))$.

165 Eqn (4) encodes the behavior of classifier $q_\phi(\mathbf{y}_{\leq T}|\mathbf{x}_{\leq T})$ when dealing with unsupervised data. Since
 166 data likelihood $\log p_\theta(\mathbf{x}_t|\mathbf{x}_{< t}, \mathbf{z}_{\leq t}, \mathbf{y}_{\leq t})$ tends to comprise the largest factor in $\mathcal{L}(\mathbf{x}_t, \mathbf{y}_t)$, taking
 167 the expectation wrt. $q_\phi(\tilde{\mathbf{y}}_t|\mathbf{x}_{< t}, \mathbf{z}_{< t}, \tilde{\mathbf{y}}_{< t})$ encourages the model to realize the highest probability
 168 mass for the latent label value $y \in \mathcal{Y}$ that incurs the smallest reconstruction loss compared to other
 169 candidates in label space. This is a desirable property that reflects our assumption of data generation.

170 3.5 Design Choices

171 As described in Section 3.2, we aim to learn a distribution over possible sequences of agent sets that
 172 may differ in share of discrete annotations (and cardinality). However, the previously derived objective
 173 functions only account for temporal dependencies, so that independence of the agent dimension
 174 is required to realize permutation-invariant models. This assumption is trivially inappropriate for
 175 interactive systems where future agent movements need to be coordinated with other agents. Thus, it
 176 is common to propose some form of graph encoding strategy [50] to account for both, equivariance
 177 and agent interactions, respectively.

178 Intuitively, optimal interaction modeling allows to adaptively accommodate structural information from
 179 different levels of granularity for each agent. In team sports, for example, agent nodes should be able
 180 to capture immediate player influences as well as holistic team strategies dependent on the interaction
 181 structure best suited for the task at-hand. Motivated by the previous considerations, we customize an
 182 GNN architecture that aims to mimic the desired behavior. Specifically, we stack attention-based
 183 GNN layers [55] with particularly designed skip-connections [57] and construct the assumed graph
 184 structure using the k (spatially) closest agents.

185 There are various possibilities to incorporate graph-based approaches into the overall scheme¹. An
 186 efficient way to model the distributions defining Eqn (6) is via hidden agent states $\mathbf{h}_t^{(a)}$ conditioned
 187 on representations $\mathbf{o}_{t-1}^{(a)}$,

$$\mathbf{h}_t^{(a)} = f_{rnn}(\mathbf{x}_t^{(a)}, \mathbf{z}_t^{(a)}, \mathbf{y}_t^{(a)}, \mathbf{o}_{t-1}^{(a)}) \quad \text{with} \quad \mathbf{o}_t = \text{GNN}(\mathbf{h}_t),$$

188 where f_{rnn} denotes an RNN transition function, GNN is the described GAT-based GNN and $\mathbf{o}_t =$
 189 $\{\mathbf{o}_t^{(a)} \mid a \in \mathcal{A}\}$ refers to the set of updated agent representations. Since the elements in set \mathbf{o}_t
 190 aggregate neighboring information of the RNN outputs \mathbf{h}_t via graph networks, the inferred feature
 191 vectors describe the entire interactive past of individual agents. While it is reasonable to additionally
 192 enforce intra-timestep dependencies on the latent variables using encoding and decoding GNN
 193 modules, we argue that capturing past spatiotemporal patterns is sufficiently informative for high-
 194 frequency data [46]. Hence, assuming conditional independence within t , the joint movement

¹We empirically evaluate different configurations in Appendix E.

195 distribution factorizes across the agent dimension and is given by

$$p_{\theta}(\mathbf{x}_t | \mathbf{x}_{<t}, \mathbf{z}_{\leq T}, \mathbf{y}_{\leq T}) = p_{\theta}(\mathbf{x}_t | \psi_t) = \prod_{a \in \mathcal{A}} p_{\theta}(\mathbf{x}_t^{(a)} | \psi_t^{(a)})$$

196 where $\psi_t = \{\psi_t^{(a)} \mid a \in \mathcal{A}\}$ are the distribution parameters and $\psi_t^{(a)} = \varphi(\mathbf{z}_t^{(a)}, \mathbf{y}_t^{(a)}, \mathbf{o}_{t-1}^{(a)})$. The
197 computational logic of the full architecture is schematically depicted in Fig. 1.

198 4 Empirical Evaluation

199 For evaluation, we focus on team sports as the coordination of players renders these tasks more
200 difficult than other domains [36].² Hence, we experiment on STATS SportVU NBA³ for comparison,
201 and tracking data from elite soccer⁴. As detailed in Section 3.2, we use agent velocities as input
202 to all models and assume linear motion between consecutive observations. For NBA, we adopt the
203 experimental setup and processing strategy from [40]⁵.

204 The STATS SportVU NBA data comprises tracking positions of offensive plays from the 2016 NBA
205 regular season covering more than 1200 different games where game segment are given by sequences
206 of length 50 and contain two-dimensional positions of all agents (10 players and ball) sampled at 5
207 frames per second. The data is split into 60% training, 20% validation, and 20% test sets. All data is
208 translated so that the origin of the underlying coordinate system is mapped onto the top-left corner.

209 The soccer data consists of 12 matches where positions of players and ball are sampled at 25 frames
210 per second. The tracking data is accompanied by manual event annotations that we will also make use
211 of in the remainder. Models are trained on eight matches, the remaining four are distributed evenly
212 into validation and test data (two matches each).

213 4.1 Baselines

214 We evaluate versus several baselines on two distinct tasks: *trajectory forecasting* (Section 4.2 and 4.4)
215 and *spatiotemporal classification* (Section 4.3 and 4.4). Since our approach constitutes the first semi-
216 supervised model for label predictions in multi-agent scenarios, we need to rely on self-constructed
217 supervised baselines for scenarios where only a few labels are available, which we discuss in more
218 detail in the relevant sections.

219 By contrast, trajectory forecasting is an active topic in different fields (see, e.g., [45] for an overview),
220 which allows comparisons with a wide range of baseline models. Fully supervised generative
221 models such as *DAG-Net* [40] and *Weak-Sup* [64] use weak labels in the form of agent objectives
222 (heuristically inferred prior to model training) to improve trajectory prediction. In these settings, an
223 agent’s objective at each timestep t is to reach the next area where the observed movement speed
224 falls below a predefined threshold. Adopting the experimental framework of [40] allows us to make
225 further comparisons with *Social-Ways* [1] and *STGAT* [21] without additional experiments. *STGAT*
226 [21] augments a standard GAT with an LSTM [19] to capture both spatial and temporal dependencies
227 of social interactions while *Social-Ways* uses Info-GAN to learn multimodal predictive distributions.

