

# **AGENTS LEARNING UNDER UNCERTAIN COMMUNICATION**

by

Kevin Corder

A dissertation submitted to the Faculty of the University of Delaware in partial fulfillment of the requirements for the degree of Doctor of Philosophy in Computer Science

Fall 2020

© 2020 Kevin Corder  
All Rights Reserved

# AGENTS LEARNING UNDER UNCERTAIN COMMUNICATION

by

Kevin Corder

Approved: \_\_\_\_\_  
Xxxx Xxxx, Highest Degree  
Chair of the Department of Xxxx

Approved: \_\_\_\_\_  
Xxxx Xxxx, Highest Degree  
Dean of the College of Xxxx

Approved: \_\_\_\_\_  
Douglas J. Doren, Ph.D.  
Interim Vice Provost for Graduate and Professional Education

I certify that I have read this dissertation and that in my opinion it meets the academic and professional standard required by the University as a dissertation for the degree of Doctor of Philosophy.

Signed: \_\_\_\_\_  
Xxxx Xxxx, Highest Degree  
Professor in charge of dissertation

I certify that I have read this dissertation and that in my opinion it meets the academic and professional standard required by the University as a dissertation for the degree of Doctor of Philosophy.

Signed: \_\_\_\_\_  
Xxxx Xxxx, Highest Degree  
Member of dissertation committee

I certify that I have read this dissertation and that in my opinion it meets the academic and professional standard required by the University as a dissertation for the degree of Doctor of Philosophy.

Signed: \_\_\_\_\_  
Xxxx Xxxx, Highest Degree  
Member of dissertation committee

I certify that I have read this dissertation and that in my opinion it meets the academic and professional standard required by the University as a dissertation for the degree of Doctor of Philosophy.

Signed: \_\_\_\_\_  
Xxxx Xxxx, Highest Degree  
Member of dissertation committee

## ACKNOWLEDGEMENTS

## TABLE OF CONTENTS

<b>LIST OF TABLES</b> . . . . .	<b>vii</b>
<b>LIST OF FIGURES</b> . . . . .	<b>viii</b>
<b>ABSTRACT</b> . . . . .	<b>x</b>
 <b>Chapter</b>	
<b>1 INTRODUCTION</b> . . . . .	<b>1</b>
<b>2 BACKGROUND AND RELATED WORK</b> . . . . .	<b>4</b>
2.1 Background . . . . .	4
2.1.1 Reinforcement Learning . . . . .	4
2.1.2 Multi-Agent Reinforcement Learning . . . . .	5
2.1.3 Generative Models . . . . .	5
2.2 Related Work . . . . .	7
<b>3 PROPOSED WORK (USE TITLE)</b> . . . . .	<b>8</b>
3.1 Context-based Generative Inference . . . . .	8
3.1.1 MADDPG with Inferred Observations . . . . .	10
3.1.1.1 Environments and Setup . . . . .	12
3.1.1.2 Results . . . . .	13
3.1.1.3 Discussion . . . . .	15
3.2 Context-history-based Generative Inference . . . . .	16

3.3	Domain Adaptation for CTDE Approaches . . . . .	16
3.4	Learning when to communicate (vs. sample) . . . . .	16
<b>4</b>	<b>EVALUATION AND TIMELINE . . . . .</b>	<b>23</b>
4.1	Preliminary Results . . . . .	23
4.2	How to evaluate rest of proposed work . . . . .	23
4.3	Timeline . . . . .	23
	<b>BIBLIOGRAPHY . . . . .</b>	<b>24</b>
	<b>Appendix</b>	
<b>A</b>	<b>TITLE OF APPENDIX 1 . . . . .</b>	<b>26</b>
<b>B</b>	<b>TITLE OF 2ND APPENDIX . . . . .</b>	<b>27</b>

## LIST OF TABLES



## LIST OF FIGURES

3.1	CC-WGAN update diagram. Describes an agent-based communication training procedure similar to random masking found in image inpainting tasks. . . . .	8
3.2	Random masking input and output used in training the CC-WGAN for an agent-distance partial observability criterion. Fig. 3.2(a) depicts $n$ agents' observation vectors fed to the masking function. Fig. 3.2(b) shows the masking function output: a set of $n$ partial observations and $n$ binary masks for agent $i$ . Masked values in observations are filled with random normal values $z \sim \mathcal{N}(0, 1)$ . . . . .	9
3.3	MADDPG with CC-WGAN learning diagram. During centralized training, the actor-critic learning updates are identical to MADDPG and CC-WGAN collects joint observations. During decentralized execution, agents sample from CC-WGAN to fill partial observations. Note that this diagram excludes the approximate policies $\mu$ ; only the observations are shared during centralized training. . . . .	17
3.4	Physical deception scenario with agent inference at the beginning of an episode. Agents 1 and 2 are near each other and can therefore observe each others' real positions. They are both too far away from Agent 0 (with partially observable distance $d_P$ ) so they each sample from the generative model to infer Agent 0's position. . . . .	18
3.5	Reward for cooperating agents with distance-based partial observations in decentralized phase. . . . .	19

3.6	Both sides of cooperating agent rewards with both distance-based partial observations and altering environment dynamics in decentralized phase. “Infer” is our approach and “None” is MADDPG. . . . .	20
3.7	Total reward for our approach with all four combinations of whether agent policies and CC-WGAN continue to update on inferred observations during decentralized phase. . . . .	21
3.8	Mean squared error CC-WGAN reconstruction between latest joint observations $o$ and inferred observations $\hat{o}$ for the predator-prey scenario with altered dynamics. Curves show varying partial observability distances $d_P$ had little effect on error. . . . .	22

## ABSTRACT

Reinforcement learning in multi-agent systems is an important problem that has increasingly many applications. Single agent reinforcement learning methods have recently been shown to be effective in low- and high-dimensional state spaces, such as game playing ( [CITE atari](#), [alphaGo](#), [others?](#)), robotics ( [CITE ddpg](#), [others?](#)), and optimization requiring exploration ( [CITE datacenter power](#), [NAS](#), [chemical discovery](#)). Multi-agent reinforcement learning has many additional challenges, including multi-agent credit assignment, curse of dimensionality, and non-stationarity in learning dynamics [7].

MARL falls on a range from independent learning (IL) to joint action learning (JAL). IL agents treat other agents as part of the stochastic (and likely non-stationary) environment, which performs poorly. JAL agents condition on all agents' action and observations spaces, but don't scale because the action and observation spaces grow exponentially with  $n$ . Decentralized MARL solves this by using separate policies for each agent without conditioning on all agents' information, where it is assumed that agents receive partial observations about the environment and other agents.

