

---

# Policy Learning with Partial Observation and Mechanical Constraints for Multi-person Modeling

---

**Keisuke Fujii**

Nagoya University, RIKEN  
fujii@i.nagoya-u.ac.jp

**Naoya Takeishi**

RIKEN  
naoya.takeishi@riken.jp

**Yoshinobu Kawahara**

Kyushu University, RIKEN  
kawahara@imi.kyushu-u.ac.jp

**Kazuya Takeda**

Nagoya University  
kazuya.takeda@nagoya-u.jp

## Abstract

Extracting the rules of real-world biological multi-agent behaviors is a current challenge in various scientific and engineering fields. Biological agents generally have limited observation and mechanical constraints; however, most of the conventional data-driven models ignore such assumptions, resulting in lack of biological plausibility and model interpretability for behavioral analyses in biological and cognitive science. Here we propose sequential generative models with partial observation and mechanical constraints, which can visualize whose information the agents utilize and can generate biologically plausible actions. We formulate this as a decentralized multi-agent imitation learning problem, leveraging binary partial observation models with a Gumbel-Softmax reparameterization and policy models based on hierarchical variational recurrent neural networks with physical and biomechanical constraints. We investigate the empirical performances using real-world multi-person motion datasets from basketball and soccer games.

## 1 Introduction

Extracting the rules of biological multi-agent behaviors in complex real world environments from data is a fundamental problem in a variety of scientific and engineering fields. For example, animals, vehicles, pedestrians, and athletes observe other's states and execute their own actions with the body constraints in complex situations. In these processes, they observe other's movements and make decisions, e.g., based on their experiences and knowledge, under spatiotemporal constraints. Pioneering works proposed various ways of rule-based modeling such as in human pedestrians [25] and animal groups [13] using hand-crafted functions (sometimes called social forces). Recent advances in machine learning have enabled data-driven modeling of such behaviors (see Section 4). These problems are formulated as imitation learning, generative adversarial learning, or (simply) sequential learning, often leveraging (deep) sequential generative models.

However, most of these works employ the following three concepts to improve prediction performance: i) they can fully utilize environmental information or that based on pre-determined rules (e.g., [13]); ii) they can optimize communication based on the centralized control; and iii) they ignore the mechanical constraints of the agent's body, resulting in biologically unrealistic behaviors. Although such ideas improve the predictability, interpretable modeling based on biological plausibility [12] is also important for scientific understanding, which is the motivation of this paper. Biological organisms in a real-world generally have limited communication and observation (i.e., they should be considered as partially observable decentralized systems for analyzing their observations). The body constraints such as inertia during motion planning can be described by the principles of robotics

and computational neuroscience [15, 53] (see Section 2.4), to generate smooth biological motions. Therefore, modeling of the agent’s observation and of biologically plausible motions would contribute to the understanding of biological multi-agent behaviors in complex real-world environments, which is one of the challenges in biological and cognitive science.

In this paper, we propose sequential generative models with partial observation and mechanical constraints, which can visualize whose information the agents utilize and can generate biologically plausible behaviors. We formulate this as a decentralized multi-agent imitation learning problem, leveraging partial observation models with a Gumbel-Softmax reparameterization (Sections 2.2 and 3.2) and hierarchical policy models based on deep sequential generative models with stochastic latent variables and mechanical constraints (Sections 2.3, 2.4, and 3.3). A team sport is an example that can be addressed with the above approach. Players observe others’ states [18, 21], while actively and/or passively ignoring less informative agents [20] regardless of the distance, and execute complex actions. We investigate the empirical performance by using real-world ballgame datasets.

In summary, our main contributions are as follows: (1) we propose sequential generative models with interpretable partial observation and mechanical constraints, contributing to the analysis of real-world multi-agent behaviors; (2) as a decentralized multi-agent imitation learning problem, our hierarchical policy models leverage binary partial observation and mechanical constraints, which can be compatible with many existing deep generative models; (3) our approach is validated by visualizing, evaluating, and counterfactually manipulating partial observation and predicted plausible behaviors using real-world basketball and soccer datasets. In the remainder of this paper, we describe the background of our problem and our method in Sections 2 and Section 3, overview related works in Section 4, present experimental results in Section 5, and conclude this paper in Section 6.

## 2 Background

We formulate our problem as a decentralized multi-agent imitation learning problem in partially observable Markov decision process (POMDP) for analyzing their observations, with the assumption that biological multi-agents may behave as decentralized agents in the short-term (e.g., several seconds). Our model includes observation and policy models based on the learning of partial observation and nonlinear dynamics with mechanical constraints, respectively. Here, we introduce imitation learning for decentralized multi-agent in POMDP, observation models with attention mechanism, sequential generative models for the policy modeling, and biomechanical constraints.

### 2.1 Imitation learning for decentralized multi-agent systems in POMDP

Here we consider a multi-agent problem that can be formulated as decentralized POMDP [33, 6, 44]. It is basically defined as a tuple  $(K, S, A, \mathcal{T}, R, O, Z, \gamma)$ , where  $K$  is the fixed number of agents;  $S$  is the set of states  $s$ ;  $A = [A_1, \dots, A_K]$  represents the set of joint action  $\vec{a}$  and  $A_k$  is the set of local action  $a_k$  that agent  $k$  can take; the transition model  $\mathcal{T}(s'|s, \vec{a}) : S \times A \times S \rightarrow [0, 1]$ ;  $R = [R_1, \dots, R_K] : S \times A \rightarrow \mathbb{R}^K$  is the joint reward function;  $O = [O_1, \dots, O_K]$  is the set of joint observation  $\vec{o}$  controlled by the observation function  $Z : S \times A \rightarrow O$ ;  $\gamma \in [0, 1]$  is the discount factor. In on-policy reinforcement learning, the agent learns a policy  $\pi_k : O_k \times A \rightarrow [0, 1]$  that can maximize  $\mathbb{E}[G_i]$  where  $G_i = \sum_{t=1}^T \gamma^t R_k^t$  is the discount return and  $T$  is the time horizon. In a decentralized multi-agent system in complex real-world environments such as team sports, transition model and reward function are generally difficult to design explicitly. Instead, since we can sometimes utilize the demonstrations of expert behaviors, we can formulate our problem as imitation learning.

The goal in imitation learning is to learn a policy  $\pi$  that imitates an expert policy  $\pi_E$  given demonstrations from that expert [50, 48]. In multi-agent imitation learning, we have  $K$  agents to achieve a common goal. Based on the notation of [48], in the case of a fully centralized multi-agent policy for clarity, let  $\vec{\pi}(\vec{o}) := \vec{a}$  denote the joint policy that maps the joint observation  $\vec{o} = [o_1, \dots, o_K]$  into  $K$  actions  $\vec{a} = [a_1, \dots, a_K]$ . Training data  $\mathcal{D}$  consists of multiple demonstrations of  $K$  agents. The decentralized setting decomposes the joint policy  $\vec{\pi} = [\pi_1, \dots, \pi_K]$  into  $K$  policies, sometimes tailored to each specific agent index or role. The loss function is then:  $\mathcal{L}_{imitation} = \sum_{k=1}^K \mathbb{E}_{s_k \sim d_{\pi_k}} [\ell(\pi_k(o_k))]$ , where  $d_{\pi_k}$  denotes the distribution of states experienced by joint policy  $\pi_k$  (again,  $o_k$  is determined by  $s_k$  and  $a_k$ ) and  $\ell$  is the imitation loss defined over the demonstrations. This formulation assumes assignment of appropriate roles for each agent (for details, see Appendix A).

## 2.2 Partial observation models using Gumbel-Softmax reparameterization

A naive approach for neural network-based partial observation is to employ soft and hard attention mechanisms. Attention is widely used in applications including natural language processing [2], computer vision [54], and multi-agent reinforcement learning (MARL) in virtual environments [44, 31, 28]. Soft attention calculates an importance distribution of elements (e.g., agents):  $w_k = \exp(p_k) / \sum_{i=1}^K \exp(p_i)$ , where  $p_k$  is a scalar variable for an agent  $k$ . Soft attention is fully differentiable and thus can be trained with back-propagation. However, it usually assigns non-zero probabilities to unrelated elements, i.e., it cannot directly reduce the number of agents.

Hard attention focuses solely on an important element but is basically non-differentiable due to the selection based on sampling. Among differentiable models with discrete variables, an approach employing Gumbel-Softmax reparameterization [30, 43] can be effectively implemented via a continuous relaxation of a discrete categorical distribution. In summary, given  $K$ -categorical distribution parameters  $p$ , a differentiable  $K$ -dimensional one-hot encoding sample  $G$  from the Gumbel-Softmax distribution can be computed as:  $G(\log p)_k = \exp((\log p_k + \epsilon)/\tau) / \sum_{i=0}^K \exp((\log p_i + \epsilon)/\tau)$ , where  $\epsilon$  are i.i.d. samples from a Gumbel(0, 1) distribution, i.e.,  $\epsilon = -\log(-\log(u))$ ,  $u \sim \mathcal{U}[0, 1]$ .  $\mathcal{U}$  is a uniform distribution and  $\tau$  is the softmax temperature parameter. However, the resulting one-hot vector  $[G(\log p)_1, \dots, G(\log p)_K]$  is insufficient to represent the observation of more than one agent.