228 It is difficult to determine the current state-of-the-art in generative baselines for sports tracking data
229 because there is no consistent setting in terms of forecasting horizon or data (see e.g., [33, 40]). While
230 some baselines are restricted to making predictions for a predetermined horizon, our method learns
231 spatiotemporal representations autoregressively and thus can be used for the full range of prediction
232 scenarios. Hence, we conduct two sets of experiments over both short (10 timesteps) and long (40
233 timesteps) and additionally benchmark against the most recent and methodologically most relevant
234 work described below. *GVRNN* [59] is a graph extension of the VRNN for explicitly modeling joint
235 agent movements for basketball and soccer data. *dNRI* [16] dynamically extracts interaction types
236 that determine the parametrization of the decoder. Finally, *GRIN* [33] recovers (static) interactions
237 via attention using a disentangled latent space.

²We report on results on Stanford Drone Data (SDD) [44] in the Appendix.

³<https://github.com/alexmonti19/dagnet/tree/master/datasets>

⁴Details omitted for anonymity

⁵The source code will be made publicly available after acceptance of the paper.

Table 1: Error for NBA in meters for a prediction interval of 10 timesteps with an observation period of 40 timesteps. Bold is the lowest avg L_2 & final L_2 for the supervised and unsupervised generative models, respectively.

	GVRNN	dNRI	GRIN	DAG-Net.	U-MAT	S-MAT
avg L_2	2.60	2.77	3.00	2.10	2.31	1.94
final L_2	5.66	5.52	6.12	4.28	4.80	3.99

Table 2: Error for NBA in meters for a prediction interval of 40 timesteps with an observation period of 10 timesteps. Top rows show results for offense, bottom rows the defense. Bold is the lowest avg L_2 & final L_2 for the supervised and unsupervised generative models, respectively.

	STGAT	Social-Ways	GVRNN	Weak-Sup.	DAG-Net	U-MAT	S-MAT
avg L_2	9.94	9.91	9.73	9.47	8.98	9.01	8.11
final L_2	15.80	15.19	15.89	16.98	14.08	13.28	12.52
avg L_2	7.26	7.31	7.29	7.05	6.87	6.88	6.21
final L_2	11.28	10.21	10.62	10.56	9.76	9.04	8.45

238 4.2 Exploring Boundary Cases for Trajectory Prediction

239 **Methods** We study the purely (un)supervised generation of trajectories by maximizing the respec-
 240 tive lower bounds on NBA. For the fully-supervised variant, the weak labels described in Section
 241 4.1 can be naturally integrated into the overall scheme via treating them as semantic concepts $y_t^{(a)}$.
 242 These target areas naturally vary over time so that the sequential label $\mathbf{y}_{\leq T}$ denotes always the next
 243 desired position for every involved agent at time t . Following related work [40, 64], we obtain 90
 244 different areas/class labels. The underlying data distribution is estimated by maximizing Eqn (3) via
 245 setting $\lambda_0 = \lambda_1 = 0$ in Eqn (6). The supervised variant of our contribution is denoted by *S-MAT*.

246 Although done in [40, 64], direct comparison of unsupervised and supervised generative models
 247 is inappropriate inasmuch as the latter use ground-truth labels over the entire observation period
 248 (which include future information) at prediction time. Thus, we additionally benchmark against a
 249 fully unsupervised instantiation of our proposed framework lacking ground-truth information for
 250 structuring the discrete latent subspace. To encourage the model to learn semantically meaningful
 251 concepts describing distinct movement patterns, we parametrize the decoder dependent on the
 252 predicted agent class $\tilde{y}_t^{(a)}$. In this context, $\tilde{y}_t^{(a)}$ may be best thought of dynamically assigned agent
 253 roles that accurately explain the observed trajectories. We maximize Eqn (4) using 3 agent types
 254 and refer to the unsupervised variant as *U-MAT*. Further details for S-MAT and U-MAT are given in
 255 Appendix C.

256 **Metrics** We measure the quality of the learned multimodal distribution using common standard
 257 metrics: the minimum over 20 generated samples of the *average* and *final* L_2 distance between
 258 predicted $\hat{\mathbf{x}}_{\leq T}^{(a)}$ and observed positions $\mathbf{x}_{\leq T}^{(a)}$.

259 **Results** Tables 1 and 2 show the results. The proposed approaches clearly outperform their peers
 260 and realize the lowest average L_2 and final L_2 distances. The *MAT* variants effectively accommodate
 261 not only mutual influences across agents, but also discrete generative factors, leading to better
 262 approximations of the underlying multi-modal distributions. Unsurprisingly, Table 1 also shows that
 263 defending players exhibit more structure and are easier to predict than offensive players. The results
 264 of an ablation study are contained in Appendix E. The right part of Fig. 2 displays exemplary rollouts
 265 that underline models’ ability to capture highly complex changes in movement directions.

266 4.3 Semi-Supervised Classification

267 We now extend the experimental protocol by targeting accurate label discovery in semi-supervised
 268 scenarios. We refer to the full model simply as *MAT* where $\lambda_0, \lambda_1 > 0$ and incorporate the fully su-
 269 pervised *S-MAT* ($\lambda_0 = 0$) as an additional baseline. We also compare to an RNN-based classification
 270 network that addresses inter-agent dependencies by inferring hidden states of agents by the GNN

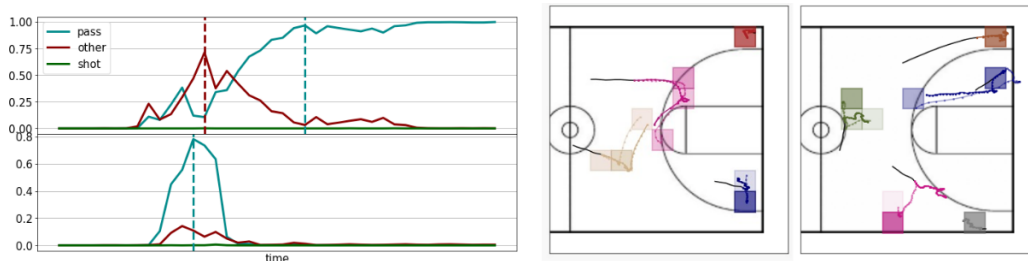


Figure 2: *Left*: Prediction probabilities for different events on soccer data. The dashed lines correspond to the human annotations. *Right*: Exemplary rollouts for offensive players on basketball data. The figure shows observed (light color) and generated player trajectories (intense color) as well as label information indicated by colored boxes, where the color intensity corresponds to the ground-truth frequency of the (weakly obtained) location-based labels.

