Adversarial Attacks on Reinforcement Learning Agents for Command and Control

James Z. Hare² Mark Mittrick² John Richardson² Ahaan Dabholkar¹ Nicholas Waytowich² Priya Narayanan² Saurabh Bagchi¹ ²DEVCOM Army Research Laboratory, USA ¹*Purdue Universitv*

Abstract

Given the recent impact of Deep Reinforcement Learning in training agents to win complex games like StarCraft and DoTA (Defense Of The Ancients) - there has been a surge in research for exploiting learning based techniques for professional wargaming, battlefield simulation and modelling. Real time strategy games and simulators have become a valuable resource for operational planning and military research. However, recent work has shown that such learning based approaches are highly susceptible to adversarial perturbations. In this paper, we investigate the robustness of an agent trained for a Command and Control task in an environment that is controlled by an active adversary. The C2 agent is trained on custom StarCraft II maps using the state of the art RL algorithms - A3C and PPO. We empirically show that an agent trained using these algorithms is highly susceptible to noise injected by the adversary and investigate the effects these perturbations have on the performance of the trained agent. Our work highlights the urgent need to develop more robust training algorithms especially for critical arenas like the battlefield.

1 Introduction

Deep Reinforcement Learning (DRL) has been successfully used to train agents in several tactical and real-time strategy games such as StarCraft [1] and DoTA [2], which involve complex planning and decision-making. These agents have demonstrated proficiency in coming up with winning strategies comparable to that of experienced human players (AlphaStar [3], OpenAI Five [4]) through techniques like selfplay, imitation learning, etc. As a result, in recent years, there has been mounting interest in the military research community in applying these RL techniques to tasks such as operational planning and command and control (C2). Simultaneously, traditional game engines have been repurposed to facilitate automated learning (pySC2 [5], SMAC [6,7], pyDoTA [8]) and new ones developed for battlefield simulation [9-13], creating what are effectively *digital wargames*. The driving force



Figure 1: Robustness Evaluation Methodology: The figure shows the difference between a benign (top) and malicious (bottom) environment at timestep t. Observations at t-1are input to a C2 agent that has been pretrained in a benign environment. The agent samples a suboptimal action as a result of the injected adversarial perturbations (orange) in the input which eventually leads to a loss for the BlueForce.

behind this research has been to improve and augment strategies used on the battlefields of the future, which are expected to be more complex and unconventional, possibly beyond the cognitive abilities of a human commander.

Recent works [14] have had considerable success in winning simulated wargames using C2 agents that have been trained through reinforcement learning techniques and synthetic data. This has been possible partly due to the scalability of RL training which has proved to be a massive advantage for exploring and exploiting different strategies when faced with difficult or complicated scenarios and only partial information about the environment. However, these evaluations are done in benign environments where information available to the C2 agent is assumed to be *uncorrupted*. Realistically, this is unlikely in battlefield situations as information there may have inherent noise because of the mode of collection (from sensors or other input source) or may be tampered with by enemy forces. In this work, we evaluate the robustness of such a trained agent when subject to potentially adversarial inputs in the context of C2.

In order to do so, first, we use the StarCraft II Learning Environment (SCLE [15]) to model conflict between two teams, the BlueForce and the RedForce. The C2 agent directs the BlueForce to win the battle by eliminating the RedForce troops. Next, we assume the an attacker present in the environment tampers with the observations collected from the battlefield before they are made available to the C2 agent. The added perturbations termed *adversarial perturbations* are constructed to be highly imperceptible to evade detection while maximally subverting the C2 agent's policy to something deleterious (Figure 1). We then evaluate the resulting drop in performance of the agent on several metrics as well as analyze the deviation in the course of action from a military perspective.

Our key contributions are summarized as follows:

• We empirically show the vulnerability of the trained C2 agent to even small adversarial perturbations in the input observations. Our studies quantify some expected trends and bring out some non-obvious trends. For example, our studies reveal that partially trained agents appear to be more resistant to noise than fully trained agents.

• For generality, we evaluate the effectiveness of the attacks on two different scenarios which correspond to the C2 agent's task to attack and defend respectively.

• We also evaluate agents trained using two state-of-the-art RL algorithms, A3C and PPO, and comment on their robustness to injected noise.

• We provide interpretability to the model's outputs by profiling the shift in action distributions predicted by the policy network caused by the attacker's perturbations.

Our evaluations demonstrate the susceptibility of vanilla RL training to adversarial perturbations and the need for robust training mechanisms and sophisticated detection and prevention techniques especially for such critical scenarios.

The structure of the paper is as follows. First we provide brief backgrounds on the use of RL for C2 followed by a description of the StarCraft Environment and two custom scenarios – TigerClaw and NTC in Section 2.1, which we use for training our agent. In Section 4, we describe the state and action space of our custom scenarios and the details of the RL agent. Section 4.2 and Section 5 contain the attack methodology and the evaluations respectively. Finally we include a discussion on the need for utilizing adversarially robust training techniques and directions for future work.