## 2.3 Sequential generative model

Agent trajectories or actions in real-world environments have been recently modeled as sequential generative models (for details, see Section 4). In this subsection, we consider single agent modeling for simplicity. Let  $a_{\leq T} = \{a_1, \dots, a_T\}$  denote a sequence of actions of length  $T$ . The goal of sequential generative modeling is to learn the distribution over sequential data  $\mathcal{D}$  consisting of multiple demonstrations. We assume that all sequences have the same length  $T$ , but in general, this does not need to be the case. A common approach is to factorize the joint distribution and then maximize the log-likelihood  $\theta^* = \arg \max_{\theta} \sum_{a_{\leq T} \in \mathcal{D}} \log p_{\theta}(a_{\leq T}) = \arg \max_{\theta} \sum_{a_{\leq T} \in \mathcal{D}} \sum_{t=1}^T \log p_{\theta}(a_t | a_{< t})$ , where  $\theta$  denotes the model’s learnable parameters, such as recurrent neural networks (RNNs).

However, RNNs with simple output distributions often struggle to capture highly variable and structured sequential data (e.g., multimodal behaviors) [59]. Recent work in sequential generative models addressed this issue by injecting stochastic latent variables into the model and optimizing using amortized variational inference to learn the latent variables [11, 17, 22] (see Section 4). Among these methods, a variational RNN (VRNN) [11] has been widely used in base models for multi-agent trajectories [56, 59]. Note that our VRNN-based approach is also compatible with other sequential generative models. A VRNN is essentially a variational autoencoder (VAE) conditioned on the hidden state of an RNN and is trained by maximizing the (sequential) evidence lower-bound (ELBO):

$$\mathbb{E}_{q_{\phi}(z_{\leq T} | a_{\leq T})} \left[ \sum_{t=1}^T \log p_{\theta}(a_t | z_{\leq t}, a_{< t}) - D_{KL}(q_{\phi}(z_t | a_{\leq t}, z_{< t}) || p_{\theta}(z_t | a_{< t}, z_{< t})) \right], \quad (1)$$

where  $z$  is a stochastic latent variable. The first term is the reconstruction term and  $p_{\theta}(a_t | z_{\leq t}, a_{< t})$  is generative model. The second term is the Kullback-Leibler (KL) divergence between the approximate posterior or inference model  $q_{\phi}(z_t | a_{\leq t}, z_{< t})$  and the prior  $p_{\theta}(z_t | a_{< t}, z_{< t})$ . Eq. (1) can be interpreted as the ELBO of VAE summed over each timestep  $t$ . For further details, see Appendix B.

## 2.4 Biomechanical constraints

In robotics and computational neuroscience, the path planning in smooth and flexible biological motions is a classical problem [15, 53]. Most research has focused on individual multi-joint motions. However, in our multi-agent cases, datasets often include only a mass-point for each agent. Thus, we can utilize the principles of biological mass-point motions such as end-point effectors. A minimum-jerk principle [15] is a simple and well-known principle for biological motor planning of voluntary motion in an end-point effector. It minimizes the motor cost:  $C = (1/2) \int_0^T \|(d^3 x / dt^3)\|_2^2 dt$ , where  $x$  is multi-dimensional position vector and  $\|\cdot\|_2$  is the Euclidean norm. We propose a simply applicable biomechanical constraint for learning a sequential model inspired by this principle.

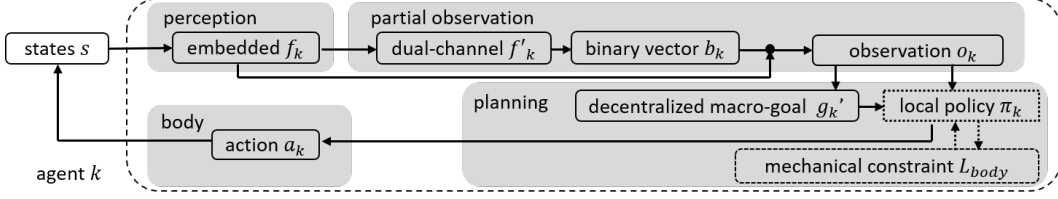


Figure 1: Diagram of our model. For clarity, we omitted the time index  $t$ . The agent  $k$  perceives the state  $s_{t-1}$  and outputs the action  $a_{t,k}$ . The detailed configuration is described in the main text.

### 3 Proposed Method

In this section, we propose sequential generative models for policy models with partial observation and mechanical constraints (Figure 1), which can visualize whose information the agents utilize and can generate plausible behaviors. We leverage binary partial observation models and a hierarchical VRNN with mechanical constraints. We first overview our model, then propose a binary partial observation model and policy models with mechanical constraints, and describe the learning method.

#### 3.1 Overview

As an imitation learning problem for decentralized multi-agent systems in POMDP, we aim to learn a policy  $\vec{\pi} = [\pi_1, \dots, \pi_K]$  that imitates an expert policy  $\vec{\pi}_E = [\pi_{E_1}, \dots, \pi_{E_K}]$  given the expert action sequences of  $K$  agents:  $a_{\leq T} = \{a_{\leq T,1}, \dots, a_{\leq T,K}\}$  under biologically realistic constraints. As shown in Figure 1, the agent perceives the state  $s_{t-1}$  and performs path planning for the action  $a_{t,k}$ . In this paper, since we consider the decentralized model, we train independent models for each agent.

Our model consists of a partial observation, local policy, macro-goal, and mechanical constraints. We first propose a binary vector  $b_{t,k}$  of the partial observation model defined in Section 3.2. The model also utilizes macro-goals [59] for long-term prediction modified for our decentralized and partially-observable problem (see Appendix C). Note that organism behaviors are generally generated by path planning including cognitive and motor factors of which rules are partially unknown [19]. For example, some cognitive factors (often higher-level factors such as tactics) remain unknown. We therefore adopt a nonlinear transition model (see Section 4) for the local policy  $\pi_k$ , specifically based on a VRNN [11], which can generate a trajectory with multiple variations [59, 56, 51]. As previously emphasized, our approach is also applicable to other sequential generative models.

In the previous works (e.g., [59, 56]), the action generated by  $\pi_k$  is usually the agent’s position. However, these papers reported difficulty in learning velocity and acceleration (e.g., change in direction of motion), which are critical to our problem. We aim to obtain policy models to generate biologically plausible actions in terms of position, velocity, and acceleration. Our solution incorporates mechanical constraints into the policy learning. We describe the details for the policy model, macro-goal, and mechanical constraints in Section 3.3.

#### 3.2 Binary partial observation model

To model partial observation in a multi-agent setting, we propose a binary partial observation model using Gumbel-Softmax reparameterization. As shown in Figure 1, the model is trained to map input state vector  $s_t = [s_{t,1}, \dots, s_{t,K}] \in \mathbb{R}^{d_s \times K}$  to a latent binary vector  $b_{t,k}$  of an agent  $k$  at each time  $t$ , where  $d_s$  is the dimension of the state for each agent. We consider a partially observable environment, where each agent  $k$  receives a binary vector  $b_{t,k} := [e_{k,1}^t, e_{k,2}^t, \dots, e_{k,K}^t] \in \{0, 1\}^K$ . That is,  $b_{t,k}$  consists of an arbitrary number of zeros and ones. In this paper,  $b_{t,k}$  for an agent  $k$  is designed to represent the importance of all agents, used as the observation coefficients described below.

To obtain the binarized differentiable vector  $b_{t,k}$ , we linearly project a state vector  $s_{t,k}$  into embedded matrices  $f_{t,k} \in \mathbb{R}^{K \times d_e}$  and dual-channels  $f'_{t,k} \in \mathbb{R}^{K \times 2}$ , where  $d_e$  is the embedding dimension. The former embedding can transform the state into a distributed representation, which can be directly weighted by  $b_{t,k}$  as described below (note that e.g., Cartesian coordinates cannot be directly weighted in principle). The latter dual-channel projection inspired by [40] allows each dimension in  $f'_{t,k}$  (and finally in  $b_{t,k}$ ) to represent multiple agents’ importance (not limited to the one-hot approach in Section 2.2). We perform a categorical reparameterization trick with Gumbel-Softmax [30] on the second dimension of the dual-channel  $f'_{t,k}$  (rather than on the first dimension as used in the attention

mechanisms in 2.2). Since the two channels are linked together by Gumbel-Softmax, it is sufficient to simply pick one of them. That is,  $e_{k,i}^t = [\text{Gumbel-Softmax}([f_{t,k}' ]_i)]_1$ , where  $[\cdot]_i$  denotes  $[\cdot]$ 's  $i$ -th element. This method is differentiable and allows direct back-propagation in our framework.

We then compute the observation vector  $o_{t,k}$  as the input to the subsequent policy learning. We concatenate the element-wise product of binary vector  $b_{t,k}$  (observation coefficients) and the embedded matrices  $f_{t,k}$  for all agents:  $o_{t,k} = [e_{k,1}^t[f_{t,k}]_1, \dots, e_{k,K}^t[f_{t,k}]_K] \in \mathbb{R}^{d_e K}$ . This allows us to eliminate the information from unrelated agents, and to focus on only the important agents.