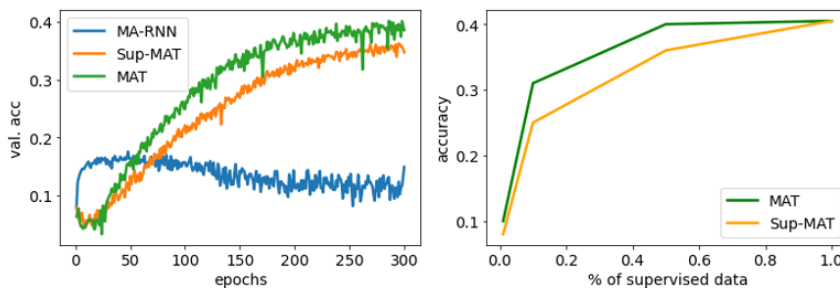


Figure 3: Semi-supervised results on NBA. *Right*: Progression of accuracy during a training run. *Left*: Test accuracies for varying amount of labeled data.

271 architecture described in Section 3.5. A final softmax layer is used for classification and the model is
 272 optimized by minimizing the negative log-likelihood. We refer to this baseline as *MA-RNN*.

273 Figure 3 shows the performance of the models on NBA when using 50% labeled and 50% unlabeled
 274 examples (left side) and when using varying portions of labeled data (right side), while the supervised
 275 baselines have access only to the labeled part. The semi-supervised *MAT* clearly emerges as the
 276 strongest classifier as it converges to the lowest generalization error during a training run and the
 277 performance benefit increases via reduction of labeled data. Comparing the numbers to its supervised
 278 peer (*S-MAT*) highlights the value of incorporating unlabeled data that is effectively utilized by *MAT*
 279 to inform classification decisions. Comparing the fully supervised approaches *S-MAT* and *MA-RNN*,
 280 the importance of an effective regularization mechanism becomes obvious; *MA-RNN* quickly overfits
 281 while the lower-bound in *S-MAT* acts as a regularizer and guarantees gradually increasing predictive
 282 accuracies on validation sets.

283 4.4 Combining Generation and Classification in Semi-Supervised Settings

284 In our last set of experiments, we study the combination of trajectory generation and classification
 285 and resort to tracking data from elite soccer. The reason for favoring this data over NBA is that
 286 the soccer tracking data is manually annotated by experts from the data provider. Since the data
 287 requires human annotations for further processing, a successful semi-supervised application bears
 288 a great deal of practical relevance by reducing the expensive labeling effort. We focus on labels
 289 $\mathcal{Y} = \{pass, other\ ball\ action, shot, none\}$, where class *none* denotes the absence of all other labels in
 290 a frame and use 20% annotated data. Every label is propagated to the previous five and the subsequent
 291 30 frames. We generate a balanced training set where half of the segments carry label *none* and study
 292 the combined task of trajectory generation and classification.

293 This set of experiments mainly serves the purpose to quantitatively validate the benefits of incorporat-
 294 ing expensive human annotations as a generation guide or for structuring the discrete posterior. This

Table 3: Results on soccer data. *Left*: Semi-supervised data generation for a prediction interval of 10 timesteps with an observation period of 40 timesteps (errors in meters). *Right*: Semi-supervised event detection.

NAME	AVG L_2	FINAL L_2	NAME	ACCURACY	F1
VRNN	2.91	7.10	MA-RNN-DIAG	0.74	0.80
VRNN + GNN	2.87	6.98	MA-RNN	0.85	0.88
GVRNN	2.89	6.97	S-MAT	0.82	0.85
S-MAT	2.87	6.91	MAT	0.88	0.92
MAT-DIAG	2.84	6.90			
MAT	2.81	6.88			

295 approach comes with little additional costs since the method operates semi-supervised and therefore
 296 requires only a small subset of multi-agent segments to be annotated. Given the higher sampling rate
 297 and less involved plays in soccer dynamics, we also intend to evaluate the extent to which interaction
 298 modules achieve empirical benefits. We design our baselines accordingly.

299 **Generation** For trajectory prediction, we compare *MAT* to *VRNN* [7] that implicitly assumes
 300 independence across agent dimensions, an interactive version of thereof with GNNs on the hidden
 301 states (*VRNN+GNN*), *GVRNN* [59], *S-MAT* with long-term goals (introduced in Section 4.2, and
 302 a diagonal *MAT* version, *MAT-Diag*. Note that the *VRNN* constitutes a strong competitor for the
 303 task at-hand, as it effectively corresponds to a diagonal *GVRNN*, which has been shown to produce
 304 competitive results on soccer data [59].

305 The left part of Table 3 shows the results. Although the results are generally more similar compared
 306 to the basketball experiments, the semi-supervised *MAT* generates trajectories that are closest in both
 307 metrics to the observed ones. We hypothesize that decreasing the frequency of the data and extracting
 308 multi-agent segments that merely consist of highly interactive plays yields an significant increase in
 309 performance gap.

310 **Classification** *MAT* learns a model over the label space \mathcal{Y} simultaneously to the generative model.
 311 To evaluate classification performance, we turn the output of the corresponding softmax for a given
 312 frame into a class label whenever it exceeded a predefined threshold. We compare the prediction
 313 performance again to *MA-RNN*, a diagonal variant *MA-RNN-Diag*, and *S-MAT*.

314 The resulting accuracies and F1 scores are summarized in Table 3. Once more, *MAT* clearly beats the
 315 baselines and stands out by the highest accuracy and F1 scores. The result impressively demonstrates
 316 the benefit of including unlabeled data for the task. Surprisingly, *MA-RNN* performs better than
 317 *S-MAT*. This finding is in contrast to the results on NBA shown in Figure 3. However, soccer is
 318 played on a much larger pitch and movements are likely more linear compared to basketball, which
 319 leads to less multimodality of the distributions involved and a consequently decreased benefit of
 320 variational methods. Together with a significantly reduced label space and sufficient amounts of data
 321 this may lead to a simpler learning task where *MA-RNN* is less prone to overfitting.

322 The left part of Fig. 2 shows exemplary prediction probabilities for two segments and three possible
 323 events. In both segments, the algorithm is highly confident no shot-action will occur in the near future.
 324 The segment on top contains an event of the class *other ball action* followed by an event of class *pass*.
 325 Both are correctly identified by *MAT*. The former is clearly indicated by a peak in the corresponding
 326 probability chart which then decreases to give rise to the following *pass* action. The segment on the
 327 bottom shows solely a single *pass* event that is also clearly identifiable by a peak at the correct frame.

328 5 Conclusion

329 We presented semi-supervised variational autoencoders for spatiotemporal multi-agent scenarios.
 330 The proposed approach effectively combined ideas from semi-supervised variational autoencoders,
 331 variational recurrent neural networks, and graph neural networks. Empirically, our approach clearly
 332 outperformed the state-of-the-art in sequential generation and (semi-supervised) classification tasks in
 333 all experiments. The performance underlines the benefit of including unlabeled data in spatiotemporal
 334 problems where labeled sequences are either scarce or assembled from weak makeshift signals.