2 Background

2.1 RL for Command and Control (C2)

Mission success in military C2 requires disseminating control actions to designated forces based on real-time Intelligence,

Surveillance, and Reconnaissance (ISR) information of the operational environment, such as kinetic and non-kinetic (e.g., weather, political, economic, cultural) variables, and terrain information. Before the battle begins, the commander and their specialized staffing officers must develop a detailed mission plan encapsulated as a Course Of Action (COA). The development of a COA requires a detailed analysis of the operational environment, predictions of the opposing force's COA (strategy), and wargaming to identify a friendly force COA that is finely tuned to meet the mission requirements [16]. Typical military planning is solely based on the commander and staffing officers, and their allocated time before the battle dictates the number of possible COAs that can be considered. Additionally, each potential COA must be wargamed and fine-tuned against a small set of opposing force COAs to identify strengths and weakness. This can result in suboptimal (heuristic-based) COA [17].

To circumvent this limitation, future military planning is envisioned to incorporate an Artificial Intelligent (AI) commander's assistant that can generate and recommend COAs to aid in the military planning process. Recent developments in deep RL for strategy games provides a promising direction to develop control policies for C2 [3]. The algorithms developed allow an AI agent to learn the best control policy that optimizes a predefined reward function by playing millions of simulated games through the exploration of many environmental state and action pairs. To extend these formulations for military C2 requires modeling, simulating, and wargaming a large number of battles faster than real-time in a virtual environment that emulates realistic combat characteristics. Furthermore, the existing RL algorithms must be adapted to handle a large number of heterogeneous actors, doctrine based control strategies, complex state and action space, and uncertain information.

Previous work on RL for C2 found that the StarCraft II gaming engine, developed by Blizzard Entertainment [1], provides a simulation environment that can be militarized and used for prototyping an AI commander's assistant [14, 18]. The following subsections provide details of the StarCraft II environment and the scenarios considered in this paper.

2.2 StarCraft II C2 Environment

StarCraft II is a multi-agent real-time strategy game developed by Blizzard Entertainment [1] that consists of multiple players competing for influence and resources with the ultimate goal of defeating the other players. As a collaboration between Deepmind and Blizzard Entertainment, the *StarCraft II Learning Environment* (SC2LE) Machine Learning API was developed to allow researcher to study many difficult challenges associated with RL [5]. For example, controlling heterogeneous assets in a complex state and action space with uncertainties. Furthermore, the *StarCraft II Editor* allows developers to construct custom scenarios, making it possible to develop RL agents for C2.

Previous work extended this framework to militarize the SC2LE and develop a C2 simulation and experimentation capability that interfaces with deep RL algorithms via RL-lib [18], an open-source industry-standard RL library [19]. The icons were re-skinned to portray standard military symbology, new StarCraft II maps were designed to emulate realistic combat characteristics, such as environment terrain and asset attributes (e.g., visibility, weapons, weapons ranges, and damage), and a custom wrapper was designed to train RL agents. For a detailed description of this framework, see [18]. In our work, we use this framework to develop baseline RL agents for C2 and study the effects of adversarial attacks on the learned policies in the following two scenarios.

2.2.1 Custom StarCraft II Scenarios



Figure 2: TigerClaw Scenario



Figure 3: TigerClaw: (**Right**) The geographical map of the scenario and (**Left**) the correspondingly designed map in StarCraft

TigerClaw: The TigerClaw melee map (Figure 3) is a highlevel recreation of the TigerClaw combat scenario (Figure 2) developed using the StarCraft II map editor. The scenario was developed by Army subject-matter experts (SMEs) at the Captain's Career Course, Fort Moore, Georgia. The BlueForce is an Armored Task Force (TF) which consists of combat armor with M1A2 Abrams, mechanized infantry with Bradley Fighting Vehicles (BFV), mortar, armored recon cavalry with BFV, and combat aviation. The RedForce is a Battalion Tactical Group (BTG) with attached artillery battery and consists of mechanized infantry with BMP, mobile artillery, armored recon cavalry, combat aviation, anti-armor with anti-tank guided missiles (ATGM), and combat infantry.

As seen in Figure 2, the BlueForce is a supporting effort with a mission to cross a dry riverbed (Wadi) and defeat the defending RedForce in preparation for a forward passage of lines by the main effort. The terrain is challenging in this scenario because there are only two viable wadi crossing points (Figure 3).



Figure 4: NTC: (**Right**) The geographical map of the scenario and (**Left**) the correspondingly designed map in StarCraft

National Training Center (NTC): The NTC map is a representation of the Ft. Irwin and Bicycle Lake area as seen in Figure 4. The Blue and RedForce units are the same as Tiger-Claw, but there is also a capability to randomize the starting forces. RedForce initially begins in the Bicycle Lake region, while BlueForce is set up in defense of Fort Irwin. Thus the RedForce will maneuver from the Bicycle Lake area and attempt to destroy the BlueForce in the Ft. Irwin area. The goal of this map was to investigate the impact of new terrain and to determine if the BlueForce would adopt a defensive strategy.

In both SC2 maps, each side has been represented down to the platoon echelon. The speed, attack range, and damage attributes of SC2 units have been scaled to estimate the capabilities of the relevant platforms in order to facilitate reinforcement learning.