The typical approach to decentralized MARL is to use agent communication: at every time step, agents send information needed by partners across a communication channel. A common approach is for agents to model each other so that less information must be shared.

In this thesis I will explore additional ways to overcome uncertainty in decentralized MARL. In particular, I propose to extend a class of decentralized MARL methods that utilize Centralized Training with Decentralized Execution (CTDE). In MADDPG, agents model each others’ deterministic policies  $\mu : O \mapsto A$ . Agents may not update without knowing other agents’ observations. When agent communication is unavailable, agents may utilize a predictive model of the most likely observations of all missing agents. During centralized training, a context-conditional generative model is trained to fill missing observations from the full joint observation for each given time step. I will also investigate the effect of using a centralized vs. decentralized generative model (therefore fully decentralized), as well as the effect of recurrent connections in several of the models being used (value and policy functions, and the generative model). We also will explore alternatives to the current formulation of the generative model as a CC-WGAN operating on the joint observation space, such as Bayesian methods (MCMC or variational inference), as well as independent marginal GANs that model each agent’s observation distribution and the relationship between agents’ distributions as a mixture model.

We also believe this context-based modeling is useful for domain adaptation / sim-to-real transfer, for which the CTDE methods are ill-suited to address directly. Agents must adapt to the target environment if deployed, however the basic CTDE methods cannot do real updates unless they use communication or some sort of inference like my proposed generative model. As such, a recurrent generative model will provide context-based history-dependent predictions that should quickly adapt to the target domain to provide useful inference to allow decentralized learning without explicit communication in a real-world scenario.

Additionally, if we allow explicit communication then the proposed context-based generative model may be used to decide when to communicate vs. use inference (as the communication network may be lossy, limited bandwidth, slow or spotty). Essentially, if we quantify the reward improvement by communicating (and get real observations) vs. sampling from the generative model (and get approximate observations), what is the behavior of the system? Intuitively, the system should occasionally choose to communicate with (non-independent) agents so that the error from inference does not outweigh the cost of communication (which must be defined).

## Chapter 1

### INTRODUCTION

The introduction should be broad in topic and conversational in tone. It could introduce a whole area rather than a single problem. Another reason to develop a substantial introduction is that a thesis is a more thorough, detailed document than is a paper. Why was the problem worth investigating in depth? How do the parts of the investigation relate to each other? What are some practical, concrete ways in which the outcomes of the work might be used? Running examples may be outlined in the introduction, to give unity to the thesis overall.

Reinforcement learning (RL) with function approximation has been used to solve difficult sequential decision making problems in high dimensional state and action spaces, such as game playing [11] and robotics [6]. Many decision making problems are best modeled as a multi-agent system in which agents learn concurrently with other agents. Two naive approaches that use single-agent RL methods in multi-agent problems are independent learning (IL) and joint action learning (JAL), but these approaches perform poorly. IL agents treat all others as part of the stochastic environment, ignoring relevant information about them. JAL agents condition on the full joint action and observation spaces for all agents, but these joint spaces grow exponentially with the number of agents and are therefore not scalable.

Most multi-agent reinforcement learning (MARL) methods rely on decentralized policies where each agent's own policy is dependent on local observations and

actions. Decentralized scenarios usually have partial observations and limited communication. In some problems, the learning must also be decentralized and rely on agent communication [17]. In other cases, a simulator may be available so that agents may learn with extra state information while assuming free, instantaneous communication. After centralized learning, the agents execute decentralized policies using only local information. This training procedure is often feasible for software systems or physical domains with dedicated simulators, such as robotics or autonomous vehicles. Many recent works have adopted this Centralized Training, Decentralized Execution (CTDE) framework; we review several such methods in the following section.

Agent communication is the main approach to learning in decentralized systems when agents have partial observations. In CTDE methods, the communication during centralized training is implicit: all observations are shared freely and instantly. However, communication networks may be lossy, delayed, or lacking coverage.

We propose to address the limited communication learning gap of CTDE methods in decentralized execution with generative modeling. Specifically, we use a modified context conditional generative adversarial network (CC-GAN) [2] to infer missing joint observations given partial observations. The task of filling in partial observations by generative inference is similar to the image inpainting problem for a missing patch of pixels: with an arbitrary number of missing observations, we would like to infer the most likely observation of the other agents.

We will extend the popular MADDPG method [10] as it appears most amenable to full decentralization. MADDPG agents require both (1) the policies of other agents and (2) other agents’ observations as input to the policies and critics. As other agents’ approximate policies can be learned, the agents only need to learn a model for the other agents’ observations. The generative model will learn this joint distribution

of agent observations by training on random combinations of missing agent communications during centralized training. During the decentralized portion, the agents may sample from this model to continue learning under partial observability.

When the decentralized environment is dissimilar to the centralized environment, both agent policies and the generative model require fine-tuning to the new dynamics. In our system, this is worsened by the co-adaptation of the generative inference and inferred observations seen by agents. We provide experimental results to show how reward differs when the decentralized environment is partially observable vs. both partially observable with altered dynamics, as well as our method’s performance with several options for decentralized updates.

Our contributions are as follows. We review the recent trend of CTDE MARL literature and identify that they are ill-suited to learn in the decentralized execution phase without explicit communication. We show how a context conditional generative model can address this problem for a popular CTDE method and provide a modified GAN that can learn the joint observation distribution. We experimentally evaluate our approach on three continuous multi-agent tasks. To the best of our knowledge, this is the first work to use generative models to overcome multi-agent partial observability or address decentralized learning in CTDE methods.



## Chapter 2

### BACKGROUND AND RELATED WORK

This is the information for this chapter.

#### 2.1 Background

##### 2.1.1 Reinforcement Learning

Formally, each task in decentralized MARL is represented by a discrete-time partially observable Markov Game [9], a multi-agent extension of the Markov decision process [15]. A Markov Game is a tuple  $\langle \mathcal{S}, \mathcal{A}, \mathcal{O}, \mathcal{R}, \mathcal{T}, n, \gamma \rangle$  where a set of  $n$  agents choose actions based on local observations to maximize their own expected cumulative reward. At each time step  $t$ , the environment has a true state  $s \in \mathcal{S}$  and each agent  $i$  simultaneously chooses an action  $a_i$  from their individual set of available actions  $\mathcal{A}_i \in \mathcal{A}$ . The environment stochastically transitions to a new state  $s'$  given by the state transition function  $\mathcal{T} : \mathcal{S} \times \mathcal{A} \mapsto \mathcal{S}$ , and each agent then receives a reward  $r_i$  according to its own reward function  $\mathcal{R}_i : \mathcal{S} \times \mathcal{A} \mapsto \mathbb{R}$ . The discount factor  $\gamma$  is used for calculating expected return  $R_i = \sum_{t=0}^T \gamma^t r_i^t$  for time horizon  $T$ .