335 References

- 336 [1] Javad Amirian, Jean-Bernard Hayet, and Julien Pettré. Social ways: Learning multi-modal
337 distributions of pedestrian trajectories with gans. In *Proceedings of the IEEE/CVF Conference*
338 *on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019.
- 339 [2] Peter W Battaglia, Razvan Pascanu, Matthew Lai, Danilo Rezende, and Koray Kavukcuoglu.
340 Interaction networks for learning about objects, relations and physics. *arXiv preprint*
341 *arXiv:1612.00222*, 2016.
- 342 [3] Justin Bayer and Christian Osendorfer. Learning stochastic recurrent networks. *arXiv preprint*
343 *arXiv:1411.7610*, 2014.
- 344 [4] Samuel R Bowman, Luke Vilnis, Oriol Vinyals, Andrew M Dai, Rafal Jozefowicz, and Samy
345 Bengio. Generating sentences from a continuous space. *arXiv preprint arXiv:1511.06349*,
346 2015.
- 347 [5] Zhengdao Chen, Soledad Villar, Lei Chen, and Joan Bruna. On the equivalence between graph
348 isomorphism testing and function approximation with gnns. *Advances in neural information*
349 *processing systems*, 32, 2019.
- 350 [6] Junyoung Chung, Caglar Gulcehre, KyungHyun Cho, and Yoshua Bengio. Empirical evaluation
351 of gated recurrent neural networks on sequence modeling. *arXiv preprint arXiv:1412.3555*,
352 2014.
- 353 [7] Junyoung Chung, Kyle Kastner, Laurent Dinh, Kratarth Goel, Aaron C Courville, and Yoshua
354 Bengio. A recurrent latent variable model for sequential data. *Advances in neural information*
355 *processing systems*, 28:2980–2988, 2015.
- 356 [8] Andrew M Dai and Quoc V Le. Semi-supervised sequence learning. *Advances in neural*
357 *information processing systems*, 28:3079–3087, 2015.
- 358 [9] Otto Fabius and Joost R Van Amersfoort. Variational recurrent auto-encoders. *arXiv preprint*
359 *arXiv:1412.6581*, 2014.
- 360 [10] Panna Felsen, Patrick Lucey, and Sujoy Ganguly. Where will they go? predicting fine-grained
361 adversarial multi-agent motion using conditional variational autoencoders. In *Proceedings of*
362 *the European conference on computer vision (ECCV)*, pages 732–747, 2018.
- 363 [11] Marco Fraccaro, Søren Kaae Sønderby, Ulrich Paquet, and Ole Winther. Sequential neural
364 models with stochastic layers. *arXiv preprint arXiv:1605.07571*, 2016.
- 365 [12] Justin Gilmer, Samuel S Schoenholz, Patrick F Riley, Oriol Vinyals, and George E Dahl. Neural
366 message passing for quantum chemistry. In *International conference on machine learning*,
367 pages 1263–1272. PMLR, 2017.
- 368 [13] Roger Girgis, Florian Golemo, Felipe Codevilla, Jim Aldon D’Souza, Martin Weiss,
369 Samira Ebrahimi Kahou, Felix Heide, and Christopher Pal. Autobots: Latent variable se-
370 quential set transformers. *arXiv preprint arXiv:2104.00563*, 2021.
- 371 [14] Dong Gong, Frederic Z Zhang, Javen Qinfeng Shi, and Anton van den Hengel. Memory-
372 augmented dynamic neural relational inference. In *Proceedings of the IEEE/CVF International*
373 *Conference on Computer Vision*, pages 11843–11852, 2021.
- 374 [15] Anirudh Goyal, Alessandro Sordoni, Marc-Alexandre Côté, Nan Rosemary Ke, and Yoshua
375 Bengio. Z-forcing: Training stochastic recurrent networks. *arXiv preprint arXiv:1711.05411*,
376 2017.
- 377 [16] Colin Graber and Alexander Schwing. Dynamic neural relational inference for forecasting
378 trajectories. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
379 *Recognition Workshops*, pages 1018–1019, 2020.
- 380 [17] David Ha and Douglas Eck. A neural representation of sketch drawings. *arXiv preprint*
381 *arXiv:1704.03477*, 2017.
- 382 [18] William L Hamilton, Rex Ying, and Jure Leskovec. Inductive representation learning on large
383 graphs. In *Proceedings of the 31st International Conference on Neural Information Processing*
384 *Systems*, pages 1025–1035, 2017.
- 385 [19] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8):
386 1735–1780, 1997.