2.2.2 Reinforcement Learning Algorithms

Asynchronous Advantage Actor-Critic (A3C) [20]: A3C is an asynchronous version of the Advantage Actor-Critic (A2C) algorithm for RL that uses multiple agents to learn a policy and an estimate of the value function in parallel. At each timestep, the agent takes an action based on the current policy and receives a reward while transitioning to a new state. It then calculates the advantage function and updates its local copies of the actor-critic networks. This process takes place independently and asynchronously for each agent. A central parameter server that stores the global parameters of

the networks is updated periodically by each agent and is used to initialize the local parameters of the agents' local networks. The asynchronous nature of the updates leads to more efficient exploration of the environment and reduces the correlation between the updates leading to more stable and efficient learning. Typically, the policy and value functions are parameterized by a shared neural network with a softmax output for the policy and a linear output for the value function.

Proximal Policy Optimization (PPO) [21]: PPO is a policy gradient method that aims to improve the performance and stability of trust-region [22] methods by introducing a *clipped surrogate* objective function. This objective function effectively restricts policy change to a small range thus reducing the variability in training of the actor. During training, the PPO algorithm iteratively updates the actor as well as the action value function and state value function using the temporal-difference (TD) method. Further, this objective function enables PPO to guarantee monotonic improvements in the objective. This allows for faster convergence without strict constraints leading to more accurate and stable performance of the agent.

3 RL Environment

In this section we provide technical details of the RL environment used to train the C2 agent.

3.1 State and Action Space

The state space observed from the StarCraft II C2 environment consists of both a *screen* or visual representation and a *nonspatial* representation. The screen representation is a size 256×256 image of the the minimap that depicts the current environmental state. The value of each pixel provides the agent with an understanding of the BlueForce and RedForce units' positions along with terrain information. The nonspatial representation consists of a vector of nonspatial features of size 287 that encodes all of the game and unit information, such as unit type, health, position, game scores etc.

The action space within the StarCraft II environment is large since it is a combination of the following three components: the number of units available to the commander, the number of possible actions that each unit can execute, and the (x, y) pixel location within the size 256×256 minimap where the action will be executed. To reduce the action space, the StarCraft II C2 environment first restricts the number of units by defining *control groups*, which lumps common units together to reduce the overall number of units needed to be controlled. In our custom scenarios, the control groups are defined as – **①** BlueForce: "AVIATION", "MECH_INF", "MORTAR", "SCOUT" and "TANK" **②** RedForce: "ANTI_ARMOR", "ARTILLERY", "AVIATION", "INFANTRY" and "MECH_INF" Within each control group, we restrict the possible actions to be either "NO_OP" or "ATTACK (x, y)", where the function of the former is to essentially do nothing, while the latter moves the control group (i.e., all associated units) to a desired (x, y) location and attacks any enemy within its firing range along the way. Additionally, the number of (x, y) pixel locations is reduced by segmenting the minimap into nine disjoint quadrants with locations defined as "LEFT", "CENTER", and "RIGHT" for the *x*-axis, and "TOP", "CENTER", and "BOTTOM" for the *y*-axis, where the exact pixel location is the center of the quadrant.

3.1.1 Reward Structure

TigerClaw: The reward function for the TigerClaw map consists of -

• *Terrain Reward:* +10 points for each BlueForce unit crossing the Wadi (a dry river bed) and -10 points for retreating back.

• Attrition Reward: +10 points for destroying a RedForce unit and -10 points if a BlueForce unit is destroyed.

The terrain reward is meant to reward an offensive strategy for the BlueForce. It is meant to reinforce crossing of the Wadi to initiate conflict with the RedForce in order to take over desired locations. To drive this scenario, the RedForce has been scripted to defend the desired locations.

NTC: The NTC map shares the same attrition reward function as TigerClaw, but the terrain rewards are not included. That is,

• Attrition Reward: +10 points for destroying a RedForce unit and -10 points if a BlueForce unit is destroyed.

The lack of terrain reward is meant to encourage the Blue-Force to defend and focus on maximizing RedForce losses. To drive this scenario the RedForce has been scripted to seek and destroy the BlueForce.

3.2 RL Agent Description

3.2.1 The Policy Network

The C2 agent in this paper uses the same policy network as our previous work [14]. The network takes as inputs, three kinds of observations and the control group

- *screen* The screen representation discussed in Section 3.1 consisting of a vector of an image of size 256×256 .
- *nonspatial* The nonspatial representation discussed in Section 3.1 consisting of a vector of size 287.
- *action mask* The action mask is used to restrict the action space to the allowed actions as described in Section 3.1.



Figure 5: **C2 Policy Network:** Computaional graph of the policy network of our C2 agent. The inputs and outputs are shown in blue and yellow respectively. Shaded rectangles represent the *concatenate* operation. Conv2D and FC layers are ReLU activated.

• *control group* - A one-hot encoding of the selected control group (specified in Section 3.1) that will be the focus of action prediction by the policy network. The encoding is then concatenated to the nonspatial input vector. At each time step, the control group is sequentially selected to identify their next actions.

The output of the network is an approximation of the value function and an 8 element vector arranged as (*action-logits*, *x-logits*, *y-logits*). The *action-logits* is a 2 element vector that determines the action. The *x*,*y-logits* are 3 dimensional entries each with the logit values corresponding to the positions (LEFT, CENTER, RIGHT) and (TOP, CENTER, BOTTOM) on the map, respectively. This output vector is used to create a probability distribution over the action space from which the next action is sampled by the agent.