### 3.3 Hierarchical VRNN with macro-goals and mechanical constraints

**Hierarchical VRNN with macro-goals.** First, we perform VRNN modeling with conditional context information including the observation and macro-goals, as illustrated in Figure 2. To model an agent's macroscopic intent during path planning, we employ weak labels for macro-goals used in [59]. This method is currently one of the best methods for long-term prediction because it utilizes the future positional information trained by weak labels (for details, see Appendix C). Here we use a macro goal as a part of the planner and modify it in a decentralized and partially observable setting (see Appendix C). We use the partial observation  $o_{t-1,k}$  and independently learn the decentralized macro-goal  $g_{t,k}'$  (not shared between agents). Overall, our VRNN-based model becomes:

$$p_{\theta_k}(a_{t,k}|o_{<t,k}) = \mathcal{N}(a_{t,k}|\mu_{dec}^{t,k}, (\sigma_{dec}^{t,k})^2) \quad (2)$$

where  $[\mu_{dec}^{t,k}, \sigma_{dec}^{t,k}] = \varphi_{dec}^k(o_{t-1,k}, z_{t,k}, h_{t-1,k}, g_{t,k}')$ ,  $\varphi_{dec}^k$  is the Gaussian VRNN decoder,  $z_{t,k}$  is the VRNN latent variable, and  $h_{t,k}$  is the hidden state of an RNN. To sample macro-goals, we train another RNN-model in a similar manner to that in [59]:  $p(g_{t,k}'|g_{<t,k}') = \mathcal{N}(g_{t,k}'|\mu_{g'}^{t,k}, (\sigma_{g'}^{t,k})^2)$ , where  $[\mu_{g'}^{t,k}, \sigma_{g'}^{t,k}] = \varphi_{g'}^k(h_{g,t-1,k}, o_{t-1,k})$ ,  $\varphi_{g'}^k$  is the Gaussian macro-goal decoder, and  $h_{g,t-1,k}$  summarizes the history of macro-goals.

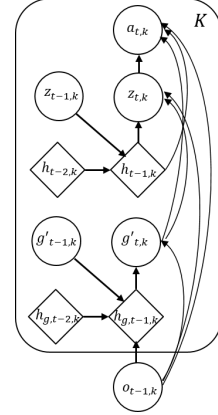


Figure 2: Graphical illustration of our model.

**Mechanical constraints.** Next, we model mechanical constraints for training of the policy model. We propose two types of mechanical constraints: a biomechanical constraint for the smooth path planning as mentioned above and physical constraints in multiple dimensions (i.e., position, velocity, and acceleration). We aim to obtain policy models to generate biologically plausible actions in terms of the above dimensions. The detailed objective function is described in Section 3.4.

For the biomechanical constraint, we consider a minimum-jerk principle [15] as mentioned in Section 2.4. However, this principle assumes that we know the start and end times of the movement. Then, we consider the penalty for minimum change in acceleration, only considering the current and next acceleration, in an analogous way of the minimum torque change principle in multi-joint motion [53]. That is, we propose the penalty which minimizes the difference between the predicted acceleration  $\hat{a}_{acc,t}$  and the true acceleration  $a_{acc,t+1}$  in the next step. For example, a negative log-likelihood (NLL)  $-\log p_{\theta}(a_{acc,t+1} | z_{\leq t}, o_{<t}, g'_{<t})$  or a Euclidean norm  $\|a_{acc,t+1} - \hat{a}_{acc,t}\|_2^2$  can be used.

Secondly, we propose physical constraints between multiple dimensions to consistently learn the relationship. Obviously, the model can effectively learn the output dimension (e.g., position) but not the others (e.g., velocity). By contrast, learning multiple dimensions results in physically-unrealistic and inconsistent predictions between the dimensions. We then propose two physical penalties for the predicted and eliminated output dimensions. For clarity, the following example uses velocity. For the predicted dimension, we propose a penalty between the directly and indirectly predicted velocity denoted by  $\hat{a}_{vel,t}$  and  $\tilde{a}_{vel,t} = (\hat{a}_{pos,t} - \hat{a}_{pos,t-1})/\Delta t$ , respectively, where the  $\hat{a}_{pos,t}$  is directly predicted position and  $\Delta t$  is a sampling interval. For example, a KL divergence between the two distributions (below) or the Euclidean norm can be proposed. It can be also computed for acceleration, but not for position in principle. For the eliminated dimension, we propose a penalty between the indirect prediction  $\tilde{a}_{m,t}$  and true  $a_{m,t}$ , where  $m = \{pos, vel, acc\}$ , e.g., a NLL between the distribution of  $\tilde{a}_{m,t}$  and  $a_{m,t}$  (below) or the Euclidean norm.

### 3.4 Learning

The objective function becomes a sequential ELBO of the hierarchical VRNN and the penalties of the mechanical constraints. The ELBO for agent  $k$  is:  $\mathcal{L}_{vrnn} = \mathbb{E}_{q_\phi(z_{\leq T} | o_{\leq t}, g'_{\leq t})} \sum_{t=1}^T [\log p_\theta(a_t | z_{\leq t}, o_{< t}, g'_{\leq t}) - D_{KL}(q_\phi(z_t | o_{\leq t}, z_{< t}, g'_{\leq t}) || p_\theta(z_t | o_{< t}, z_{< t}, g'_{\leq t}))]$ . For brevity, the symbol of the agent  $k$  is now shown here. The penalties for mechanical constraints  $\mathcal{L}_{body}$  are:

$$\sum_{t=2}^T \left[ D_{KL}(p_\theta(\hat{a}_{m',t} | z_{\leq t}, o_{< t}, g'_{\leq t}) || p_\theta(\tilde{a}_{m',t} | z_{\leq t}, o_{< t}, g'_{\leq t})) - \log p_\theta(a_{acc,t+1} | z_{\leq t}, o_{< t}, g'_{\leq t}) \right], \quad (3)$$

where  $m' = \{vel, acc\}$  indicates velocity and acceleration as action outputs, depending on experimental conditions. The first term is the penalty for consistently learning the relationship between the directly and indirectly estimated dimensions. If we eliminate some dimensions (e.g., position), we can add NLLs as the penalty between the distribution of the indirect prediction and ground truth:  $-\log p_\theta(a_{m,t} | z_{\leq t}, o_{< t}, g'_{\leq t})$ . The second term is the biomechanical smoothing penalty, which minimizes the difference between the distribution of the directly or indirectly predicted ( $\hat{a}_{acc,t}$  or  $\tilde{a}_{acc,t}$ ) and the observed acceleration  $a_{acc,t+1}$  in the next step. We jointly maximize  $\mathcal{L}_{vrnn} + \lambda \mathcal{L}_{body}$  with respect to all of model parameters (for the specific regularization parameter  $\lambda$ , see Appendix D).

## 4 Related Work

**Deep generative models for sequential data.** There has been recently increasing interest in deep generative models for sequential data, because of the flexibility of deep learning and (often probabilistic) generative models. In particular, many researchers intensively developed the methods for the physical systems with governing equations, such as in a bouncing ball and a pendulum. These works modeled latent linear state space dynamics (e.g., [34, 16]) or ordinary differential equations (e.g., [9, 57]) mainly without RNNs. Meanwhile, in biological systems with partially unknown governing equations (e.g., external forces and/or non-trivial interactions), RNN-based models are still in use (see below). Among the RNN-based generative models [11, 17, 22], we incorporate a VRNN [11] with a stochastic latent state into the observation model and prior knowledge about mechanics.

**Biological multi-agent trajectory prediction.** In existing researches in various biological multi-agent trajectories, pedestrian prediction problems including rule-based models (e.g., [25]) have been widely investigated for a long time. Recent works using RNN-based models (e.g., [1, 24]) effectively aggregated information across multiple persons using specialized pooling modules, but most of which predicted trajectories in a centralized manner. Vehicle trajectory prediction has been intensively researched, often using deep generative models based on RNNs [4, 47, 52], whereas little research has focused on animals [14, 32]. For sports multi-agent trajectory prediction such as basketball and soccer, most methods have leveraged RNNs [60, 38, 42, 37, 29] including VRNNs [59, 56], although some have utilized generative adversarial networks [8, 27] without RNNs. Most of these works assumed full observation to achieve long-term prediction (e.g., [59, 56]), except for an image-based work on partial observation [51]. In contrast, we aim to model decentralized biological multi-agent systems and visualize their partial observations for behavioral analyses in real-world agents.

**Observation in multi-agent systems.** In pure rule-based models, researchers investigating animals and vehicles conventionally proposed pre-defined observation rules based on specific distances and visual angles (e.g., [13]), the specific number of the nearest agents (e.g., [3]), or other environments (e.g., [58]). However, it is difficult to define interactions between agents using pre-defined rules in general large-scale multi-agent systems. Thus, methods for learning interaction between agents have been proposed, particularly for MARL in virtual environments [44, 31, 28, 41]. In real-world multi-agent systems, a few applications in vehicles [39] and physical and biological systems [23, 51] existed. However, binary observation in a decentralized setting has not previously been considered.

## 5 Experiments

We quantitatively compared our models to various baselines using basketball and soccer game datasets. The former dataset includes observation of smaller players and more frequent fine-grained acceleration due to the smaller playing area than the latter. We mainly visualize the results of the basketball dataset, but quantitatively validated our methods using both datasets to demonstrate

versatility. In the experiments, we focused on learning team defense policies such as [38] because defensive players can be regarded as decentralized agents for short-term periods (several seconds). We provided the states of the offense players and the ball as hierarchical inputs to our models, and updated the states using methods of [38]. Modeling the ball and offense was left for future work.

## 5.1 Common setups

**Baselines.** We compared our approach with four baselines: (1) Velocity, (2) RNN-gauss, (3) VRNN, and (4) VRNN-macro. First, as a sanity check, we used velocity extrapolation as a simple baseline, i.e., each of the agent’s predictions was linearly extrapolated from its previously observed velocity. The second baseline is an RNN implemented using a gated recurrent unit (GRU) [10] and a decoder with Gaussian distribution for prediction [5]. The third baseline is a VRNN [11] as our base-model. The last is VRNN with macro-goals as weak supervision [59], which is a state-of-the-art method in ballgame trajectory prediction (but implemented here as a decentralized version for fair comparison).