- 387 [20] Yedid Hoshen. Vain: Attentional multi-agent predictive modeling. *arXiv preprint*
388 *arXiv:1706.06122*, 2017.
- 389 [21] Yingfan Huang, HuiKun Bi, Zhaoxin Li, Tianlu Mao, and Zhaoqi Wang. Stgat: Modeling
390 spatial-temporal interactions for human trajectory prediction. In *Proceedings of the IEEE/CVF*
391 *International Conference on Computer Vision*, pages 6272–6281, 2019.
- 392 [22] Franziska Hübl, Sreten Cvetojevic, Hartwig Hochmair, and Gernot Paulus. Analyzing refugee
393 migration patterns using geo-tagged tweets. *ISPRS International Journal of Geo-Information*, 6
394 (10):302, 2017.
- 395 [23] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax.
396 *arXiv preprint arXiv:1611.01144*, 2016.
- 397 [24] Tom Joy, Sebastian Schmon, Philip Torr, N Siddharth, and Tom Rainforth. Capturing label
398 characteristics in vaes. In *International Conference on Learning Representations*, 2020.
- 399 [25] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
400 *arXiv:1412.6980*, 2014.
- 401 [26] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
402 *arXiv:1312.6114*, 2013.
- 403 [27] Diederik P Kingma, Shakir Mohamed, Danilo Jimenez Rezende, and Max Welling. Semi-
404 supervised learning with deep generative models. In *Advances in neural information processing*
405 *systems*, pages 3581–3589, 2014.
- 406 [28] Thomas Kipf, Ethan Fetaya, Kuan-Chieh Wang, Max Welling, and Richard Zemel. Neural
407 relational inference for interacting systems. In *International Conference on Machine Learning*,
408 pages 2688–2697. PMLR, 2018.
- 409 [29] Thomas N Kipf and Max Welling. Semi-supervised classification with graph convolutional
410 networks. *arXiv preprint arXiv:1609.02907*, 2016.
- 411 [30] John Krumm and Eric Horvitz. Predestination: Inferring destinations from partial trajectories.
412 In *International Conference on Ubiquitous Computing*, pages 243–260. Springer, 2006.
- 413 [31] Jiachen Li, Fan Yang, Masayoshi Tomizuka, and Chiho Choi. Evolvegraph: Multi-agent
414 trajectory prediction with dynamic relational reasoning. *arXiv preprint arXiv:2003.13924*,
415 2020.
- 416 [32] Jiachen Li, Fan Yang, Hengbo Ma, Srikanth Malla, Masayoshi Tomizuka, and Chiho Choi.
417 Rain: Reinforced hybrid attention inference network for motion forecasting. In *Proceedings of*
418 *the IEEE/CVF International Conference on Computer Vision*, pages 16096–16106, 2021.
- 419 [33] Longyuan Li, Jian Yao, Li Wenliang, Tong He, Tianjun Xiao, Junchi Yan, David Wipf, and
420 Zheng Zhang. Grin: Generative relation and intention network for multi-agent trajectory
421 prediction. *Advances in Neural Information Processing Systems*, 34, 2021.
- 422 [34] Lars Maaløe, Casper Kaae Sønderby, Søren Kaae Sønderby, and Ole Winther. Auxiliary deep
423 generative models. In *International conference on machine learning*, pages 1445–1453. PMLR,
424 2016.
- 425 [35] Chris J Maddison, Andriy Mnih, and Yee Whye Teh. The concrete distribution: A continuous
426 relaxation of discrete random variables. *arXiv preprint arXiv:1611.00712*, 2016.
- 427 [36] Osama Makansi, Julius von Kügelgen, Francesco Locatello, Peter Gehler, Dominik Janzing,
428 Thomas Brox, and Bernhard Schölkopf. You mostly walk alone: Analyzing feature attribution
429 in trajectory prediction. *arXiv preprint arXiv:2110.05304*, 2021.
- 430 [37] Karttikeya Mangalam, Harshayu Girase, Shreyas Agarwal, Kuan-Hui Lee, Ehsan Adeli, Jitendra
431 Malik, and Adrien Gaidon. It is not the journey but the destination: Endpoint conditioned
432 trajectory prediction. In *European Conference on Computer Vision*, pages 759–776. Springer,
433 2020.
- 434 [38] Karttikeya Mangalam, Yang An, Harshayu Girase, and Jitendra Malik. From goals, waypoints &
435 paths to long term human trajectory forecasting. In *Proceedings of the IEEE/CVF International*
436 *Conference on Computer Vision*, pages 15233–15242, 2021.
- 437 [39] Patrick L McDermott, Christopher K Wikle, and Joshua Millspaugh. Hierarchical nonlinear
438 spatio-temporal agent-based models for collective animal movement. *Journal of Agricultural,*
439 *Biological and Environmental Statistics*, 22(3):294–312, 2017.

- 440 [40] Alessio Monti, Alessia Bertugli, Simone Calderara, and Rita Cucchiara. Dag-net: Double
441 attentive graph neural network for trajectory forecasting. In *2020 25th International Conference*
442 *on Pattern Recognition (ICPR)*, pages 2551–2558. IEEE, 2021.
- 443 [41] Siddharth Narayanaswamy, Timothy Paige, Jan-Willem Van de Meent, Alban Desmaison, Noah
444 Goodman, Pushmeet Kohli, Frank Wood, and Philip Torr. Learning disentangled representations
445 with semi-supervised deep generative models. 2018.
- 446 [42] Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito,
447 Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in
448 pytorch. 2017.
- 449 [43] Danilo Jimenez Rezende, Shakir Mohamed, and Daan Wierstra. Stochastic backpropagation
450 and approximate inference in deep generative models. In *International conference on machine*
451 *learning*, pages 1278–1286. PMLR, 2014.
- 452 [44] Alexandre Robicquet, Amir Sadeghian, Alexandre Alahi, and Silvio Savarese. Learning social
453 etiquette: Human trajectory understanding in crowded scenes. In *European conference on*
454 *computer vision*, pages 549–565. Springer, 2016.
- 455 [45] Andrey Rudenko, Luigi Palmieri, Michael Herman, Kris M Kitani, Darius M Gavrilă, and Kai O
456 Arras. Human motion trajectory prediction: A survey. *The International Journal of Robotics*
457 *Research*, 39(8):895–935, 2020.
- 458 [46] Yannick Rudolph and Ulf Brefeld. Modeling conditional dependencies in multiagent trajectories.
459 In *International Conference on Artificial Intelligence and Statistics*, pages 10518–10533. PMLR,
460 2022.
- 461 [47] Amir Sadeghian, Vineet Kosaraju, Agrim Gupta, Silvio Savarese, and A Alahi. Trajnet: Towards
462 a benchmark for human trajectory prediction. *arXiv preprint*, 2018.
- 463 [48] Amir Sadeghian, Vineet Kosaraju, Ali Sadeghian, Noriaki Hirose, Hamid Reza Tofighi, and
464 Silvio Savarese. Sophie: An attentive gan for predicting paths compliant to social and physical
465 constraints. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern*
466 *Recognition*, pages 1349–1358, 2019.
- 467 [49] Tim Salzmann, Boris Ivanovic, Punarjay Chakravarty, and Marco Pavone. Trajectron++:
468 Dynamically-feasible trajectory forecasting with heterogeneous data. In *Computer Vision–*
469 *ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part*
470 *XVIII 16*, pages 683–700. Springer, 2020.
- 471 [50] Franco Scarselli, Marco Gori, Ah Chung Tsoi, Markus Hagenbuchner, and Gabriele Monfardini.
472 The graph neural network model. *IEEE transactions on neural networks*, 20(1):61–80, 2008.
- 473 [51] Kihyuk Sohn, Honglak Lee, and Xinchen Yan. Learning structured output representation using
474 deep conditional generative models. *Advances in neural information processing systems*, 28:
475 3483–3491, 2015.
- 476 [52] Chen Sun, Per Karlsson, Jiajun Wu, Joshua B Tenenbaum, and Kevin Murphy. Stochastic pre-
477 diction of multi-agent interactions from partial observations. *arXiv preprint arXiv:1902.09641*,
478 2019.
- 479 [53] Karl Tuyls, Shayegan Omidshafiei, Paul Muller, Zhe Wang, Jerome Connor, Daniel Hennes, Ian
480 Graham, William Spearman, Tim Waskett, Dafydd Steel, et al. Game plan: What ai can do for
481 football, and what football can do for ai. *Journal of Artificial Intelligence Research*, 71:41–88,
482 2021.
- 483 [54] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez,
484 Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information*
485 *processing systems*, 30, 2017.
- 486 [55] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Lio, and Yoshua
487 Bengio. Graph attention networks. *arXiv preprint arXiv:1710.10903*, 2017.
- 488 [56] Bing Xu, Naiyan Wang, Tianqi Chen, and Mu Li. Empirical evaluation of rectified activations
489 in convolutional network. *arXiv preprint arXiv:1505.00853*, 2015.
- 490 [57] Keyulu Xu, Chengtao Li, Yonglong Tian, Tomohiro Sonobe, Ken-ichi Kawarabayashi, and
491 Stefanie Jegelka. Representation learning on graphs with jumping knowledge networks. In
492 *International Conference on Machine Learning*, pages 5453–5462. PMLR, 2018.