In the partially observable setting typical in decentralized MARL, each agent receives a private observation  $o_i \in \mathcal{O}_i$  correlated with the state  $s$ . Agents choose actions using a stochastic policy  $\pi_i : \mathcal{O}_i \times \mathcal{A}_i \mapsto [0, 1]$ , where  $\pi_i$  is a learned function with parameters  $\theta_i$ .

The three main approaches to RL are action-value methods, policy gradient methods, and the actor-critic hybrid approach [15]. Q-learning estimates the action-value function  $Q^\pi(s, a)$ : the future discounted reward when taking action  $a$  from state  $s$  while following policy  $\pi$ . Deep Q-Networks (DQN) used Q-learning with neural network function approximation to play Atari games from pixels to superhuman performance [11]. DQN also introduced two stability improvements: target  $Q$  functions that are updated less frequently, and an experience replay buffer that stores environment transitions  $(s, a, r, s')$  for decorrelated batch updates.

Instead of learning a value function, policy gradient methods learn a parameterized policy directly [15]. This approach is often more efficient but tends to have high variance. To reduce variance, actor-critic algorithms combine an action-value function  $Q$  along with the parameterized policy  $\pi$  to guide policy updates.

Policy gradient methods normally learn a stochastic policy  $\pi : \mathcal{S} \times \mathcal{A} \mapsto [0, 1]$ . Deterministic policy gradient (DPG) methods learn a policy  $\pi : \mathcal{S} \mapsto \mathcal{A}$  that returns a single action [14]. Deep DPG (DDPG) is an off-policy model-free actor-critic algorithm that combines the DQN value function with a deterministic policy [8]. DDPG is a popular single-agent RL method that works with continuous state and action spaces. Like DQN, DDPG uses experience replay and target networks for both value and policy networks. Random noise is added to the policy’s output for better exploration. From here on, all policies  $\pi$  are assumed deterministic.

### 2.1.2 Multi-Agent Reinforcement Learning

#### 2.1.3 Generative Models

Generative models learn a data distribution and can generate new samples similar to the learned distribution. The most popular class of models is the generative adversarial network (GAN) [4]. A GAN is composed of two neural networks

with opposing goals: a generator network  $\mathcal{G}$  that receives noise as input and produces samples similar to the data distribution, and a discriminator network  $\mathcal{D}$  that tries to determine real data points from those sampled from  $\mathcal{G}$ . While GANs have largely been applied to image generation, they should be able to learn any joint data distribution.

Wasserstein GANs (WGANs) are a variant of GANs that have been shown to have more reliable convergence and less mode collapse [1]. WGAN uses a critic rather than a discriminator (outputs are not probabilities), trains using a simple loss metric that approximates the Wasserstein distance when the network enforces a 1-Lipshitz constraint, and allows pre-training the critic to optimality. To avoid confusing the WGAN critic with the critic  $Q$ , we will continue to refer to it as the discriminator  $\mathcal{D}$  as this distinction makes no difference in our work. For sample batch  $\mathbf{x}$  and noise batch  $z \sim \mathcal{N}(0, 1)$  with batch size  $m$ , the WGAN discriminator and generator have the following losses:

$$\nabla_{\theta_{\mathcal{D}}} \frac{1}{m} \sum_{i=1}^m [\mathcal{D}(\mathbf{x}^{(i)}) - \mathcal{D}(\mathcal{G}(z^{(i)}))] \quad (2.1)$$

$$\nabla_{\theta_{\mathcal{G}}} \frac{1}{m} \sum_{i=1}^m \mathcal{D}(\mathcal{G}(z^{(i)})) \quad (2.2)$$

Our work uses the context-conditional generative model, where the model takes a partial input and must generate a complete data sample. The closest computer vision task to our problem is image inpainting, where a patch of pixels from an image is removed and the model must fill the missing patch based on its learned model of the pixels' joint distribution. The CC-GAN objective function is given by

$$\min_{\mathcal{G}} \max_{\mathcal{D}} \mathbb{E}_{\mathbf{x} \sim \mathcal{X}} [\log \mathcal{D}(\mathbf{x})] + \mathbb{E}_{\substack{\mathbf{x} \sim \mathcal{X}, \\ \mathbf{m} \sim \mathcal{M}}} [\log (1 - \mathcal{D}(\mathbf{x}_{\mathbf{I}}))] \quad (2.3)$$

where  $\mathbf{m}$  denotes a binary mask used to drop a patch from image  $\mathbf{x}$ , and  $\mathbf{x_I} = (1 - \mathbf{m}) \odot \mathbf{x_G} + \mathbf{m} \odot \mathbf{x}$  is the combined inpainted image where  $\odot$  is element-wise multiplication.

Other generative models for image inpainting include autoregressive models and context encoders, but they are not suitable for our approach. Autoregressive models, such as PixelRNN [16], require a pre-specified ordering over the pixels and thus will not work for arbitrary missing data. Context encoders use a variational autoencoder coupled with adversarial loss [12], but results tend to be less accurate compared to CC-GANs.