4 Adversarial Attacks on RL agents

Prior works on the robustness of RL training have focused on evaluating the algorithms from a perspective of sensitivity to environment dynamics [23, 24] or the ability to train adaptive adversarial policies against them [25, 26]. In this work, we focus on the former approach. A lot of previous research has shown that neural network predictions are highly sensitive to perturbations in their input space [27–29]. As DRL approaches typically rely on parameterizing policies with neural networks, they suffer from the same vulnerabilities.

4.1 Adversarial Attacks on Image Classifiers

An adversarial perturbation is a small perturbation that is added to a benign input to fool a trained network into predicting an incorrect output. Typically the perturbation is constrained to be small enough to avoid detection. For example, in the case of image classifiers, the ideal adversarial perturbations in pixel space would be imperceptible to a human observer but would cause the neural network classifier to predict an incorrect class with high confidence. Formally, for a trained classifier f(w;.) and an input image x, computing an adversarial sample $x' = x + \sigma$ involves computing σ such that $f(w,x') \neq f(w,x)$ while minimizing d(x',x), where d is a distance metric such as Euclidean distance.

4.1.1 Fast Gradient Sign Method (FGSM) [28]

FGSM provides an efficient method for generating adversarial samples given whitebox access to the model. Given a trained classifier model f(w; .), a first order attacker can generate adversarial perturbations for the benign sample x by first computing the gradient of the classifier loss (L) with respect to the input. The weighted perturbation when added to x, creates the adversarial sample x'.

$$x' = x + \varepsilon \cdot \operatorname{sgn}(\nabla_x L(f(w; x), y_t))$$

where y_t is the ground truth label of the sample x and ε is the perturbation budget that controls the amount of distortion to the original sample. In the case of FGSM, it bounds the l_{∞} norm of the perturbation added to the original image. Intuitively, the attack moves the sample x in the direction of $\nabla_x L(f(w;x), y_t)$ which maximizes the classifier loss L. In this work, we prefer FGSM to more powerful attacks such as PGD [29] and C&W [27] because of its lower computational cost, which leads to a more efficient attack.

4.2 Inference Time Attacks on Policies

Inference time attacks are used against pre-trained agents which are deployed in the wild. The goal of the attack is to induce the trained policy network into predicting a suboptimal action distribution by surreptitiously perturbing the input observations. This inevitably leads to the agent losing its expected reward.

In the wild, such an attack could be realized as a cyberattack where sensors collecting data from the battlefield could be compromised and transmit corrupted data. A high level overview of the attack is presented in Figure 1.

4.2.1 Threat Model

First we assume a C2 agent that has been trained in a benign environment and has learned an optimum policy. This agent is subsequently deployed in an unsafe setting to direct BlueForce troops in a battle against the RedForce. Next, we consider an attacker that has the ability to intercept and modify observations coming from the environment before they are received by the C2 agent to select the next actions. Through



(a) PPO/TigerClaw: Episode Re-(b) PPO/TigerClaw: Reward trend (c) PPO/NTC: Episode Rewards (d) PPO/NTC: Reward trend w.r.t wards over 100 episodes w.r.t perturbation budget (ϵ). over 100 episodes perturbation budget (ϵ).

Figure 6: Inference time attack on an agent in the TigerClaw (a),(b) and NTC (c),(d) scenario trained using PPO.



(a) A3C/TigerClaw: Episode Re-(b) A3C/TigerClaw: Reward trend (c) A3C/NTC: Episode Rewards (d) A3C/NTC: Reward trend w.r.t wards over 100 episodes w.r.t perturbation budget (ϵ). over 100 episodes perturbation budget (ϵ).

Figure 7: Inference time attack on an agent in the TigerClaw (a),(b) and NTC (c),(d) scenario trained using A3C.

such modifications, the attacker hopes to influence the agent to sample an incorrect action that leads to a low reward for the BlueForce. To generate such input perturbations efficiently, we also assume that the attacker has white-box access to the policy network used by the C2 agent.

4.2.2 Attack Methodology

We use FGSM as a basis for generating perturbations at inference time and use the modification presented in [23] to target policies instead of classifiers. Unlike supervised learning, we do not know the "ground truth" action at any given timestep in RL and we assume that the action predicted by the policy network with the highest likelihood is optimal.

It should be noted that since the output of our policy network has three different components (action logits, x-logits, y-logits), ideally the ground truth vector should be computed for each component separately. With this in mind we construct the ground truth vector with a *degenerate distribution*. That is, for an input observation x, for each component of the output, we take the element with the highest value in the vector $y = [y^{(1)}, y^{(2)}, y^{(3)}] = f(w; x)$ and assign weight 1 to it's logit value and 0's everywhere else. In other words,

$$y_i^{\prime(j)} = \begin{cases} 1, & \text{if } y_i^{(j)} = \max(y^{(j)}) \\ 0, & \text{otherwise} \end{cases}$$

where y_i is the i^{th} element of vector y.

However, we find that treating the components as a single unit and calculating the degenerate distribution on the entire output is more efficient while making for an effective attack. After making this relaxation, the malicious perturbation is then calculated for *x* using the gradient $\nabla_x L(f(w;x), y')$, where *L* is the cross-entropy (CE) loss function. We use this technique to craft perturbations for two of the three input components namely – the *screen* and *nonspatial* components. We do not perturb the action mask as it is simply used to mask out invalid actions predicted by the network. As a result, perturbing this component does not conform to any realistic setting.