**Our models.** We validated our approach with three variants: (1) VRNN-Mech, (2) VRNN-Bi, and (3) VRNN-macro-Bi-Mech. First, we evaluated the effectiveness of our mechanical constraints (denoted by Mech). Since our proposed binary partial observation model (denoted by Bi) does not necessarily improve the prediction performance, we evaluated the prediction performance of VRNN-Bi and VRNN-macro-Bi-Mech as a validation.

**Training.** We trained all the models using the Adam optimizer [35] with default parameters using teacher forcing [55]. To prevent over-fitting, dropout and batch normalization layers were used (dropout rate was set to 0.5), and we selected the model with the best performance on the validation set. To obtain a long-term prediction for validation and test, we evaluated 60 timesteps after an initial burn-in period of 20 timesteps with ground-truth states (with data sampled at 10 Hz for both datasets). The duration of the prediction was similar to that of the previous works [59, 56]. We selected xy position, velocity, and acceleration as the input states (i.e.,  $d_s = 6$ ). The reason and the analysis of various inputs and outputs are described in Appendix F. For other training details, see Appendix D.

**Prediction metrics.** Due to the difficulty in evaluating a generative model with a single criterion [56], we evaluated several different metrics in terms of prediction error and distribution to demonstrate the effectiveness of our approach on the test set. For prediction error, we used the two basic metrics: the mean and best  $L_2$ -error between prediction and ground truth. Due to the multimodal nature of the system, we randomly sampled  $N = 10$  trajectories for each test case. More precisely, we compute the  $L_2$ -error between ground truth and the  $n$ ’th generated sample  $\hat{a}^{n,m}$  using:  $L_2^{n,m} = (1/(T \cdot K)) \sum_{t=1}^T \sum_{k=1}^K \|\hat{a}_{m,t,k}^n - a_{m,t,k}\|_2$ . We then reported the average and best result from the samples  $N$ , as is standard practice (e.g., [49, 7]). Note that although we used the ELBO for training, the tighter bounds do not necessarily lead to better performance [46]. To examine the plausibility of acceleration and its temporal change, we compared boxplots (distributions) of their  $L_2$  norms [56].

## 5.2 Basketball data

We used the basketball dataset from the NBA 2015-2016 season (<https://www.stats.com/data-science/>), which contains tracking trajectories for professional basketball players and the ball. The data was pre-processed such that the offense team always moved towards the left-side of the court. We chose 100 games so that the amount of data was similar to the subsequent soccer dataset. In total the dataset contained 19968 training, 2235 validation, and 2608 test sequences.

**Prediction performance.** Table 1 (left) shows prediction performances. We confirmed the effectiveness of our proposed mechanical constraints (VRNN-Mech). Our binary observation models (VRNN-Bi, VRNN-macro-Bi-Mech) show similar prediction performances to the baseline models (VRNN, VRNN-macro) regardless of the partial observation. Figure 3 (top) shows the distributions of acceleration and its temporal change. Our VRNN-macro-Bi-Mech (red) was closer to the true distribution (blue) than conventional methods (VRNN-macro: green and VRNN: orange). Figure 4D shows acceleration examples for defender #1. Our model (red) generated

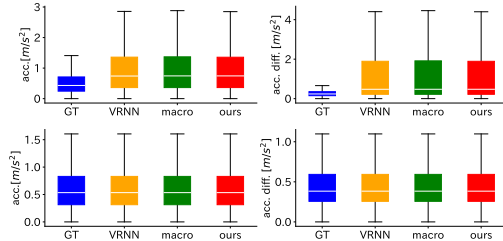


Figure 3: Acceleration and its temporal change for basketball (top) and soccer (bottom) datasets.

smoother accelerations than VRNN-macro (green). The detailed analyses, discussion about VRNN-macro, and the overall mean  $L_2$  prediction errors are described in Appendices E, C, and G, respectively.

	Basketball data			Soccer data		
	position	velocity	acceleration	position	velocity	acceleration
Velocity	$1.41 \pm 0.34$	$1.08 \pm 0.21$	$10.90 \pm 2.09$	$4.83 \pm 1.27$	$2.72 \pm 0.46$	$27.22 \pm 4.56$
RNN-Gauss	$1.31 \pm 0.32$	$1.05 \pm 0.13$	$1.88 \pm 0.30$	$2.97 \pm 0.89$	$1.65 \pm 0.23$	$2.22 \pm 0.33$
VRNN	$0.71 \pm 0.17$	$0.68 \pm 0.10$	$1.43 \pm 0.20$	$1.23 \pm 0.35$	$1.06 \pm 0.21$	$1.57 \pm 0.26$
VRNN-macro	$0.71 \pm 0.17$	$0.68 \pm 0.10$	$1.43 \pm 0.20$	$1.37 \pm 0.40$	$1.11 \pm 0.21$	$1.57 \pm 0.24$
VRNN-Mech	<b><math>0.69 \pm 0.17</math></b>	$0.68 \pm 0.10$	$1.37 \pm 0.20$	$1.22 \pm 0.34$	$1.05 \pm 0.20$	$1.59 \pm 0.26$
VRNN-Bi	$0.72 \pm 0.19$	<b><math>0.66 \pm 0.10</math></b>	$1.36 \pm 0.19$	<b><math>1.20 \pm 0.37</math></b>	<b><math>1.02 \pm 0.20</math></b>	$1.50 \pm 0.24$
VRNN-macro-Bi-Mech	$0.73 \pm 0.18$	$0.68 \pm 0.10$	<b><math>1.34 \pm 0.19</math></b>	$1.31 \pm 0.42$	$1.06 \pm 0.21$	<b><math>1.47 \pm 0.23</math></b>

Table 1: Best  $L_2$  prediction errors for the basketball and soccer datasets. Units in positions are meters.

**Evaluation of observation model.** We visually and quantitatively evaluated our VRNN-macro-Bi-Mech. The averaged observation coefficient (i.e., binary vector  $b_{t,k}$ ) was  $4.42 \pm 0.31$  for each defender (maximum: 11). Figure 4E (at the moment with larger marks in Figure 4B) shows that the defenders observed both near players and far players (the sequences for  $b_{t,k}$  are shown in Appendix H). Quantitatively, we compared the two distances between the subject player and the furthest player observed in our and rule-based models. The latter considered the same-order nearest player as the number of players observed in our model, which is a compatible and famous pre-determined rule [3]. The mean distance from the furthest player in our model was larger than that of the same-order nearest player ( $7.56 \pm 1.30$  and  $3.39 \pm 0.81$  [m]), indicating that our model reflected far agent information. For example, defender #1 (arrow) in Figures 4A and B adopted the balanced position between the attacker #1 and others to help teammates near the ball (black). Finally, our model can create a counterfactual example when  $b_{t,k}$  is artificially set to a one-hot vector at each moment (only the player highest in  $b_{t,k}$ ) in Figures 4C and F. For example, like sports beginners, defender #1 went to the nearest attacker #1 (ignored the ball). Our model can analyze such behaviors in real-world agents.

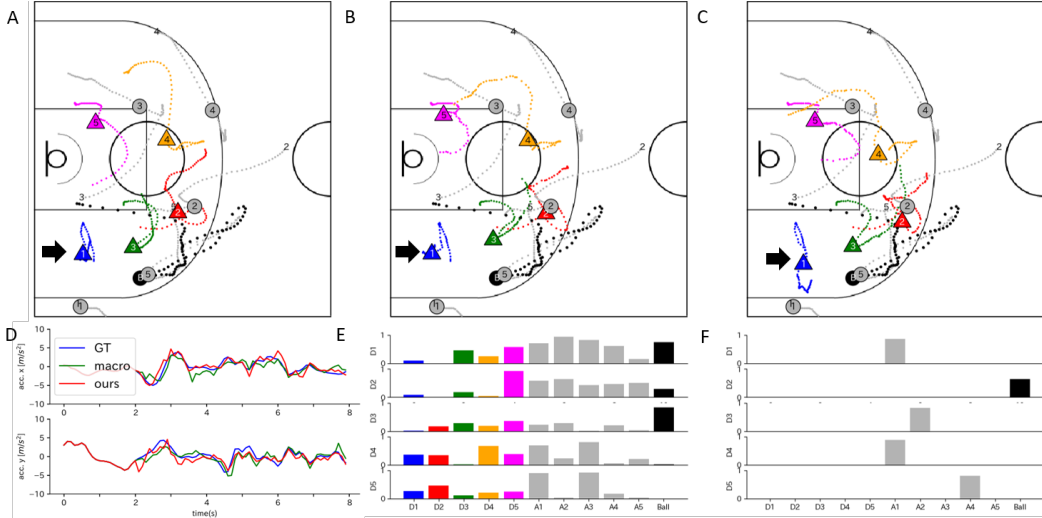


Figure 4: Example results using our method. (A) Ground truth, (B) a normal and (C) a counterfactual prediction are shown. Colored triangles, gray circles, the black circle are defenders, attackers, and the ball. (D) defender #1's acceleration ( $x, y$ ), and five defenders  $b_{t,k}$  in (E) a normal and (F) a counterfactual prediction are shown.