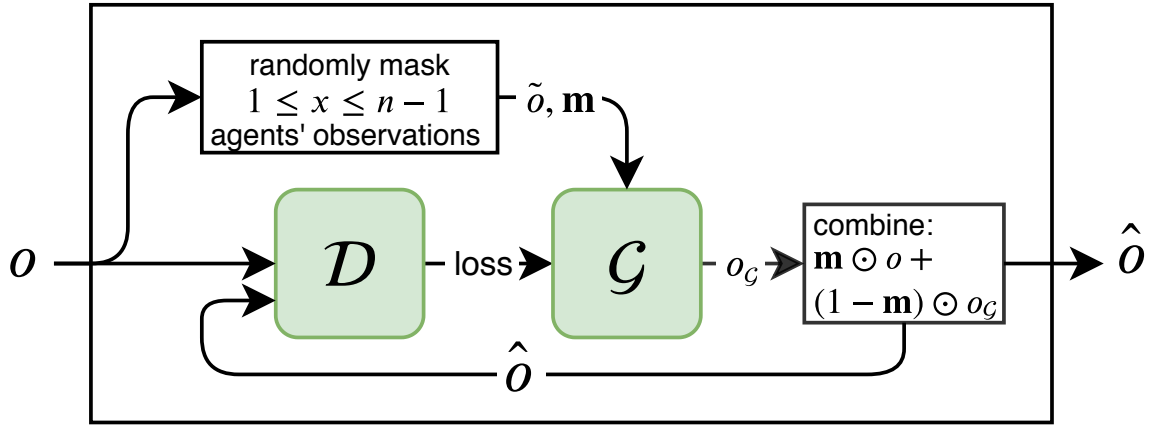
## 2.2 Related Work

## Chapter 3

### PROPOSED WORK (USE TITLE)

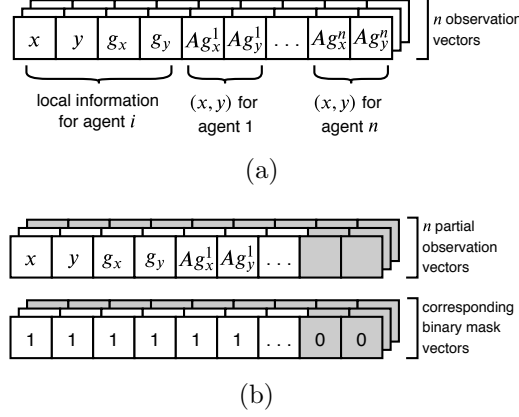
This is the information for this chapter.

#### 3.1 Context-based Generative Inference



**Figure 3.1:** CC-WGAN update diagram. Describes an agent-based communication training procedure similar to random masking found in image inpainting tasks.

We approach the problem of inferring missing information from partial observations as a generative sampling problem, similar to the task of image inpainting. We use a modified CC-GAN as our generative model [2]. Specifically, we train a WGAN with gradient penalty constraints [1, 5] with the CC-GAN random data masking training procedure. We refer to our modified model as the context-conditional



**Figure 3.2:** Random masking input and output used in training the CC-WGAN for an agent-distance partial observability criterion. Fig. 3.2(a) depicts  $n$  agents’ observation vectors fed to the masking function. Fig. 3.2(b) shows the masking function output: a set of  $n$  partial observations and  $n$  binary masks for agent  $i$ . Masked values in observations are filled with random normal values  $z \sim \mathcal{N}(0, 1)$ .

WGAN (CC-WGAN). In our experiments, the CC-WGAN was more reliable than regular CC-GAN for low-dimensional data.

Unlike the standard image generation task for GANs, we have no training data. We store joint observations  $o = \langle o_1, \dots, o_n \rangle$  in a replay buffer  $\mathcal{B}_G$  just as MADDPG to stabilize learning when batch training.

Fig. 3.1 shows the update procedure for training the CC-WGAN. When the model updates, it randomly samples joint observations  $o \sim \mathcal{B}_G$  from the joint observation replay buffer  $\mathcal{B}_G$  to randomly mask  $1 \leq x \leq n - 1$  agent observations. During centralized training, all joint observations are available. If the CC-WGAN is updated in the decentralized phase, inferred observations are mixed into the updates. This is a form of semi-supervised learning because the model updates on its own predictions [3].

For each joint observations  $o$ , we randomly mask combinations of missing agents from  $o$  with a binary mask  $\mathbf{m}$ . Masked elements in  $o$  are replaced with random normal noise  $z \sim \mathcal{N}(0, 1)$  to get the partial observation  $\tilde{o}$ . The masking procedure requires some knowledge of the conditions when observations will become partial, e.g., inter-agent distance greater than communications allow. Fig. 3.2 shows the masking procedure for distance-based partial observability based on  $(x, y)$  coordinates. In the diagram,  $\mathcal{G}$  takes joint partial observation  $\tilde{o}$  and binary mask vector  $\mathbf{m}$ , and produces the generated output  $o_{\mathcal{G}} = \mathcal{G}(\tilde{o}, \mathbf{m})$ . We then replace the masked portion in  $o$  with that portion of the generated output  $o_{\mathcal{G}}$  to get a combined observation  $\hat{o} = \mathbf{m} \odot o + (1 - \mathbf{m}) \odot o_{\mathcal{G}}$ .

Where CC-GAN passes only the inpainted patch to the discriminator, this assumes fixed size image patches. Since any number of agents may be missing from the joint observation  $\hat{o}$ , we instead pass the combined observation to the discriminator.  $\mathcal{D}$  is trained to distinguish batches of size  $m$  of real joint observations  $o$  and inferred observations  $\hat{o}$  by minimizing the empirical Wasserstein distance:

$$\frac{1}{m} \sum_{i=1}^m [\mathcal{D}(o^{(i)}) - \mathcal{D}(\hat{o}^{(i)})] \quad (3.1)$$

where  $\hat{o} = \mathbf{m} \odot o + (1 - \mathbf{m}) \odot \mathcal{G}(\tilde{o}, \mathbf{m})$ . Similarly,  $\mathcal{G}$  is updated by maximizing:

$$\frac{1}{m} \sum_{i=1}^m \mathcal{D}(\hat{o}^{(i)}) \quad (3.2)$$

### 3.1.1 MADDPG with Inferred Observations

As stated before, we augment MADDPG [10] with the CC-WGAN because it appears the most flexible CTDE method for decentralization. Each MADDPG agent  $i$  learns a deterministic policy  $\pi_i$ , a centralized critic  $Q_i$ , and a set of approximate policies  $\mu_i^j$  for each other agent  $j$ .

Fig. 3.3 shows the MADDPG method updates along with the CC-WGAN for both centralized and decentralized phases. During centralized training, the critics  $Q_i$  and policies  $\pi_i$  are updated exactly as MADDPG. In addition, the CC-WGAN is collecting joint observations in its replay buffer and updating as described above.

In the decentralized phase, local observations may be missing information about other agents. At each time step each agent  $i$  receives a partial observation  $\tilde{o}_i$  which consists of the agent’s local information and possibly information about other agents. When updating the centralized critics, if an agent  $i$  has information about agent  $j$  in its local partial observation  $\tilde{o}_i$ , then it can also see agent  $j$ ’s partial observation  $\tilde{o}_j$ . This is because we assume agents within range are “communicating” all local information. Just as in training, the joint partial observation  $\tilde{o}$  is passed to the generator to get  $o_G$  and combined via binary mask  $\mathbf{m}$  with  $\tilde{o}$  to get the inferred observation  $\hat{o}$ .