- 493 [58] Andy Yuan Xue, Rui Zhang, Yu Zheng, Xing Xie, Jin Huang, and Zhenghua Xu. Destination
494 prediction by sub-trajectory synthesis and privacy protection against such prediction. In *Data*
495 *Engineering (ICDE), 2013 IEEE 29th International Conference on*, pages 254–265. IEEE, 2013.
- 496 [59] Raymond A Yeh, Alexander G Schwing, Jonathan Huang, and Kevin Murphy. Diverse genera-
497 tion for multi-agent sports games. In *Proceedings of the IEEE/CVF Conference on Computer*
498 *Vision and Pattern Recognition*, pages 4610–4619, 2019.
- 499 [60] Jiaxuan You, Zhitao Ying, and Jure Leskovec. Design space for graph neural networks. *Advances*
500 *in Neural Information Processing Systems*, 33, 2020.
- 501 [61] Cunjun Yu, Xiao Ma, Jiawei Ren, Haiyu Zhao, and Shuai Yi. Spatio-temporal graph transformer
502 networks for pedestrian trajectory prediction. In *European Conference on Computer Vision*,
503 pages 507–523. Springer, 2020.
- 504 [62] Ye Yuan, Xinshuo Weng, Yanglan Ou, and Kris M Kitani. Agentformer: Agent-aware transform-
505 ers for socio-temporal multi-agent forecasting. In *Proceedings of the IEEE/CVF International*
506 *Conference on Computer Vision*, pages 9813–9823, 2021.
- 507 [63] Manzil Zaheer, Satwik Kottur, Siamak Ravanbakhsh, Barnabas Poczos, Ruslan Salakhutdinov,
508 and Alexander Smola. Deep sets. *arXiv preprint arXiv:1703.06114*, 2017.
- 509 [64] Eric Zhan, Stephan Zheng, Yisong Yue, Long Sha, and Patrick Lucey. Generating multi-agent
510 trajectories using programmatic weak supervision. *arXiv preprint arXiv:1803.07612*, 2018.
- 511 [65] Stephan Zheng, Yisong Yue, and Patrick Lucey. Generating long-term trajectories using deep
512 hierarchical networks. *arXiv preprint arXiv:1706.07138*, 2017.

513 Checklist

- 514 1. For all authors...
- 515 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
516 contributions and scope? [Yes]
- 517 (b) Did you describe the limitations of your work? [Yes] Since this work constitutes the
518 first semi-supervised approach for spatiotemporal multi-agent problems, it is difficult
519 to discuss technical limitations at this stage. Application-wise, there are clear limitations
520 e.g., by using weak makeshift labels. We discuss these issues across the paper (e.g.,
521 Section 4.)
- 522 (c) Did you discuss any potential negative societal impacts of your work? [No] We have
523 a clear focus on team sports. Though dual use of the proposed technique is certainly
524 possible (e.g., military domains), it may be a bit far fetched as we have no expertise in
525 these domains and would only be able to deliver uneducated guesses about whether
526 this is realistic or not.
- 527 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
528 them? [Yes]
- 529 2. If you are including theoretical results...
- 530 (a) Did you state the full set of assumptions of all theoretical results? [Yes] See Section
531 4.1. – 4.3.
- 532 (b) Did you include complete proofs of all theoretical results? [Yes] All proofs are
533 contained in the Appendix.
- 534 3. If you ran experiments...
- 535 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
536 mental results (either in the supplemental material or as a URL)? [Yes] The NBA data
537 set is publicly available and we provide a URL in the paper. The soccer data is private
538 and must not be shared with third parties. We will release the code and include a URL
539 in the paper as well after the paper has been accepted.
- 540 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
541 were chosen)? [Yes] See Section 4, more details are provided in the Appendix.
- 542 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
543 ments multiple times)? [No]

- 544 (d) Did you include the total amount of compute and the type of resources used (e.g., type
545 of GPUs, internal cluster, or cloud provider)? [Yes] See Appendix
- 546 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 547 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 548 (b) Did you mention the license of the assets? [Yes] See supplemental material
- 549 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
- 550 (d) Did you discuss whether and how consent was obtained from people whose data you're
551 using/curating? [Yes] The NBA data is publicly available and free, for soccer we have
552 the permission of the data provider.
- 553 (e) Did you discuss whether the data you are using/curating contains personally identifiable
554 information or offensive content? [Yes] The soccer and basketball players are probably
555 identifiable in the data, but they/their contract/their clubs agreed in recording the data.
556 The soccer data is private and cannot be shared.
- 557 5. If you used crowdsourcing or conducted research with human subjects...
- 558 (a) Did you include the full text of instructions given to participants and screenshots, if
559 applicable? [N/A]
- 560 (b) Did you describe any potential participant risks, with links to Institutional Review
561 Board (IRB) approvals, if applicable? [N/A]
- 562 (c) Did you include the estimated hourly wage paid to participants and the total amount
563 spent on participant compensation? [N/A]

564 **A Proofs of the Theorems**

565 We begin with Theorem 2 for simplicity and sketch the proof of Theorem 1 afterwards.

566 **Theorem 2.** Let $\mathcal{H}(\beta)$ be the entropy of quantity β . A lower bound on $\log p_\theta(\mathbf{x}_{\leq T})$ in Eqn (1) is
 567 given by

$$\log p_\theta(\mathbf{x}_{\leq T}) \geq \sum_t \left(\mathcal{H}(q_\phi(\tilde{\mathbf{y}}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})) - \mathbb{E}_{q_\phi(\tilde{\mathbf{y}}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \tilde{\mathbf{y}}_{< t})} [\mathcal{L}(\mathbf{x}_t, \tilde{\mathbf{y}}_t)] \right) \equiv \sum_t -\mathcal{U}(\mathbf{x}_t)$$

568 *Proof.* For unlabeled training instances, defining the marginal likelihood according to the generative
 569 structure introduced in Section 3.3 yields

$$\begin{aligned} p_\theta(\mathbf{x}_{\leq T}) &= \int_{\mathbf{z}_{\leq T}} \int_{\mathbf{y}_{\leq T}} p_\theta(\mathbf{x}_{\leq T}, \mathbf{z}_{\leq T}, \mathbf{y}_{\leq T}) d\mathbf{y}_{\leq T} d\mathbf{z}_{\leq T} \\ &= \int_{\mathbf{z}_{\leq T}} \sum_{\mathbf{y}_{\leq T}} \prod_{t=1}^T p_\theta(\mathbf{x}_t | \mathbf{x}_{< t}, \mathbf{z}_{\leq t}, \mathbf{y}_{\leq t}) p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) d\mathbf{z}_{\leq T}. \end{aligned}$$