### 5.3 Soccer data

To demonstrate the generality of our approach, we also used a soccer dataset [38, 56], containing trajectories of soccer players and the ball from 45 professional soccer league games. We randomly split the dataset into 21504 training, 2165 validation, and 2452 test sequences. We did not model the goalkeepers since they tend to move very little. Table 1 (right) indicate the similar tendencies of prediction performance to that of the basketball dataset. We confirmed the effectiveness of our mechanical constraints and similar prediction performance of our binary observation models to the baselines. The distributions in Figure 3 (bottom) were similar to the ground truth in all methods. For



VRNN-macro-Bi-Mech, the averaged  $b_{t,k}$  was  $8.04 \pm 1.54$  for each defender (maximum: 23). The model reflected far agent information (the mean distance in our model:  $26.59 \pm 5.44$  and that in the same-order nearest player:  $18.09 \pm 6.52$  [m]). An illustrative example is shown in Appendix I.

## 6 Conclusions

We proposed sequential generative models for policy models with partial observation and mechanical constraints. We visualized partial observation and predicted biologically plausible actions using real-world basketball and soccer datasets. One possible future research direction is to incorporate other physiological constraints into the models such as visuomotor delays and body loads, which could contribute to the understanding of the biological multi-agent behaviors in complex environments.

## Acknowledgments

This work was supported by JSPS KAKENHI (Grant Numbers 19H04941 and 20H04075).

## References

- [1] A. Alahi, K. Goel, V. Ramanathan, A. Robicquet, L. Fei-Fei, and S. Savarese. Social lstm: Human trajectory prediction in crowded spaces. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 961–971, 2016.
- [2] D. Bahdanau, K. Cho, and Y. Bengio. Neural machine translation by jointly learning to align and translate. *arXiv preprint arXiv:1409.0473*, 2014.
- [3] M. Ballerini, N. Cabibbo, R. Candelier, A. Cavagna, E. Cisbani, I. Giardina, V. Lecomte, A. Orlandi, G. Parisi, A. Procaccini, et al. Interaction ruling animal collective behavior depends on topological rather than metric distance: Evidence from a field study. *Proceedings of the National Academy of Sciences*, 105(4):1232–1237, 2008.
- [4] M. Bansal, A. Krizhevsky, and A. Ogale. Chauffeurnet: Learning to drive by imitating the best and synthesizing the worst. *arXiv preprint arXiv:1812.03079*, 2018.
- [5] S. Becker, R. Hug, W. Hübner, and M. Arens. Red: A simple but effective baseline predictor for the trajnet benchmark. In *European Conference on Computer Vision*, pages 138–153. Springer, 2018.
- [6] D. S. Bernstein, R. Givan, N. Immerman, and S. Zilberstein. The complexity of decentralized control of markov decision processes. *Mathematics of Operations Research*, 27(4):819–840, 2002.
- [7] A. Bhattacharyya, B. Schiele, and M. Fritz. Accurate and diverse sampling of sequences based on a best of many sample objective. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 8485–8493, 2018.
- [8] C.-Y. Chen, W. Lai, H.-Y. Hsieh, W.-H. Zheng, Y.-S. Wang, and J.-H. Chuang. Generating defensive plays in basketball games. In *Proceedings of the 26th ACM International Conference on Multimedia*, pages 1580–1588, 2018.
- [9] T. Q. Chen, Y. Rubanova, J. Bettencourt, and D. K. Duvenaud. Neural ordinary differential equations. In *Advances in Neural Information Processing Systems 31*, pages 6571–6583, 2018.
- [10] K. Cho, B. Van Merriënboer, D. Bahdanau, and Y. Bengio. On the properties of neural machine translation: Encoder-decoder approaches. *arXiv preprint arXiv:1409.1259*, 2014.
- [11] J. Chung, K. Kastner, L. Dinh, K. Goel, A. C. Courville, and Y. Bengio. A recurrent latent variable model for sequential data. In *Advances in Neural Information Processing Systems 28*, pages 2980–2988, 2015.
- [12] R. M. Cichy and D. Kaiser. Deep neural networks as scientific models. *Trends in Cognitive Sciences*, 23(4):305–317, 2019.
- [13] I. D. Couzin, J. Krause, R. James, G. D. Ruxton, and N. R. Franks. Collective memory and spatial sorting in animal groups. *Journal of Theoretical Biology*, 218(1):1–11, 2002.
- [14] E. Eyjolfsson, K. Branson, Y. Yue, and P. Perona. Learning recurrent representations for hierarchical behavior modeling. In *International Conference on Learning Representations*, 2017.
- [15] T. Flash and N. Hogan. The coordination of arm movements: an experimentally confirmed mathematical model. *Journal of Neuroscience*, 5(7):1688–1703, 1985.
- [16] M. Fraccaro, S. Kamronn, U. Paquet, and O. Winther. A disentangled recognition and nonlinear dynamics model for unsupervised learning. In *Advances in Neural Information Processing Systems 30*, pages 3601–3610, 2017.
- [17] M. Fraccaro, S. K. Sønderby, U. Paquet, and O. Winther. Sequential neural models with stochastic layers. In *Advances in Neural Information Processing Systems 29*, pages 2199–2207, 2016.

- [18] K. Fujii, T. Isaka, M. Kouzaki, and Y. Yamamoto. Mutual and asynchronous anticipation and action in sports as globally competitive and locally coordinative dynamics. *Scientific Reports*, 5, 2015.
- [19] K. Fujii, N. Takeishi, B. Kibushi, M. Kouzaki, and Y. Kawahara. Data-driven spectral analysis for coordinative structures in periodic human locomotion. *Scientific Reports*, 9(1):1–14, 2019.
- [20] K. Fujii, K. Yokoyama, T. Koyama, A. Rikukawa, H. Yamada, and Y. Yamamoto. Resilient help to switch and overlap hierarchical subsystems in a small human group. *Scientific Reports*, 6, 2016.
- [21] K. Fujii, S. Yoshioka, T. Isaka, and M. Kouzaki. The preparatory state of ground reaction forces in defending against a dribbler in a basketball 1-on-1 dribble subphase. *Sports Biomechanics*, 14(1):28–44, 2015.
- [22] A. G. A. P. Goyal, A. Sordoni, M.-A. Côté, N. R. Ke, and Y. Bengio. Z-forcing: Training stochastic recurrent networks. In *Advances in Neural Information Processing Systems 30*, pages 6713–6723, 2017.
- [23] L. Guanyu, J. Bo, Z. Hao, C. Zhengping, and L. Yan. Generative attention networks for multi-agent behavioral modeling. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- [24] A. Gupta, J. Johnson, L. Fei-Fei, S. Savarese, and A. Alahi. Social gan: Socially acceptable trajectories with generative adversarial networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2255–2264, 2018.
- [25] D. Helbing and P. Molnar. Social force model for pedestrian dynamics. *Physical Review E*, 51(5):4282, 1995.
- [26] S. Hochreiter and J. Schmidhuber. Long short-term memory. *Neural Computation*, 9(8):1735–1780, 1997.
- [27] H.-Y. Hsieh, C.-Y. Chen, Y.-S. Wang, and J.-H. Chuang. Basketballgan: Generating basketball play simulation through sketching. In *Proceedings of the 27th ACM International Conference on Multimedia*, pages 720–728, 2019.
- [28] S. Iqbal and F. Sha. Actor-attention-critic for multi-agent reinforcement learning. In *International Conference on Machine Learning*, pages 2961–2970, 2019.
- [29] B. Ivanovic, E. Schmerling, K. Leung, and M. Pavone. Generative modeling of multimodal multi-human behavior. In *2018 IEEE/RSJ International Conference on Intelligent Robots and Systems*, pages 3088–3095. IEEE, 2018.
- [30] E. Jang, S. Gu, and B. Poole. Categorical reparametrization with gumbel-softmax. In *International Conference on Learning Representations*. OpenReview. net, 2017.
- [31] J. Jiang and Z. Lu. Learning attentional communication for multi-agent cooperation. In *Advances in Neural Information Processing Systems 31*, pages 7254–7264, 2018.
- [32] M. J. Johnson, D. K. Duvenaud, A. Wiltchko, R. P. Adams, and S. R. Datta. Composing graphical models with neural networks for structured representations and fast inference. In *Advances in Neural Information Processing Systems 29*, pages 2946–2954, 2016.
- [33] L. P. Kaelbling, M. L. Littman, and A. R. Cassandra. Planning and acting in partially observable stochastic domains. *Artificial Intelligence*, 101(1-2):99–134, 1998.
- [34] M. Karl, M. Soelch, J. Bayer, and P. van der Smagt. Deep variational bayes filters: Unsupervised learning of state space models from raw data. In *International Conference on Learning Representations*, 2017.
- [35] D. P. Kingma and J. Ba. Adam: A method for stochastic optimization. In *International Conference on Learning Representations*, 2015.
- [36] D. P. Kingma and M. Welling. Auto-encoding variational bayes. In *International Conference on Learning Representations*, 2014.
- [37] T. Kipf, E. Fetaya, K.-C. Wang, M. Welling, and R. Zemel. Neural relational inference for interacting systems. In *International Conference on Machine Learning*, pages 2688–2697, 2018.
- [38] H. M. Le, Y. Yue, P. Carr, and P. Lucey. Coordinated multi-agent imitation learning. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pages 1995–2003. JMLR. org, 2017.
- [39] E. Leurent and J. Mercat. Social attention for autonomous decision-making in dense traffic. *arXiv preprint arXiv:1911.12250*, 2019.
- [40] A. T. Liu, P.-c. Hsu, and H.-Y. Lee. Unsupervised end-to-end learning of discrete linguistic units for voice conversion. *Proc. Interspeech 2019*, pages 1108–1112, 2019.
- [41] Y. Liu, W. Wang, Y. Hu, J. Hao, X. Chen, and Y. Gao. Multi-agent game abstraction via graph attention neural network. In *Thirty-Fourth AAAI Conference on Artificial Intelligence*, 2020.
- [42] Y. Liu, R. Yu, S. Zheng, E. Zhan, and Y. Yue. Naomi: Non-autoregressive multiresolution sequence imputation. In *Advances in Neural Information Processing Systems 32*, pages 11236–11246, 2019.
- [43] C. J. Maddison, A. Mnih, and Y. W. Teh. The concrete distribution: A continuous relaxation of discrete random variables. In *International Conference on Learning Representations*, 2017.