Following the derivation in [10], the deterministic policy loss with inferred observations is:

$$\nabla_{\theta_i} J(\pi_i) = \mathbb{E}_{\hat{o}, a \sim \mathcal{B}} [\nabla_{\theta_i} \pi_i(a_i | \hat{o}_i) \nabla_{a_i} Q_i^\pi(\hat{o}, a)], \quad (3.3)$$

where  $a = a_1, \dots, a_n$  are taken from approximate policies such that  $\mu_i^j(\hat{o}_j) = a_j$ . The approximate policies are updated with:

$$\mathcal{L}(\theta_i^j) = -\mathbb{E}_{\hat{o}_j, a_j} [\log \pi_i^j(a_j | \hat{o}_j) + \lambda H(\pi_i^j)] \quad (3.4)$$

where  $H(\pi_i^j)$  is the entropy of the policy distribution and  $\lambda$  is a small weight (0.001 in experiments). The centralized critics  $Q_i$  are updated with:

$$\mathcal{L}(\theta_i) = \mathbb{E}_{\mathbf{x}, a, r, \mathbf{x}'} [(Q_i^\mu(\mathbf{x}, a_1, \dots, a_N) - y)^2], \quad (3.5)$$

$$y = r_i + \gamma Q_i^\mu(\hat{o}', \mu_i^1(\hat{o}_1), \dots, \mu_i^N(\hat{o}_N)) \quad (3.6)$$



where  $\mu'_i$  is a target policy and  $\hat{o}'_i$  is an inferred next observation following observation  $o$ .

### 3.1.1.1 Environments and Setup

We evaluate our method under three continuous scenarios, two competitive and one cooperative, of the Multi-agent Particles Environment (MPE)<sup>1</sup>. MPE was introduced in the original MADDPG paper [10]. In order to directly compare to their results, we use the physical deception and predatory-prey competitive scenarios also used by [10]. We additionally test on a cooperative navigation scenario where agents share a reward function to show both competitive and cooperative settings. The communication-based scenarios tested by Lowe et al. have no clear way to make partially observable and still be solvable so we did not include them. In contrast, the physical scenarios we evaluate on are easy to make partially observable by using a partially observable distance  $d_P$ : when agents' are farther from each other than  $d_P$ , they cannot observe each others' positions, velocity, etc. (see Fig. 3.4). We use  $d_P = 1$  in all experiments, where the width of the 2D square environment is 2.

When using agent-distance partial observability, learning the coordination of predator-prey appears hardest, followed by physical deception, and lastly cooperative navigation. In predator-prey, three agents must coordinate to catch one faster agent, so there is no stable strategy with limited view. In physical deception, two agents should learn to deceive an adversary agent by covering two landmarks to hide which is the correct goal. If the adversary is out of range, this strategy should not change. In cooperative navigation, each agent must move to and remain near different landmarks. With limited range observations, an agent can still tell if another

---

<sup>1</sup> MPE code: <http://github.com/openai/multiagent-particle-envs>

agent is covering the same landmark within distance  $d_P$  and move to a different one.

In addition to making the decentralized phase partially observable, we approximate real-world deployment by changing the transition function by modifying simulation dynamics. Dynamics are changed by adding scaled random normal noise and translation to both actions and observations. Combined with partial observability, the added noise makes learning more difficult for both the policies and the CC-WGAN and requires fine-tuning to the new distribution.

Each episode has 200 steps with no early termination. All models are updated every 100 steps, and are represented with a three-layer, feed-forward neural network with 64 hidden units. The models use the Adam optimizer and each non-output layer uses a ReLU activation function. Each plot shows the mean and standard deviation shading over 30 independent trials for each algorithm. Each algorithm within a single plot receives the same set of 30 random seeds for accurate comparisons with random exploration. In all plots, a vertical dashed line marks the episode in which the environment becomes decentralized.

### 3.1.1.2 Results

The following plots compare our approach of augmenting MADDPG with CC-WGAN inference against regular MADDPG and DDPG. We chose DDPG because it performed the best among all IL methods in [10] and we use the same environment. Agents learn approximate policies for all other agents in MADDPG with and without generative inference. MADDPG and our version are identical during the centralized training because agents only use the CC-WGAN inference in the decentralized phase. After centralized training, we let MADDPG continue learning while treating the partial observability as random noise, whereas our approach infers the

missing data. Except for Fig. 3.7, the policies and CC-WGAN continue updating in the decentralized phase.

We did not test against the other CTDE methods discussed in Related Work because these methods use recurrent critics or policies. As such, they condition on past trajectories and would likely overcome the partial observability implicitly. Since MADDPG may also be used with recurrent units, we would like to compare recurrent MADDPG with and without generative inference against some of the recurrent CTDE methods in future work.

In Fig. 3.5, we show the cooperating agents’ reward for all agents using MADDPG, MADDPG with CC-WGAN inference, and DDPG. In this plot we only introduce partial observability when agents are farther than the partially observable distance  $d_P = 1$ .

Fig. 3.6 shows the reward for cooperating agents with both partial observability and altered environment dynamics in the decentralized phase to evaluate the capability of fine-tuning to another environment. The overall reward is much lower here than Fig. 3.5 which suggests the CC-WGAN is not well-suited to switching its modeled observation distribution. Deploying into a real world scenario is an important application, which is addressed by sim-to-real transfer [13].

In Fig. 3.7, we compare the total reward for our method with four options of decentralized updates, where the decentralized phase has both partial observability and altered environment dynamics. All agents use the CC-WGAN inferred observations when choosing actions. The curves show the difference in whether the policy and CC-WGAN update on inferred observations in the decentralized phase.

Lastly, in Fig. 3.8, we show the CC-WGAN’s reconstruction MSE during training over several partial observability distances  $d_P \in [0.0, 2.0]$ . When  $d_P = 0$ , the model is effectively using IL in the decentralized phase; when  $d_P = 2$ , the model

is usually using complete observations. We show this reconstruction plot because the agents’ observations are low-dimensional (i.e., not images as GANs are usually used for). MSE may not be a good metric for this error, however it was more informative than cosine similarity. We initially expected the error to be lower for larger  $d_P$ , but it is clear that the CC-WGAN reconstruction error has little to do with the partial observability distance. We would like to investigate the conditions under which the CC-WGAN gives better predictions.

### 3.1.1.3 Discussion

The results shown here reveal properties about context-based modeling of observations in MARL and the scenarios under which our current version is appropriate.