5 Evaluations

5.1 Experimental Setup

Our experiments were conducted on a cluster node with two AMD Epyc 7763 "Milan" CPUs 128 cores with 256GBs of main memory. We used the SC2 framework described in Section 2.2 to train our RL agents. Each training run was performed using 90 workers, each occupying a single core (90 parallel instances of the game) and consumed about 170GBs of memory. The PPO and A3C agents were trained for 5K/25K iterations respectively corresponding to around 40M/100M timesteps or 134K/376K episodes respectively. The rollouts were visualized using a custom pygame interface.

5.2 Evaluation of Inference Time Attacks on Agent Reward

For evaluating effectiveness of inference time attacks, first we train our C2 agent in a benign environment until it learns a policy that consistently achieves a high reward on the given scenario. The trend of the attained reward when the agent is deployed in the presence of an attacker is then studied over the course of a 100 episodes or *rollouts*.

For a comprehensive evaluation, we consider two scenarios or tasks – **1** TigerClaw (Attacking BlueForce) and **2** NTC (Defending BlueForce). We also consider two different state-of-the-art training algorithms, A3C and PPO. In both scenarios the C2 agent controls the movements of the Blue-Force while the RedForce follows a fixed policy. A detailed description of these scenarios was given in Section 2.2.

To further analyze and understand the shift in agent strategy we use several quantitative metrics as well as observe multiple rollouts. Our insights are presented below.

5.2.1 Vulnerability to Adversarial Perturbations

Figures 6 and 7 show the resulting reward trends for an agent trained using PPO and A3C, respectively. We present boxplot statistics aggregated over 100 episodes, for the rewards attained by the agent under attack with different perturbation levels in Figures 6b and 6d (PPO) and Figures 7b and 7d (A3C). $\varepsilon = 0.0$ corresponds to the benign case when no perturbations are made. In most cases, we observe a steep decrease in the median reward even for minute perturbations ($\varepsilon = 0.05, 0.08$). Increasing ε also shows diminishing returns for attack effectiveness.

To maintain the secrecy of the attack, the perturbation budget needs to be small enough to be imperceptible to a human auditor, especially for the screen component of the input. Figure 8 shows a visual comparison of perturbations of the screen component at different levels. From our evaluations we observe that $\varepsilon = 0.1$ is optimal in lowering the reward while keeping perturbations to a minimum.



Figure 8: Visualizing perturbations of the screen component at different ε levels in the TigerClaw scenario. Higher ε levels result in a greater amount of noise.

Figure 9 provides explainability into the utility of the attack by showing how it changes the actions taken by the agent. We first plot the action distribution predicted by the policy network on benign observations collected at certain timesteps during a PPO/TigerClaw rollout. These are shown in Figure 9a, 9c, and 9e, respectively. We then compare it to the distribution predicted by the same policy network after maliciously perturbing the observations with our attack for $\varepsilon = 0.1$ in Figures 9b, 9d, and 9f, respectively. The probabilities are plotted on a log scale. As can be observed the actions with the highest likelihood shift from $(1,1,2) \rightarrow (1,2,1), (1,1,2) \rightarrow$ $(0,2,1) \& (0,0,0) \rightarrow (1,2,1)$ at timestep 10, 21, and 100.



Figure 9: Shift in the action distributions predicted by the policy network at different timesteps during a PPO/TigerClaw rollout. (Left) Action distribution on benign observations. (**Right**) Action distribution predicted by network after perturbing observation with $\varepsilon = 0.1$. Each caption represents the action sampled with highest likelihood.

5.2.2 Analyzing Agent Behavior under Attack

We observe a definite change in the strategy used by the C2 agent under an attacker's influence. Over a number of rollouts in either scenario, we frequently observe artifacts like erratic troop movements where the BlueForce troops keep oscillating about a single position for certain number of timesteps, straying off course of the aviation units which are crucial for winning the TigerClaw scenario, etc. Similarly in the NTC scenario, the attacks cause the BlueForce to retreat towards the bottom of the map, away from the advancing RedForce as opposed to following the original strategy – to



Figure 10: Sampled actions of a PPO agent over 5 TigerClaw rollouts. We can observe a shift in the distribution in the presence of an attacker.

aggressively pursue and eliminate the RedForce units.

To better understand this shift in behavior, we compare actions taken by the agent in a benign environment to those taken by the same agent when an adversary is present. That is, at each timestep we compute the actions that would be taken by the agent if the observations were maliciously modified and compare these "subverted actions" to the actions that are actually taken by the agent at that timestep. We do so by perturbing the observations using our attack (with $\varepsilon = 0.1$) and using the agent to predict the next action. The normalized action frequencies are plotted in Figure 10 for a PPO/TigerClaw agent over 5 episodes. A clear difference is observed in the actions taken in the two cases. Actions (0,2,1), (1,1,1), (1,2,1)are sampled frequently when perturbed observations are presented to the agent. These actions misdirect the BlueForce as they target areas on the map that do not contain any RedForce troops.



Figure 11: Actions taken by the BlueForce during a TigerClaw rollout in a benign setting compared to those taken in the presence of an inference time (FGSM) attacker.