- [44] H. Mao, Z. Zhang, Z. Xiao, and Z. Gong. Modelling the dynamic joint policy of teammates with attention multi-agent ddpg. In *Proceedings of the 18th International Conference on Autonomous Agents and MultiAgent Systems*, pages 1108–1116. International Foundation for Autonomous Agents and Multiagent Systems, 2019.
- [45] C. H. Papadimitriou and K. Steiglitz. *Combinatorial optimization*, volume 24. Prentice Hall Englewood Cliffs, 1982.
- [46] T. Rainforth, A. Kosiorek, T. A. Le, C. Maddison, M. Igl, F. Wood, and Y. W. Teh. Tighter variational bounds are not necessarily better. In *International Conference on Machine Learning*, volume 80, pages 4277–4285, 2018.
- [47] N. Rhinehart, R. McAllister, K. Kitani, and S. Levine. Precog: Prediction conditioned on goals in visual multi-agent settings. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 2821–2830, 2019.
- [48] S. Ross, G. Gordon, and D. Bagnell. A reduction of imitation learning and structured prediction to no-regret online learning. In *Proceedings of the fourteenth International Conference on Artificial Intelligence and Statistics*, pages 627–635, 2011.
- [49] C. Rupprecht, I. Laina, R. DiPietro, M. Baust, F. Tombari, N. Navab, and G. D. Hager. Learning in an uncertain world: Representing ambiguity through multiple hypotheses. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3591–3600, 2017.
- [50] S. Schaal. Is imitation learning the route to humanoid robots? *Trends in Cognitive Sciences*, 3(6):233–242, 1999.
- [51] C. Sun, P. Karlsson, J. Wu, J. B. Tenenbaum, and K. Murphy. Predicting the present and future states of multi-agent systems from partially-observed visual data. In *International Conference on Learning Representations*, 2019.
- [52] C. Tang and R. R. Salakhutdinov. Multiple futures prediction. In *Advances in Neural Information Processing Systems 32*, pages 15398–15408, 2019.
- [53] Y. Uno, M. Kawato, and R. Suzuki. Formation and control of optimal trajectory in human multijoint arm movement. *Biological Cybernetics*, 61(2):89–101, 1989.
- [54] X. Wang, R. Girshick, A. Gupta, and K. He. Non-local neural networks. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 7794–7803, 2018.
- [55] R. J. Williams and D. Zipser. A learning algorithm for continually running fully recurrent neural networks. *Neural Computation*, 1(2):270–280, 1989.
- [56] R. A. Yeh, A. G. Schwing, J. Huang, and K. Murphy. Diverse generation for multi-agent sports games. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4610–4619, 2019.
- [57] C. Yildiz, M. Heinonen, and H. Lahdesmaki. Ode2vae: Deep generative second order odes with bayesian neural networks. In *Advances in Neural Information Processing Systems 32*, pages 13412–13421, 2019.
- [58] Y. Yoshihara, Y. Morales, N. Akai, E. Takeuchi, and Y. Ninomiya. Autonomous predictive driving for blind intersections. In *2017 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, pages 3452–3459, 2017.
- [59] E. Zhan, S. Zheng, Y. Yue, L. Sha, and P. Lucey. Generating multi-agent trajectories using programmatic weak supervision. In *International Conference on Learning Representations*, 2019.
- [60] S. Zheng, Y. Yue, and J. Hobbs. Generating long-term trajectories using deep hierarchical networks. In *Advances in Neural Information Processing Systems 29*, pages 1543–1551, 2016.

## A Role assignment problem

In most real-world data, training (i.e., demonstration) and test data include sequences from different types of agents (e.g., teams and games). Learning policies based on player’s positions or roles, instead of the identity, seems to be more natural and data-efficient in general. Among several approaches [38, 56], for the behavioral modeling in POMDP, we separate the problem into role assignment problem and policy learning based on [38].

The unstructured set of demonstrations is denoted by  $u = \{u_1, \dots, u_K\}$ , where  $u_k = \{u_{t,k}\}_{t=1}^T$  is the sequence of actions by agent  $k$  at time  $t$ . Let  $c = \{c_t\}_{t=1}^T$  be the context associated with each demonstration sequence (e.g., an opponent team and ball). In a role assignment problem, the indexing mechanism is formulate as an assignment function  $\mathcal{A}$  which maps the unstructured set  $c$  and some probabilistic structured model  $q$  to an indexed set of states  $s$  rearranged from  $u$ , i.e.,

$$\mathcal{A} : \{u_1, \dots, u_K, c\} \times q \mapsto [s_1, \dots, s_K, c],$$

where the set  $\{s_1, \dots, s_K\} \equiv \{u_1, \dots, u_K\}$ . We consider  $q$  as a latent variable model that infers the role assignments for each set of demonstrations. The role assignment with a latent structured model addresses two main issues: (1) unsupervised learning a probabilistic role assignment model  $q$ ; and (2) the indexing with  $q$  so that unstructured sequences can be mapped to structured sequences.

First, in the unsupervised learning of the stochastic role assignment model, we use a Gaussian Hidden Markov Model (Gaussian HMM) according to [38]. The Gaussian HMM inputs state feature vector defined by all agents’ information and learns transition probabilities and Gaussian mixture distributions as output probabilities of hidden states.

Second, for indexing based on  $q$  to map  $u$  to  $s$ , we solve a well-known linear assignment problem [45]. Concretely, the distance between the mixed Gaussian distribution obtained above and the state feature vector is computed for each time and each player. The linear assignment problem is then solved using the distance as a cost function (i.e., roles are assigned in order of players whose distance is closer to each Gaussian distribution).

## B Variational recurrent neural networks

In this section, we briefly overview recurrent neural networks (RNNs), variational autoencoders (VAEs), and variational RNNs (VRNNs).

From the perspective of probabilistic generative model, a RNN models the conditional probabilities with a hidden state  $h_t$  that summarizes the past history in the first  $t - 1$  timesteps:

$$p_\theta(a_t | a_{<t}) = \varphi(h_{t-1}), \quad h_t = f(a_t, h_{t-1}), \quad (4)$$

where  $\varphi$  maps the hidden state to a probability distribution over states and  $f$  is a deterministic function such as LSTMs [26] or GRUs [10]. RNNs with simple output distributions often struggle to capture highly variable and structured sequential data. Recent work in sequential generative models addresses this issue by injecting stochastic latent variables into the model and using amortized variational inference to infer latent variables from data. VRNNs [11] is one of the methods using this idea and combining RNNs and VAEs.

VAE [36] is a generative model for non-sequential data that injects latent variables  $z$  into the joint distribution  $p_\theta(a, z)$  and introduces an inference network parametrized by  $\phi$  to approximate the posterior  $q_\phi(z | a)$ . The learning objective is to maximize the evidence lower-bound (ELBO) of the log-likelihood with respect to the model parameters  $\theta$  and  $\phi$ :

$$\mathbb{E}_{q_\phi(z|a)} [\log p_\theta(a|z)] - D_{KL}(q_\phi(z | a) || p_\theta(z)) \quad (5)$$

The first term is known as the reconstruction term and can be approximated with Monte Carlo sampling. The second term is the Kullback-Leibler divergence between the approximate posterior and the prior, and can be evaluated analytically if both distributions are Gaussian with diagonal covariance. The inference model  $q_\phi(z | a)$ , generative model  $p_\theta(a | z)$ , and prior  $p_\theta(z)$  are often implemented with neural networks.

VRNNs combine VAEs and RNNs by conditioning the VAE on a hidden state  $h_t$ :

$$p_\theta(z_t | a_{<t}, z_{<t}) = \varphi_{\text{prior}}(h_{t-1}) \quad (\text{prior}) \quad (6)$$

$$q_\phi(z_t | a_{\leq t}, z_{<t}) = \varphi_{\text{enc}}(a_t, h_{t-1}) \quad (\text{inference}) \quad (7)$$

$$p_\theta(a_t | z_{\leq t}, a_{<t}) = \varphi_{\text{dec}}(z_t, h_{t-1}) \quad (\text{generation}) \quad (8)$$

$$h_t = f(a_t, z_t, h_{t-1}). \quad (\text{recurrence}) \quad (9)$$

VRNNs are also trained by maximizing the ELBO, which can be interpreted as the sum of VAE ELBOs over each timestep of the sequence:

$$\mathbb{E}_{q_\phi(z_{\leq T} | a_{\leq T})} \left[ \sum_{t=1}^T \log p_\theta(a_t | z_{\leq T}, a_{<t}) - D_{KL} \left( q_\phi(z_t | a_{\leq T}, z_{<t}) || p_\theta(z_t | a_{<t}, z_{<t}) \right) \right] \quad (10)$$

Note that the prior distribution of latent variable  $z_t$  depends on the history of states and latent variables (Eq. (6)).

## C Decentralized and partially observed macro-goals

Our model utilizes macro-goal [59] for the long-term prediction by modifying to our decentralized and partially-observable setting. In this section, we briefly overview the original macro-goal and describe our modification to the decentralized and partially-observable setting. We also discuss the related results of our experiments.