As CC-WGAN learns a joint observation distribution by sampling joint observations from its replay buffer  $\mathcal{B}_G$ , it has no temporal coherence: its inference is independent from the previous step given the current context. Without a model of observation trajectories, it is ill-suited for dynamic tasks with no clear stable behaviors under the partial observability.

This problem can be seen in Fig. 3.5 depending on the scenario’s need for non-local information and the stability of optimal behavior. Inferring other agents’ observations is useful when the task requires non-local coordination like the physical deception and predator-prey scenarios. Cooperative navigation agents can move to a different landmark if another agent is covering the same landmark, but may take slightly longer. Without temporal coherence, the CC-WGAN has trouble modeling non-stationary observation distributions like predator-prey. In contrast, physical deception and cooperative navigation have a stable optimal policy. In summation, our approach works best with a stable observation distribution (physical deception

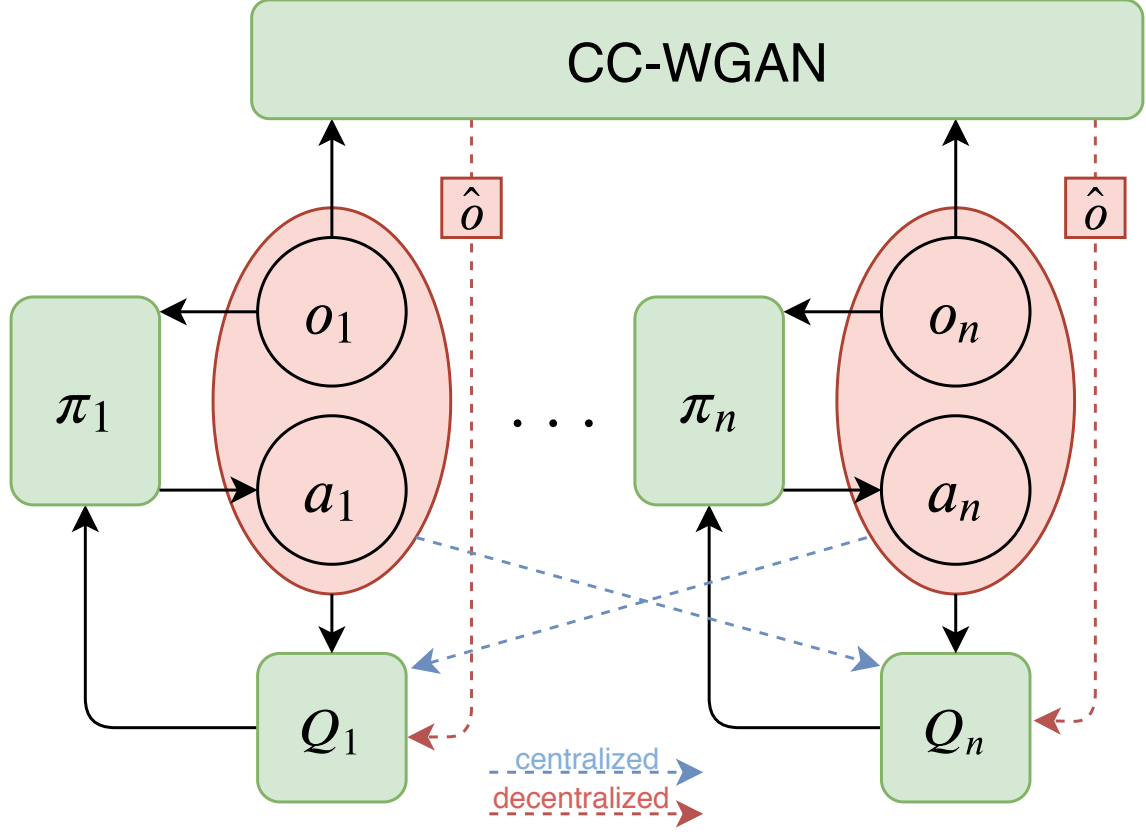
and cooperative navigation) and is useful in tasks requiring non-local coordination (physical deception and predator-prey). This is why our reward is significantly higher in physical deception, slightly higher in cooperative navigation, and slightly lower in predator-prey.

As seen in Fig. 3.7, when the CC-WGAN and policy update on inferred observations the decentralized reward dropoff is more drastic than without updating on the inference. This is due to the co-adaptation between the agents' policies and CC-WGAN inference: the CC-WGAN must learn based on new observations being generated by agents choosing actions based on the inferred observations from the CC-WGAN. Also it appears that having either policy updates or GAN updates on inferred observations gives roughly the same benefit.

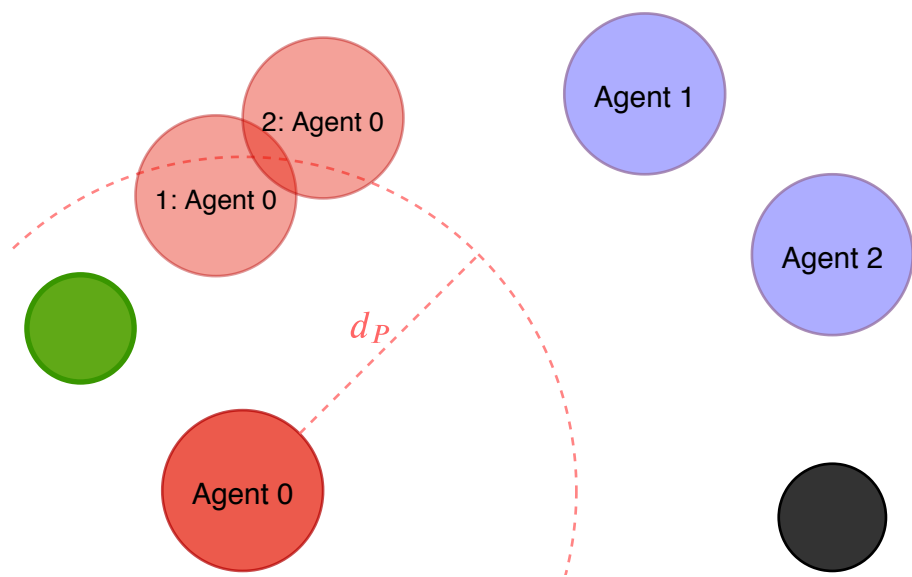
### **3.2 Context-history-based Generative Inference**