In order to examine the BlueForce movements more thoroughly, we select one episode and for each timestep plot the subverted actions and the actual actions taken (Figure 11) by the agent. A rollout in a benign environment sees the BlueForce enter conflict and destroy most of the RedForce troops in the first 100 timesteps. This corresponds to the actions (1,1,2), (1,1,1), (1,1,0) which are taken frequently within the first 50 timesteps. The (0,1,1) actions correspond to NO_OPs that are taken after exiting conflict when most of the RedForce are killed. In the presence of an attacker, however, in the first 50 timesteps the perturbations cause the C2 agent to sample actions like (1,2,1) with high frequency. As a result, the BlueForce are misdirected and end up being killed by the RedForce.



Figure 12: Comparing the impact of inference time attack on additional game metrics for a PPO agent in TigerClaw. The results are aggregated over 100 rollouts.

This is also reflected in the casualty metrics of both teams in Figure 12a. There are greater BlueForce casualties and a fewer RedForce casualties in the presence of an $\varepsilon = 0.1$ attacker. To get an idea of the impact of each unit on their respective team, we look at the health remaining percentage of each control group at the end of an episode. Figures 12c and 12d shows the health statistics for the BlueForce and RedForce respectively, aggregated over 100 episodes. We can see in Figure 12c that the health remaining of the *aviation* units drops significantly to 2% in the presence of an attacker. Rollouts show that these units are critical to winning the TigerClaw scenario for the BlueForce. On the other hand the RedForce *aviation* units actually see an increase in remaining health (77%). Other groups of the BlueForce also see a noticeable decrease in their health metrics.

Similarly in the NTC scenario, the health remaining for the RedForce increase greatly on attack (Figure 13d), while the BlueForce health left only decreases marginally. Observing multiple rollouts reveals the attacks cause the BlueForce units to get misdirected and as a result, they do not get into conflict with the RedForce. This is also supported by the casualty metrics for the two teams (Figure 13a) with the BlueForce only having a slightly higher number of casualties when compared to the benign case.

Finally, the attack causes a sharp drop in the number of games that end with higher health levels for the BlueForce than the RedForce (*termed a partial win*). This can be seen in Figures 12b and 13b.



Figure 13: Comparing the impact of inference time attack on additional game metrics for a PPO agent in NTC. The results are aggregated over 100 rollouts.

5.2.3 Strength & Reliability of the Attack

For small perturbations, the attack is not very reliable as evidenced by the large dispersion in attained reward. High variance can be observed in Figures 6a, 6c, 7a, and 7c where we plot the (EMA smoothed) episode reward over 100 episodes. This can be explained partly by the stochastic nature of action sampling and the reward structure associated with the map. This dispersion is reduced when increasing ε , leading to greater reliability but trading off the secrecy of the attack. A surprising instance is that of an A3C/TigerClaw agent where the attack remains ineffective even for larger perturbations. We analyze this in more detail below.

Attacking the A3C/TigerClaw agent: As can be seen in Figure 7b, the attack appears to fail on the C2 agent trained using A3C on the TigerClaw map. Even for large perturbations ($\varepsilon > 0.1$), the attack fails to degrade the agent's reward as extensively as the other cases. Interestingly, we note that even in a benign environment, the A3C/TigerClaw agent achieves a significantly lower reward than the PPO/TigerClaw agent.

First, to rule out the training algorithm as a potential factor and for a fair comparison, we partially train a PPO agent on the TigerClaw scenario until it achieves a similar reward to the A3C agent and study the rewards achieved by this agent in the presence of an attacker.

Our observations are shown in Figure 14. In Figure 14a we see that the partially trained PPO agent (*PPO-partial*) gets similar rewards as the A3C agent over 100 rollouts. Notably, its reward trend when attacked (Figure 14b) is comparable to that of the A3C agent (Figure 7b). This empirically shows



(a) Episode rewards of the agent (b) PPO-partial/TigerClaw: Reunder benign conditions. ward trend w.r.t ε

Figure 14: Comparing the effect of the inference time attack on the A3C agent with the partially trained PPO agent (PPOpartial). This agent was trained for 1M timesteps.

that the attack is largely unaffected by the training algorithm. Coupled with the previous observation, we hypothesize that the effectiveness of the attack is correlated to the quality of the trained agent.



Figure 15: Action distribution shifts over 5 episodes for A3C/TigerClaw and PPO-Partial/TigerClaw

To test this hypothesis in greater detail we turn to the frequencies of actions taken by both agents over multiple rollouts. Similar to the analysis in previous sections, to visualize the shift induced by the attack we plot the subverted and actual ($\varepsilon = 0.1$) actions sampled over 5 episodes in Figure 15 for both agents. Surprisingly, we observe an almost complete overlap of the two action plots for the A3C agent (Figure 15a) and only a marginal difference for the PPO-partial agent. This indicates that the attack is in most cases incapable of flipping the actions sampled by the agent.