As an illustrative example, Figure C.1 shows macro-goals for a basketball defender as specific areas on the court (boxes). After reaching the macro-goal in the center, the blue player moves towards the next macro-goal in the top-middle. The macro-goals provide a compact summary of the players' sequences over a long time to encode long-term intent.

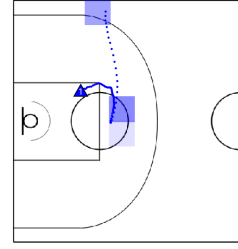


Figure C.1: A macro-goal (boxes) for a player.

**Shared macro-goals [59].** The original macro-goal or macro-intent [60, 59], obtained via some labeling functions, is defined as low-dimensional and spatiotemporal representations of the data for the learning of multi-agent long-term coordination such as basketball (Figure C.1). The original macro-goal assumes that: i) provide a tractable way to capture coordination between agents; ii) encode long-term goals of agents and enable long-term planning at a higher-level timescale; and iii) compactly represent some low-dimensional structure in an exponentially large multi-agent state space.

Specifically, the modeling assumptions for the original macro-goals are: 1) agent states  $s_t$  in a time period  $[t_1, t_2]$  are conditioned on some shared macro-goal  $g_t$ ; 2) the start and end times  $[t_1, t_2]$  of episodes can vary between sequences; 3) macro-goals change slowly over time relative to the agent states:  $dg_t/dt \ll 1$ ; and 4) due to their reduced dimensionality, (near-) arbitrary dependencies between macro-goals (e.g., coordination) can be modeled by a neural network approach.

Among various labeling approaches (see [59]) for the macro-goal, programmatic weak supervision is a method requiring low labor cost, effectively learning the underlying structure of large unlabeled datasets, and allowing users to incorporate domain knowledge into the model. For example, the labeling function to obtain macro-goals for basketball sequences computes the regions on the court in which players remain stationary; this integrates the idea that players aim to set up specific formations on the court.

Specifically, the previous work labeled the macro-goal independently among agents, and learning the shared macro-goal model via supervised learning by maximizing the log-likelihood of macro-goal labels obtained programmatically. The model finally returns the one-hot encoding of the box that contains the position information.

**Specific setup in our experiments.** Our decentralized and partially observed macro-goals i) use partial observation  $o_{t-1,k}$  to obtain the macro-goal model and ii) independently learn the decentralized macro-goal  $g'_{t,k}$  (not shared between agents).

Among several labeling functions, we adopted the stationary labeling function computing the macro-goal based on stationary positions because it reflects the important information about the structure of the data and the better performance was confirmed [59].

For basketball data, according to the previous work [59], we define the macro-goal by segmenting the left half-court into a  $10 \times 9$  grid of  $5 \times 5$  feet boxes (Figure C.1). For soccer data, we segmented the all-court into a  $34 \times 22$  grid of approximately  $3 \times 3$  m boxes.

**Related discussion in our experiments.** We confirmed the decentralized and partially observed macro-goals did not improve the prediction performances. There would be mainly two reasons. One is obviously the decentralized (i.e., not shared) setting, but it is a necessary assumption for our modeling. Second would be the improvement of the VRNN baseline by adding dropout and batch normalization layers to avoid overfitting (note that all models added them for fair comparisons). Specifically in the soccer experiment, the resolution of the grid might be larger than that of prediction (but the smaller grid might be difficult to learning the macro-goal).

## D Training details

In this section, we describe a specific objective function in our experiments, and other training details. For other implementation details such as pre-processing including role assignment in Appendix A, see the available code at <https://github.com/PO-MC-model-NeurIPS2020/PO-MC-model>.

### D.1 Specific objective function in our experiments.

We designed the objective function basically based on that described in Section 3.4, but we specially weighted each term of the objective function. Note that, in our experiments, we selected velocity and acceleration as the output actions (see in Appendix F). In this case, the penalty for the predicted and eliminated dimension can be computed only for acceleration and only for position, respectively. That is, the weighted penalties of mechanical constraints are:

$$\begin{aligned} \mathcal{L}_{w-body} = & \mathbb{E}_{\theta} \sum_{t=2}^T \left[ \lambda_{acc} D_{KL}(p_{\theta}(\hat{a}_{acc,t} \mid z_{\leq t}, o_{< t}, g'_{\leq t}) \parallel p_{\theta}(\hat{a}_{acc,t} \mid z_{\leq t}, o_{< t}, g'_{\leq t})) \right. \\ & \left. - \lambda_{pos} \log p_{\theta}(a_{pos,t} \mid z_{\leq t}, o_{< t}, g'_{< t}) - \lambda_{jrk} \log p_{\theta}(a_{acc,t+1} \mid z_{\leq t}, o_{< t}, g'_{< t}) \right], \quad (11) \end{aligned}$$

where  $\lambda_{acc}$ ,  $\lambda_{pos}$ ,  $\lambda_{jrk}$  are regularization parameters. As shown in Appendix E, the first and second terms improved the prediction performance of velocity and position, respectively (the third term basically contributed to smoothing in acceleration shown in Figure 4D). Results also show that the learning of acceleration was relatively difficult due to the above imbalance between the dimensions. We then added the reconstruction term  $\mathcal{L}_{acc}$  in  $\mathcal{L}_{vrnn}$  only for acceleration. In summary, the specific objective function in our experiments was  $\mathcal{L}_{vrnn} + \mathcal{L}_{w-body} + \lambda_{rec} \mathcal{L}_{acc}$ . We set  $\lambda_{acc} = 0.1$ ,  $\lambda_{pos} = 0.01$ ,  $\lambda_{jrk} = 0.1$ ,  $\lambda_{rec} = 0.2$  in the basketball experiment and  $\lambda_{acc} = 0.01$ ,  $\lambda_{pos} = 0.001$ ,  $\lambda_{jrk} = 0.01$ ,  $\lambda_{rec} = 0.02$  in the soccer experiment, because the soccer dataset was more difficult to be predicted than the basketball (i.e., the reconstruction was prioritized).

### D.2 Other training details.

We trained all the models using the Adam optimizer [35] with default parameters using teacher forcing [55]. To prevent over-fitting, dropout and batch normalization layers were used (dropout rate was set to 0.5), and we selected the model with the best performance on the validation set. We set the embedding dimension to  $d_e = 32$  for each agent (and the ball). The embedding layer, and macro-goal decoder  $\varphi_{g'_k}^k$ , prior  $\varphi_{prior}^k$ , encoder  $\varphi_{enc}^k$ , and decoder  $\varphi_{dec}^k$  in VRNN were implemented by a 2-layer neural network with a hidden layer of size 64. We modeled each latent variable  $z$  as a multivariate Gaussian with diagonal covariance of dimension 64. All GRUs were implemented by a 2-layer neural network with a hidden layer of size 100. Other implementations were based on [59]. We selected xy position, velocity, and acceleration as the input states (i.e.,  $d_s = 6$ ). The reason and the analysis of various input states and output actions are described in Appendix F. For the temperature  $\tau$  in a Gumbel-softmax distribution, we set it to 1 in both experiments.

VRNN decoder  $\varphi_{dec}^k$  were implemented by a 2-layer neural network returning a multivariate Gaussian with diagonal covariance. The original VRNN [11] has no constraint in the learning of the variance. It may cause that the NLL  $\log p_{\theta}(a_t \mid z_{\leq t})$  tends toward infinity when the variance approaches to zero. Most of cases including our experiments did not happen such problems (in such difficult case, setting the decoder variance as a global hyperparameter will be practically effective). In our model, the KL

divergence in the first term of Eq.(3) includes  $-2\log(\hat{\sigma}_{dec}^{t,k})$ , where  $\hat{\sigma}_{dec}^{t,k}$  is the directly estimated variances, and the term may prevent the variance from approaching to zero.

## E Analysis of mechanical constraints

In this section, we verified various mechanical constraints in our model. The detailed objective function was described in Appendix D. For clarity, we evaluated the prediction performance of VRNN using the basketball data among various options: (1) VRNN, (2) VRNN- $C_{pos}$ , (3) VRNN- $C_{pos,acc}$ , (4) VRNN- $C_{pos,acc,jrk}$ , and (5) VRNN- $C_{pos,acc,jrk}-\mathcal{L}_{acc}$  (VRNN-Mech in the main text). The second, third, and fourth options added the penalty of the second, first, third terms in Eq.(11), respectively. The fifth option added the acceleration reconstruction term  $\mathcal{L}_{acc}$ .

Table E.5 indicates the results of our analysis. The second terms in Eq.(11) improved the prediction performance of velocity (the third term basically contributed to smoothing in acceleration shown in Figure 4D). Results also show that the learning of acceleration was relatively difficult due to the above imbalance between the dimensions. We thus added the reconstruction term  $\mathcal{L}_{acc}$ . Our VRNN- $C_{pos,acc,jrk}-\mathcal{L}_{acc}$  (VRNN-Mech in the main text) shows better prediction performance in all dimensions.