### **3.3 Domain Adaptation for CTDE Approaches**

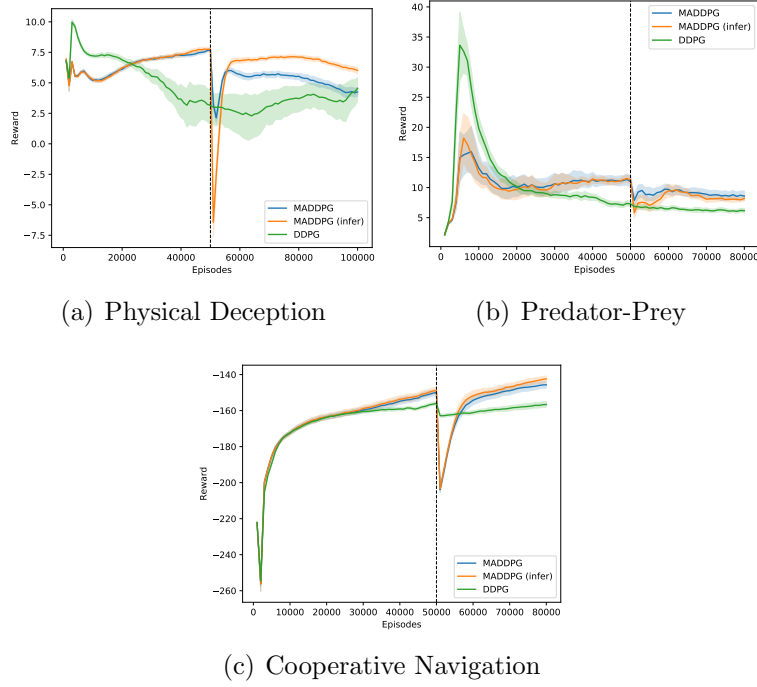
### **3.4 Learning when to communicate (vs. sample)**



**Figure 3.3:** MADDPG with CC-WGAN learning diagram. During centralized training, the actor-critic learning updates are identical to MADDPG and CC-WGAN collects joint observations. During decentralized execution, agents sample from CC-WGAN to fill partial observations. Note that this diagram excludes the approximate policies  $\mu$ ; only the observations are shared during centralized training.

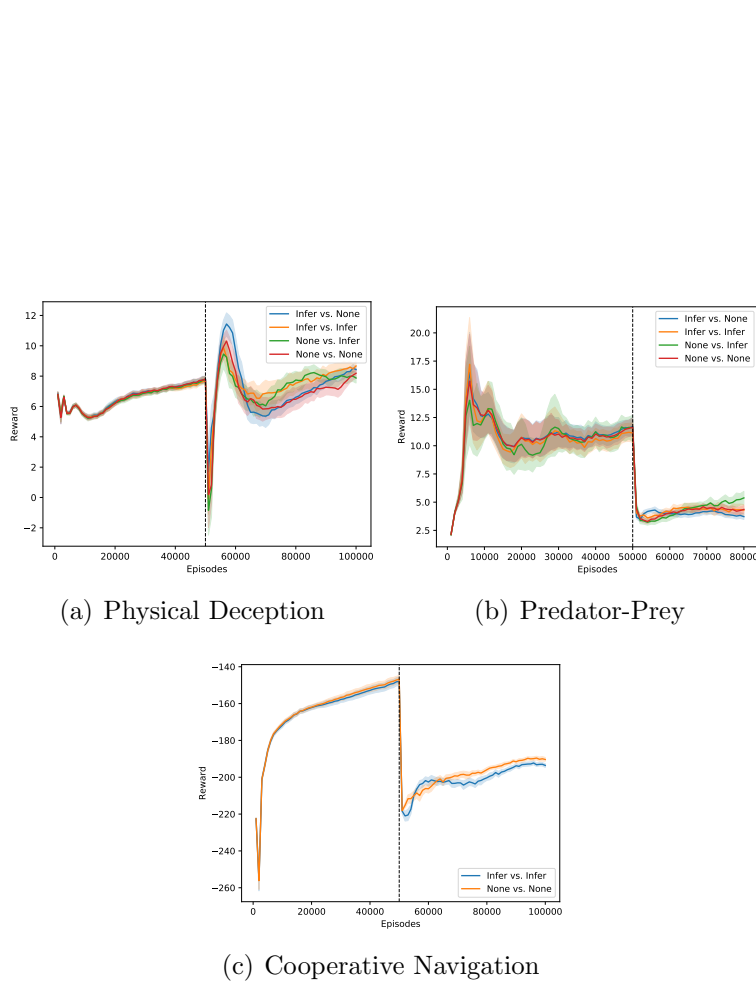


**Figure 3.4:** Physical deception scenario with agent inference at the beginning of an episode. Agents 1 and 2 are near each other and can therefore observe each others' real positions. They are both too far away from Agent 0 (with partially observable distance  $d_P$ ) so they each sample from the generative model to infer Agent 0's position.