As the malicious perturbations are constructed using a loss gradient, we try to understand the loss landscape of the poorly trained policy networks to explain the attack's ineffectiveness. To do so, we plot the loss value for 10^4 randomly sampled points in the l_{∞} -ball of ε radius around a fixed observation O_c . The loss is computed as the component-wise sum of the CE loss between the predicted action distribution on that observation and a fixed ground-truth action distribution. The ground truth action distribution is computed corresponding to the actual action taken by the agent on O_c (Section 4.2). The CE loss is calculated for each output component (action



(a) PPO/TigerClaw Agent mean re-(b) PPO/NTC Agent mean reward (c) A3C/TigerClaw Agent mean (d) A3C/NTC Agent mean reward ward comparison.

Figure 16: Component wise impact of the input on mean episode reward.



Figure 17: Loss Value plots for the agent policy networks. The *y* axis represent the normalized frequencies (10^4 trials) for the loss value on the *x* axis.

logit, x-value, y-value) and summed to get the final loss. Figure 17 shows the results for $\varepsilon = 0.1$ for both the A3C and PPO-Partial agents when compared to the PPO/TigerClaw agent as a baseline. Compared to the PPO agent, both the A3C and PPO-Partial agent's predictions are highly similar to the ground truth as evidenced by the highly frequent ~ 0 loss value. This corresponds to a *flatness* in the prediction space of both (A3C and PPO-Partial) policy networks where inputs in neighbourhood result in the same prediction. Consequently this leads to greater robustness to injected noise – benign and adversarial as small perturbations are not enough to significantly change the predicted action distributions. Investigating the relation between this perceived robustness and the quality of training is left as future work.

5.2.4 Component-wise Impact of Input on the Attack

In this section we analyze the impact of the screen and nonspatial components on the effectiveness of the attack. To do so, we perform the attack by restricting malicious perturbations to only the screen or nonspatial components respectively. We track the mean reward achieved by the C2 agent in each case and compare them to the baseline attack which perturbs both components. We record the rewards gained over 100 rollouts and show the results in Figure 16. Figures 16a, 16b, and 16d indicate that for smaller perturbations levels ($\varepsilon < 0.5$), modifications to the screen component are primarily responsible for a successful attack. Further, we only see a minor degradation in the reward for small nonspatial-only perturbations in both scenarios. On the other hand, the drop in mean reward because of screen-only perturbations is similar to that in the baseline attack. Intuitively, this skewed importance makes sense as the screen component is a much larger input component (256×256 versus 287) when compared to the nonspatial component. The screen component further encodes important spatial information on the terrain and location of troops that has a significant impact on the policy network's prediction.

We also note the anomalous trends in Figure 16c for the A3C/TigerClaw agents. It appears that the nonspatial component is more important to the policy network, however since this agent is poorly trained, we cannot conclusively reason about the resulting trends.

5.2.5 Adversarial Robustness of the Training Algorithm

To compare the robustness of an agent trained using PPO to one trained using A3C we compare the trend followed by the *mean relative reward* (R_r) defined as the ratio of the mean reward obtained by an agent under attack (R_a) to the mean reward earned in a benign environment (R_b). That is, $R_r = R_a/R_b$. A larger value of R_r signifies more robustness. Figure 18 presents the results of our comparisons. In both cases, we see that the A3C agent seems to perform marginally better than the PPO agent for small perturbations and worse for higher perturbation levels. As in the previous section, for a fair comparison we compare the partially trained PPO agent (PPO-partial) with the A3C agent (Figure 18a) and see a similar decreasing trend for R_r .

While we cannot conclusively comment on the robustness of A3C vs PPO, the robustness of the PPO-partial agent compared to a fully trained PPO agent makes for a compelling argument. This has been explored in more detail in the preceding section. It should be noted however that PPO does seem to be functionally more performant than A3C as evidenced



(a) Relative Rewards of PPO vs (b) Relative Rewards of PPO vs A3C on TigerClaw A3C on NTC

Figure 18: Comparison of Algorithm Robustness: Mean Relative Reward with respect to the perturbation budget

by the higher rewards in both scenarios.

6 Discussions and Future Work

In this work we have focused on evaluating the robustness of RL agents used for C2 through the lens of an inference time attacker. Our investigations supported by existing literature reveal that even the latest RL training algorithms cannot be trusted to train agents that can be reliably deployed in unsafe environments. Our evaluations show that well-trained agents are highly sensitive to even minute perturbations in their input space and act suboptimally as a result. In an arena such as the battlefield where observations received can be prone to noise either benign or malicious, this raises serious questions on the use of such agents.

This directly leads us to two avenues for future work. The first is developing explainable and controllable approaches to COA generation using RL training and developing robustness mechanisms for RL agents that work both during inference and training time. Such mechanisms can be a combination of adversarial perturbation detection and prevention mechanisms that can be deployed on top of pre-existing agents. The second avenue for future work is existing training algorithms can be augmented to train agents that are certifiably robust to malicious noise. In this vein, adversarial training offers a promising alternative to train robust agents. In this work, while we present a preliminary study relating the flatness of the loss landscape of the policy network to it's apparent robustness, further research is required to quantify the susceptibility of a policy network to noise.

7 Conclusions

Our evaluations reveal the fragile nature of vanilla RL agents trained for C2 when deployed in insecure environments where even minute perturbations to the input introduced by a malicious actor are sufficient to introduce a large variability in the agent's prediction. We analyze this susceptibility of a C2 agent when commanding the BlueForce in two custom scenarios and discuss the reasons behind such behavior and the implications from a strategic perspective. Finally, we emphasize the need to develop robust training algorithms for RL which would be critical for reliable mission planning on battlefields of the future.