	Average $L_2$			Best $L_2$		
	position	velocity	acceleration	position	velocity	acceleration
VRNN	$0.90 \pm 0.18$	$0.73 \pm 0.10$	$1.52 \pm 0.20$	$0.71 \pm 0.17$	$0.68 \pm 0.10$	$1.43 \pm 0.20$
VRNN- $C_{pos}$	$0.90 \pm 0.18$	$0.74 \pm 0.10$	$1.49 \pm 0.20$	$0.72 \pm 0.17$	$0.69 \pm 0.10$	$1.40 \pm 0.20$
VRNN- $C_{pos,acc}$	$0.88 \pm 0.18$	$0.71 \pm 0.10$	$1.54 \pm 0.21$	$0.70 \pm 0.17$	$0.66 \pm 0.10$	$1.45 \pm 0.21$
VRNN- $C_{pos,acc,jrk}$	$0.89 \pm 0.18$	$0.71 \pm 0.10$	$1.54 \pm 0.21$	$0.71 \pm 0.17$	$0.66 \pm 0.10$	$1.45 \pm 0.21$
VRNN- $C_{pos,acc,jrk}-\mathcal{L}_{acc}$	$0.87 \pm 0.18$	$0.72 \pm 0.10$	$1.46 \pm 0.21$	$0.69 \pm 0.17$	$0.68 \pm 0.10$	$1.37 \pm 0.20$

Table E.5: Average and best  $L_2$  prediction errors for various constraints using the basketball dataset.

## F Analysis of input state and output action dimensions

In this section, we verified various options to select input state and output action dimensions among position, velocity, and acceleration. The action in the previous works (e.g., [59, 56]) is usually the agent’s position, but these papers reported the difficulty in learning velocity and acceleration (e.g., change in direction of movement), which is critical in our problem. We aim to obtain policy models to generate biologically realistic actions in terms of position, velocity, and acceleration.

For simplicity, we evaluated the prediction performance of VRNN using the basketball data among various options: (1) VRNN-pos, (2) VRNN-vel, (3) VRNN-acc, (4) VRNN-pos-vel-acc, and (5) VRNN-vel-acc (in the main text). The first one uses only positional information as the input and output as the baseline used in most of the previous papers (e.g., [59, 56]). The second one uses position and velocity as the input and only velocity as the output to compare with VRNN-pos and our main VRNN-vel-acc. The remaining uses full information (position, velocity, and acceleration) as the input. The third and fourth use only acceleration and the full information as the output, respectively, to compare with our main VRNN-vel-acc. The last uses velocity and acceleration as output. The second, third, and the last compute the future positions using current position and velocity (the third similarly computes velocity).

Table F.5 indicates the results of the analysis. The model effectively learned the dimension of the output (e.g., position in VRNN-pos) but the indirect learning of the differential value (e.g., velocity) was difficult. In contrast, for the indirect learning of the integral value, the learning of position using velocity was effective, but that of velocity using acceleration was difficult. This may be caused by the data property having noisy accelerations. The learning of all dimensions (VRNN-pos-vel-acc) did not show better performance in all dimensions. VRNN-vel-acc shows better prediction performance in all dimensions. Therefore, we chose VRNN-vel-acc as our base-model.

## G Average prediction error

Overall, average  $L_2$  prediction errors in Table G.5 indicate similar tendencies the best  $L_2$  prediction errors in Table 1.

	Average $L_2$			Best $L_2$		
	position	velocity	acceleration	position	velocity	acceleration
VRNN-pos	$1.60 \pm 0.33$	$16.01 \pm 3.26$	$160.12 \pm 32.59$	$1.50 \pm 0.31$	$15.05 \pm 3.10$	$150.46 \pm 31.01$
VRNN-vel	$0.75 \pm 0.17$	$0.58 \pm 0.09$	$5.76 \pm 0.88$	$0.63 \pm 0.16$	$0.54 \pm 0.09$	$5.40 \pm 0.85$
VRNN-acc	$1.61 \pm 0.34$	$0.80 \pm 0.13$	$1.21 \pm 0.17$	$1.23 \pm 0.31$	$0.68 \pm 0.12$	$1.13 \pm 0.17$
VRNN-pos-vel-acc	$0.90 \pm 0.18$	$0.73 \pm 0.10$	$1.52 \pm 0.20$	$0.71 \pm 0.17$	$0.68 \pm 0.10$	$1.43 \pm 0.20$
VRNN-vel-acc (main)	$1.08 \pm 0.23$	$0.75 \pm 0.11$	$1.32 \pm 0.19$	$0.85 \pm 0.21$	$0.70 \pm 0.11$	$1.23 \pm 0.19$

Table F.5: Average and best  $L_2$  prediction errors for various inputs and outputs using the basketball dataset.

	Basketball data			Soccer data		
	position	velocity	acceleration	position	velocity	acceleration
Velocity	$1.41 \pm 0.34$	$1.08 \pm 0.21$	$10.90 \pm 2.09$	$4.83 \pm 1.27$	$2.72 \pm 0.46$	$27.22 \pm 4.56$
RNN-Gauss	$1.57 \pm 0.34$	$1.12 \pm 0.14$	$1.95 \pm 0.31$	$3.19 \pm 0.90$	$1.71 \pm 0.23$	$2.28 \pm 0.34$
VRNN	$0.90 \pm 0.18$	$0.73 \pm 0.10$	$1.52 \pm 0.20$	$1.46 \pm 0.39$	$1.11 \pm 0.21$	$1.63 \pm 0.27$
VRNN-macro	$0.90 \pm 0.18$	$0.73 \pm 0.10$	$1.53 \pm 0.21$	$1.63 \pm 0.44$	$1.17 \pm 0.21$	$1.63 \pm 0.25$
VRNN-Mech	$0.87 \pm 0.18$	$0.72 \pm 0.10$	$1.46 \pm 0.21$	$1.44 \pm 0.38$	$1.10 \pm 0.21$	$1.65 \pm 0.27$
VRNN-Bi	$0.90 \pm 0.20$	$0.71 \pm 0.11$	$1.44 \pm 0.19$	$1.41 \pm 0.41$	$1.07 \pm 0.20$	$1.55 \pm 0.24$
VRNN-macro-Bi-Mech	$0.93 \pm 0.20$	$0.73 \pm 0.11$	$1.42 \pm 0.19$	$1.55 \pm 0.46$	$1.11 \pm 0.22$	$1.53 \pm 0.24$

Table G.5: Average  $L_2$  prediction errors for basketball and soccer datasets.

## H Example sequences of the observation coefficients

We show example sequences of the observation coefficients in Figures H.4 A and B, which correspond to defender #1’s observation coefficients in Figures 4B and C (a normal and counterfactual prediction with one-hot observation using the basketball dataset), respectively. We confirmed the observation coefficients in the ball and the nearest attacker #1 were higher than other players in Figures H.4A and B.

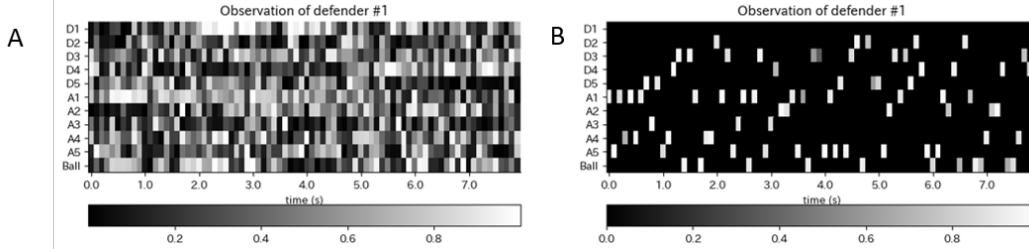


Figure H.4: Example sequences of the observation coefficients for the defender #1 in (B) a normal and (C) a counterfactual prediction are shown. These correspond to Figures 4B and C (a normal and counterfactual prediction with one-hot observation using the basketball dataset), respectively.

## I An illustrative example using the soccer dataset

We obtained an example of VRNN-macro-Bi-Mech using the soccer dataset in Figure I.4. Unlike the basketball example in Figure 4, the ground truth, a normal and a counterfactual prediction in Figures I.4 A, B, and C were similar, possibly because the soccer pitch ( $105 \times 68$  m) was larger than the basketball half court ( $14 \times 15$  m). For the observation model, the results in a normal and a counterfactual prediction in Figures I.4 D and E indicate that the defender #1 observed specific players such as the attacker #3 and the ball.



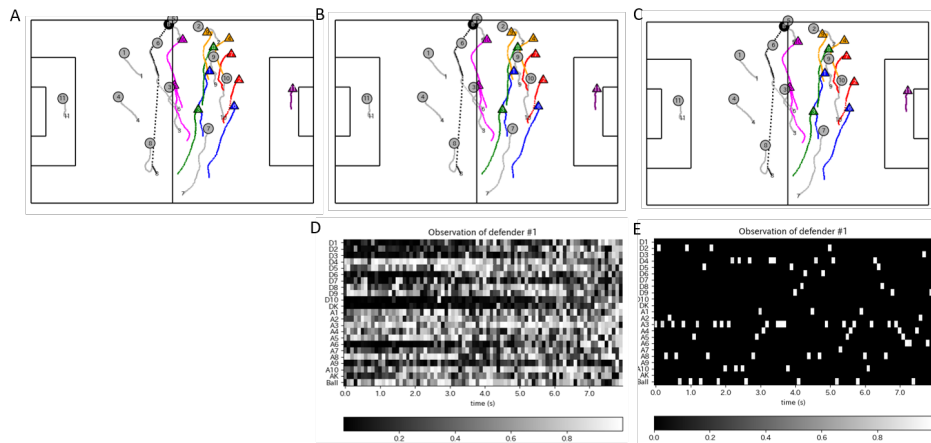


Figure I.4: Example results using our method. (A) Ground truth, (B) a normal and (C) a counterfactual prediction are shown. Colored triangles, gray circles, the black circle are defenders, attackers, and the ball. Defender #1's observations in (D) a normal and (E) a counterfactual prediction are shown.