570 To derive the lower bound on the log likelihood, we incorporate the variational information into the
 571 above definition and apply Jensen's inequality:

$$\begin{aligned} &\int \int q_\phi(\mathbf{z}_{\leq T}, \mathbf{y}_{\leq T} | \mathbf{x}_{\leq T}) \log \prod_{t=1}^T \frac{p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})}{q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq t}) q_\phi(\mathbf{y}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{\leq t}, \mathbf{y}_{< t})} d\mathbf{z}_{\leq T} d\mathbf{y}_{\leq T} \\ &= \sum_{t=1}^T \int \int q_\phi(\mathbf{z}_{\leq T}, \mathbf{y}_{\leq T} | \mathbf{x}_{\leq T}) \log \frac{p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})}{q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq t}) q_\phi(\mathbf{y}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{\leq t}, \mathbf{y}_{< t})} d\mathbf{z}_{\leq T} d\mathbf{y}_{\leq T} \\ &= \sum_{t=1}^T \int \int q_\phi(\mathbf{z}_{\leq t}, \mathbf{y}_{\leq t} | \mathbf{x}_{\leq t}) \log \frac{p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})}{q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq t}) q_\phi(\mathbf{y}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{\leq t}, \mathbf{y}_{< t})} d\mathbf{z}_{\leq t} d\mathbf{y}_{\leq t} \\ &= \sum_{t=1}^T \int \int q_\phi(\mathbf{z}_{< t}, \mathbf{y}_{< t} | \mathbf{x}_{< t}) \left(-\mathbb{E}_{q_\phi(\mathbf{y}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})} [\mathcal{KL}[q_\phi(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq t} \parallel p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})]] \right. \\ &\quad \left. - \mathcal{KL}[q_\phi(\mathbf{y}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) \parallel p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})] \right) d\mathbf{z}_{< t} d\mathbf{y}_{< t} \\ &= \mathbb{E}_{q_\phi(\mathbf{z}_{\leq T}, \mathbf{y}_{\leq T} | \mathbf{x}_{\leq T})} \left[\sum_{t=1}^T -\mathcal{KL}[q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq T}) \parallel p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})] \right. \\ &\quad \left. - \mathcal{KL}[q_\phi(\mathbf{y}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) \parallel p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})] \right]. \end{aligned}$$

572 Thus, we can write

$$\begin{aligned} \log p_\theta(\mathbf{x}_{\leq T}) &\geq \mathbb{E}_{q_\phi(\mathbf{z}_{\leq T}, \mathbf{y}_{\leq T} | \mathbf{x}_{\leq T})} \left[\sum_{t=1}^T \log p_\theta(\mathbf{x}_t | \mathbf{x}_{< t}, \mathbf{z}_{\leq t}, \mathbf{y}_{\leq t}) \right. \\ &\quad \left. - \mathcal{KL}[q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq t}) \parallel p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})] \right. \\ &\quad \left. - \mathcal{KL}[q_\phi(\mathbf{y}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) \parallel p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})] \right], \end{aligned}$$

573 which is identical to Eqn 3. □

574 **Theorem 1.** A lower bound on $\log p_\theta(\mathbf{x}_{\leq T}, \mathbf{y}_{\leq T})$ in Eqn (1) is given by

$$\begin{aligned} \log p_\theta(\mathbf{x}_{\leq T}, \mathbf{y}_{\leq T}) &\geq \sum_t \log p_\theta(\mathbf{y}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t}) + \mathbb{E}_{q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq t})} \left[\log p_\theta(\mathbf{x}_t | \mathbf{x}_{< t}, \mathbf{z}_{\leq t}, \mathbf{y}_{\leq t}) \right. \\ &\quad \left. - \mathcal{KL}[q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, \mathbf{y}_{\leq t}) \parallel p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, \mathbf{y}_{< t})] \right] \equiv \sum_{t=1}^T -\mathcal{L}(\mathbf{x}_t, \mathbf{y}_t). \end{aligned}$$

Proof.

$$\begin{aligned} \log p_\theta(\mathbf{x}_{\leq T}, y_{\leq T}) &= \log \int_{\mathbf{z}_{\leq T}} \prod_{t=1}^T \frac{q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, y_{\leq t})}{q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, y_{\leq t})} p_\theta(\mathbf{x}_t | \mathbf{x}_{< t}, \mathbf{z}_{\leq t}, y_{\leq t}) p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, y_{< t}) \\ &\geq \mathbb{E}_{q_\phi(\mathbf{z}_{\leq T} | \mathbf{x}_{\leq T}, y_{\leq T})} \left[\sum_{t=1}^T \log p_\theta(\mathbf{x}_t | \mathbf{x}_{< t}, \mathbf{z}_{\leq t}, y_{\leq t}) \right. \\ &\quad \left. - \mathcal{KL}[q_\phi(\mathbf{z}_t | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, y_{\leq t}) \parallel p_\theta(\mathbf{z}_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, y_{< t})] + \log p_\theta(y_t | \mathbf{x}_{< t}, \mathbf{z}_{< t}, y_{< t}) \right]. \end{aligned}$$

575

□

576 B Implementation Details

577 All models are implemented using PyTorch [42]. For our experiments, we
578 define the model components as follows. The encoders are modeled as

579 $q_\phi(\tilde{y}^{(a)} | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, y_{< t}) = \text{Cat}(y_t^{(a)} | f_{enc}^{(a)}([\mathbf{o}_{t-1}^{(a)}, \varphi^x(\mathbf{x}_t^{(a)})])$ and $q_\phi(\mathbf{z}_t^{(a)} | \mathbf{x}_{\leq t}, \mathbf{z}_{< t}, y_{\leq t}) =$
580 $\mathcal{N}(\mathbf{z}_t^{(a)} | f_{enc}^{(a)}([\mathbf{o}_{t-1}^{(a)}, \varphi^x(\mathbf{x}_t^{(a)})])$ for uncovering discrete and continuous latent information,
581 respectively. The prior distributions are computed analogously omitting input \mathbf{x}_t . Agent movements
582 are represented as $p_\theta(\mathbf{x}_t^{(a)} | \mathbf{x}_{< t}, \mathbf{z}_{\leq t}, y_{\leq t}) = \mathcal{N}(\mathbf{z}_t^{(a)} | f_{dec}^{(a)}([\mathbf{o}_{t-1}^{(a)}, \varphi^x(\mathbf{x}_t^{(a)})])$.