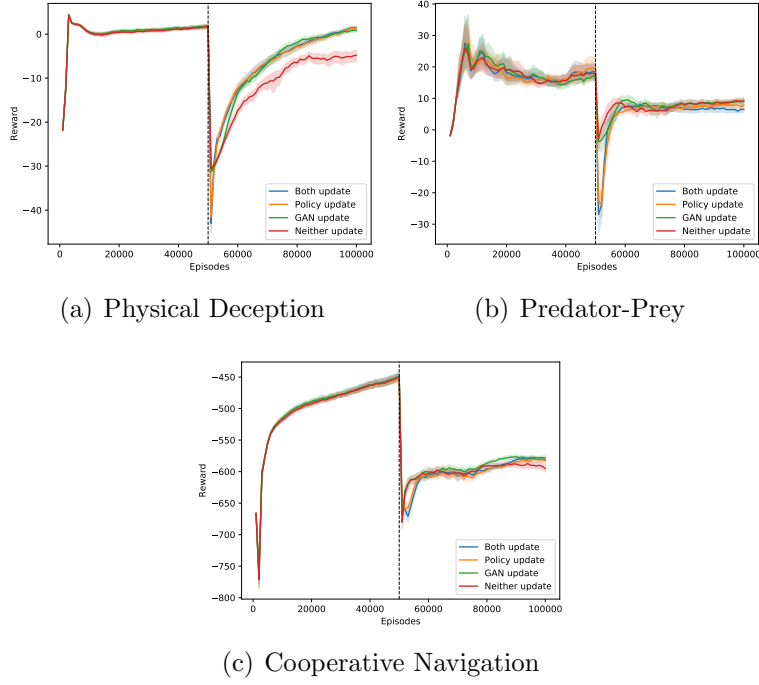


**Figure 3.5:** Reward for cooperating agents with distance-based partial observations in decentralized phase.

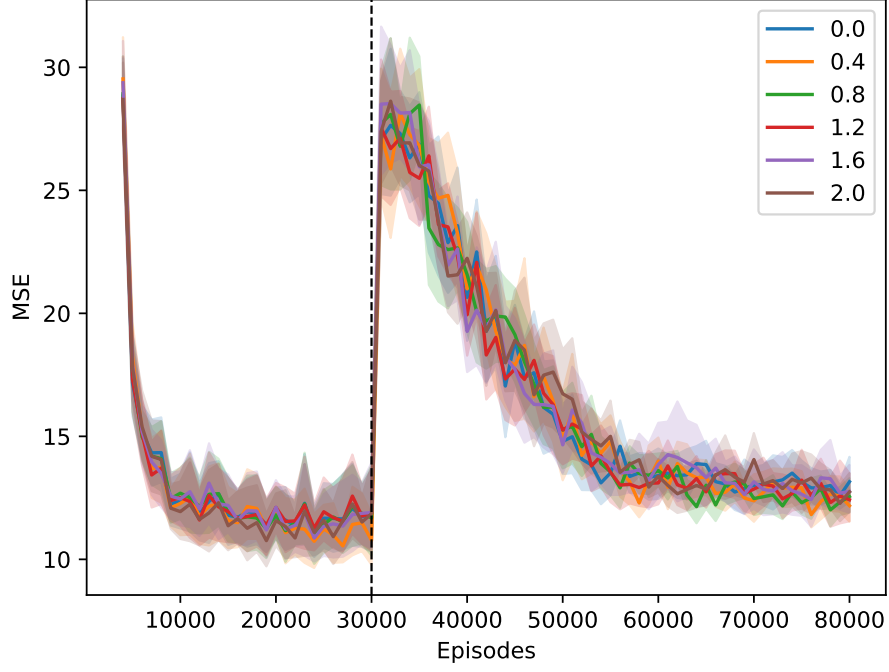




**Figure 3.6:** Both sides of cooperating agent rewards with both distance-based partial observations and altering environment dynamics in decentralized phase. “Infer” is our approach and “None” is MADDPG.



**Figure 3.7:** Total reward for our approach with all four combinations of whether agent policies and CC-WGAN continue to update on inferred observations during decentralized phase.



**Figure 3.8:** Mean squared error CC-WGAN reconstruction between latest joint observations  $o$  and inferred observations  $\hat{o}$  for the predator-prey scenario with altered dynamics. Curves show varying partial observability distances  $d_P$  had little effect on error.

## **Chapter 4**

### **EVALUATION AND TIMELINE**

#### **4.1 Preliminary Results**

This is the information for the first section of the chapter.

#### **4.2 How to evaluate rest of proposed work**

#### **4.3 Timeline**

## BIBLIOGRAPHY

- [1] Martin Arjovsky, Soumith Chintala, and Léon Bottou. Wasserstein GAN. *arXiv*, 2017.
- [2] Emily Denton, Sam Gross, and Rob Fergus. Semi-Supervised Learning with Context-Conditional Generative Adversarial Networks. *arXiv*, 2016.
- [3] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. *Deep Learning*. MIT press, 2016.
- [4] Ij Goodfellow, J Pouget-Abadie, and Mehdi Mirza. Generative Adversarial Networks. *arXiv preprint arXiv: ...*, pages 1–9, 2014.
- [5] Ishaan Gulrajani, Faruk Ahmed, Martin Arjovsky, Vincent Dumoulin, and Aaron Courville. Improved Training of Wasserstein GANs. In *Advances in Neural Information Processing Systems*, pages 5767–5777, 2017.
- [6] Tuomas Haarnoja, Sehoon Ha, Aurick Zhou, Jie Tan, George Tucker, and Sergey Levine. Learning to Walk via Deep Reinforcement Learning. In *International Conference on Machine Learning*, 2018.
- [7] Pablo Hernandez-Leal, Bilal Kartal, and Matthew E Taylor. Is multiagent deep reinforcement learning the answer or the question? A brief survey. *arXiv*, 2018.
- [8] Timothy P Lillicrap, Jonathan J Hunt, Alexander Pritzel, Nicolas Heess, Tom Erez, Yuval Tassa, David Silver, and Daan Wierstra. Continuous control with deep reinforcement learning. In *International Conference on Learning Representations*, 2016.
- [9] Michael L Littman. Markov games as a framework for multi-agent reinforcement learning. In *Machine Learning*, 1994.
- [10] Ryan Lowe, Yi Wu, Aviv Tamar, Jean Harb, Pieter Abbeel, and Igor Mordatch. Multi-Agent Actor-Critic for Mixed Cooperative-Competitive Environments. In *Neural Information Processing Systems (NIPS)*, 2017.

- [11] Volodymyr Mnih, David Silver, and Martin Riedmiller. Playing Atari with Deep Reinforcement Learning. *arXiv*, pages 1–9, 2013.
- [12] Deepak Pathak, Philipp Krähenbühl, Jeff Donahue, Trevor Darrell, and Alexei A Efros. Context Encoders: Feature Learning by Inpainting. *arXiv*, 2016.
- [13] Xue Bin Peng, Marcin Andrychowicz, Wojciech Zaremba, and Pieter Abbeel. Sim-to-Real Transfer of Robotic Control with Dynamics Randomization. In *International Conference on Robotics and Automation*, 2018.
- [14] David Silver, Nicolas Heess, Thomas Degris, Daan Wierstra, and Martin Riedmiller. Deterministic Policy Gradient Algorithms. In *International Conference on Machine Learning*, 2014.
- [15] Richard S Sutton and Andrew G Barto. *Reinforcement Learning: An Introduction*. MIT press, 2nd edition, 2018.
- [16] Aaron van Oord, Nal Kalchbrenner, and Koray Kavukcuoglu. Pixel Recurrent Neural Networks. In *International Conference on Machine Learning*, pages 1747–1756, 2016.
- [17] Kaiqing Zhang, Zhuoran Yang, Han Liu, Tong Zhang, and Tamer Basar. Fully Decentralized Multi-Agent Reinforcement Learning with Networked Agents. In *International Conference on Machine Learning*, 2018.

## **Appendix A**

### **TITLE OF APPENDIX 1**

This is the information for the first appendix, Appendix A. Copy the base file, appA.tex, for each additional appendix needed such as appB.tex, appC.tex, etc. Modify the main base file to include each additional appendix file.

If there is only one appendix, then modify the main file to only use app.tex instead of appA.tex.

**Appendix B**  
**TITLE OF 2ND APPENDIX**

Text for 2nd appendix.