583 All functions f and feature extractors φ are two-layer MLPs with LeakyRelu [56] activations. We use
584 2-layer GRU networks [6] for recurrence. Moreover, we define the graph network using 3 GAT layers
585 implemented according to [55] and the graph structure wrt. the $k = 2$ and $k = 5$ spatially closest
586 agents for basketball and soccer, respectively. We implement skip-connections as described in [57]
587 and use concatenation followed by a linear transformation to aggregate intermediate layer embeddings
588 to the GNN model output. For generative tasks, we use $\lambda_1 = 0$, however, including the auxiliary loss
589 yields negligible deterioration. For classification, we use $\lambda_0 = 1$ and $\alpha = 10 * (1/\%$ of labeled data).

590 The model operates solely on agent velocities (input and output). However, we inject position tuples
591 to the all model components via teacher forcing. We use the gumbel-softmax trick [35, 23] when
592 sampling from categorical distributions. Training was executed on an Nvidia V100 GPU and took
593 about 16 hours to complete 300 epochs while consuming $\sim 15\text{Gb}$ for basketball. All models are
594 optimized using Adam [25] with a learning rate of 0,001 and gradient clipping using a max norm of
595 10.

596 **Baselines** As stated in Section 4.1, we report against the values reported in [40]. For [16, 33, 37],
597 we adapt the source code from their official repositories to our experimental setting⁶ For GVRNN
598 [59], we re-implemented the model according to the descriptions in their paper and designed the
599 overall architecture such that it is comparable in parameter number to our method.

600 **Soccer application** Since the labels denote ball-centric events, we use the output of the ball node
601 for loss computation and evaluation. The $F1$ score is computed as follows. We annotate a multi-
602 agent segment when the derived probability estimates exceed an externally defined threshold value.
603 We obtain TP values (FP values) when the predicted event coincides (disagrees) with the ground
604 truth annotation. FN values are defined by annotated segments that remain undetected. We compute
605 F1-scores for 100 distinct threshold values in the range between 0.5 and 0.98 and only report the
606 maximum F1-score. However, threshold optimization yields only negligible improvement over simply
607 using 0.5.

608 C More Details on the MAT Models

609 In Section 3.1, we define $\mathbf{y}_{< T}$ as arbitrary discrete (possibly latent) behavioral indicators. To validate
610 this general formulation (and our proposed architecture), we vary the specific definition of the variable

⁶For [37], we used https://github.com/crowdbotp/OpenTraj/blob/master/datasets/SDD/estimated_scales.yaml to map between real-world and pixel coordinate. s

611 across experiments. Accordingly, we propose different framework instantiations that are described in
 612 more detail in this section.

613 **S-MAT** First, we observe in Section 4.2 that existing SOTA trajectory prediction approaches
 614 (Section 4.1) use heuristically generated labels for trajectory prediction that encode agents’ intents
 615 or goals over a discretized position space. We then note that our formalization allows to naturally
 616 integrate such long-term goals into the overall scheme via treating them as discrete semantic concepts.
 617 Since these labels are generated heuristically based on the trajectory input prior to model training,
 618 this model instantiation is fully-supervised and is referred to as *S-MAT*. We refer to [64] for more
 619 details on how to produce the weak labels used for training S-MAT.

620 **U-MAT** However, we found (although done in [64, 40]) that benchmarking against unsupervised
 621 baselines is inappropriate for the reasons described in Section 4.1. Since most SOTA approaches
 622 are unsupervised generative models, we additionally propose a fully-unsupervised instantiation of
 623 our framework (*U-MAT*) employing no external guidance for discrete latent structuring. From the
 624 predicted (latent) categorical distribution q_ϕ , we sample a label value and exploit separate motion
 625 predictors (different decoder parameterizations) p_θ dependent on the realized value. Intuitively, this
 626 encourages the model to learn categories describing fundamentally different movement patterns,
 627 which can be interpreted as dynamic “agent roles”. Thus, the concept of agent roles here is merely
 628 an intuitive explanation and realized via an inductive bias that increases the “scope” of latent
 629 information encoded without utilizing any supervision. To the best of our knowledge, parameterizing
 630 the generation module based on inferred agent categories is novel and could provide valuable insights
 631 for practitioners.

632 D Experiments with Drone Data

Table 4: Results for SDD (observation phase of 8 timesteps and a prediction horizon with 12 timesteps) expressed in real-world coordinates.

NAME	$avgL_2$	FINAL L_2
STGAT	0.58	1.11
SOCIAL-WAYS	0.62	1.16
DAG-NET	0.53	1.04
PECNET	0.67	1.03
S-MAT	0.51	1.03

633 To showcase applicability to scenarios that exhibit variable numbers of agents, we also report results
 634 on the Stanford Drone Data (SDD) [44]. SDD is a collection of videos recorded by drones at eight
 635 locations at Stanford. While pedestrians predominate as interacting agents, cyclists, skateboarders,
 636 cars, buses, and golf carts are also present. We use the TrajNet benchmark [47] of the data, providing
 637 sequential two-dimensional real-world coordinates at a frame rate of 2.5 frames per second. We
 638 follow [40] regarding the data processing strategy⁷. We incorporate *PECNet* [37] for comparison as
 639 it is considered state-of-the-art for SDD [36]. *PECNet* models stochasticity in the final position of the
 640 pedestrians conditioned on the past motion history.

641 To generate weak labels, we move a time window through each trajectory, with the respective end
 642 cells acting as agent targets $y_t^{(a)}$. Static time windows allow us to use the fully-supervised variant for
 643 all comparisons (including fully unsupervised models). The scene at hand is thereby discretized into
 644 960 areas and results are shown in Table 4. We observe the same general pattern as for the basketball
 645 experiments.

Table 5: Ablation study on NBA.

NAME	L_2	FINAL L_2
S-MAT-DIAG	8.93	13.92
S-MAT-FULL	8.87	13.87
S-MAT-GVRNN	9.78	14.27
S-MAT-GVRNN-HIDDEN	9.90	15.20
S-MAT	8.11	12.52

646 E Ablation Study

647 Table 5 validates the proposed architecture by showing results of an ablation study. We test a fully
648 connected graph (*S-MAT-Full*), an independent version with diagonal adjacency matrices (*S-MAT-*
649 *Diag*), a variant that employs GNNs for variational, generative, and prior distribution parameters, but
650 ignores interactive updates for the hidden states (*S-MAT-GRVNN*) similar to [59], and a variant that is
651 identical to the latter but additionally includes a GNN for the hidden state (*S-MAT-GRVNN-Hidden*).
652 The table provides supporting evidence for our design choices: though the GVRNNs are theoretically
653 able to capture intra-timestep dependencies, we observe significant drops in performance compared to
654 the other competitors. The S-MAT-Diag and MAT-Full experiments suggest that our model captures
655 interaction patterns among agents very well; *S-MAT* denotes a valuable contribution to the large body
656 of research that explicitly addresses modeling multi-agent data accurately.

⁷Data and preprocessing can be accessed at <https://github.com/alexmonti19/dagnet/tree/master/datasets>