Funding

This material is based in part upon work supported by the Army Research Lab (ARL) under Contract number W911NF-2020-221. Any opinions, findings, and conclusions or recommendations expressed in this material are those of the authors and do not necessarily reflect the views of the sponsors.

References

- [1] Blizzard, "Starcraft ii." [Online]. Available: https://starcraft2.blizzard.com
- [2] Valve, "Dota 2." [Online]. Available: https://www.dota2. com/home
- [3] O. Vinyals, I. Babuschkin, W. M. Czarnecki, M. Mathieu, A. Dudzik, J. Chung, D. H. Choi, R. Powell, T. Ewalds, P. Georgiev, J. Oh, D. Horgan, M. Kroiss, I. Danihelka, A. Huang, L. Sifre, T. Cai, J. P. Agapiou, M. Jaderberg, A. S. Vezhnevets, R. Leblond, T. Pohlen, V. Dalibard, D. Budden, Y. Sulsky, J. Molloy, T. L. Paine, C. Gulcehre, Z. Wang, T. Pfaff, Y. Wu, R. Ring, D. Yogatama, D. Wünsch, K. McKinney, O. Smith, T. Schaul, T. Lillicrap, K. Kavukcuoglu, D. Hassabis, C. Apps, and D. Silver, "Grandmaster level in StarCraft II using multiagent reinforcement learning," *Nature*, 2019.
- [4] OpenAI, :, C. Berner, G. Brockman, B. Chan, V. Cheung, P. Dębiak, C. Dennison, D. Farhi, Q. Fischer, S. Hashme, C. Hesse, R. Józefowicz, S. Gray, C. Olsson, J. Pachocki, M. Petrov, H. P. d. O. Pinto, J. Raiman, T. Salimans, J. Schlatter, J. Schneider, S. Sidor, I. Sutskever, J. Tang, F. Wolski, and S. Zhang, "Dota 2 with large scale deep reinforcement learning," 2019.
- [5] google deepmind, "pysc2." [Online]. Available: https: //github.com/google-deepmind/pysc2
- [6] M. Samvelyan, T. Rashid, C. S. de Witt, G. Farquhar, N. Nardelli, T. G. J. Rudner, C.-M. Hung, P. H. S. Torr, J. Foerster, and S. Whiteson, "The starcraft multi-agent challenge," 2019.
- [7] oxwhirl, "smacv2." [Online]. Available: https://github. com/oxwhirl/smacv2
- [8] pydota2, "pydota2." [Online]. Available: https://github. com/pydota2/pydota2

- [9] P. Narayanan, M. Vindiola, S. Park, A. Logie, N. Waytowich, M. Mittrick, J. Richardson, D. Asher, and A. Kott, "First-year report of arl directors strategic initiative (fy20-23): artificial intelligence (ai) for command and control (c2) of multi-domain operations (mdo)," US Army Combat Capabilities Development Command, Army Research Laboratory, 2021.
- [10] S. J. Park, M. M. Vindiola, A. C. Logie, P. Narayanan, and J. Davies, "Deep reinforcement learning to assist command and control," in *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications IV.* SPIE, 2022.
- [11] S. Soleyman and D. Khosla, "Multi-agent mission planning with reinforcement learning," in AAAI Symposium on the 2nd Workshop on Deep Models and Artificial Intelligence for Defense Applications: Potentials, Theories, Practices, Tools, and Risks. AAAI, 2020.
- [12] L. Zhang, J. Xu, D. Gold, J. Hagen, A. K. Kochhar, A. J. Lohn, and O. A. Osoba, "Air dominance through machine learning," *Santa Monica, CA: RAND Corporation*, 2020.
- [13] A. Basak, E. G. Zaroukian, K. Corder, R. Fernandez, C. D. Hsu, P. K. Sharma, N. R. Waytowich, and D. E. Asher, "Utility of doctrine with multi-agent rl for military engagements," in *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications IV.* SPIE, 2022.
- [14] N. Waytowich, J. Hare, V. G. Goecks, M. Mittrick, J. Richardson, A. Basak, and D. E. Asher, "Learning to guide multiple heterogeneous actors from a single human demonstration via automatic curriculum learning in starcraft ii," in *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications IV*. SPIE, 2022.
- [15] O. Vinyals, T. Ewalds, S. Bartunov, P. Georgiev, A. S. Vezhnevets, M. Yeo, A. Makhzani, H. Küttler, J. Agapiou, J. Schrittwieser *et al.*, "Starcraft ii: A new challenge for reinforcement learning," *arXiv preprint arXiv:1708.04782*, 2017.
- [16] J. J. Marr, The military decision making process: Making better decisions versus making decisions better. School of Advanced Military Studies, US Army Command and General Staff College, 2001.
- [17] W. A. Shoffner, *The Military Decision Making Process: Time for a Change*. School of Advanced Military Studies, US Army Command and General Staff College, 2000.

- [18] V. G. Goecks, N. Waytowich, D. E. Asher, S. Jun Park, M. Mittrick, J. Richardson, M. Vindiola, A. Logie, M. Dennison, T. Trout *et al.*, "On games and simulators as a platform for development of artificial intelligence for command and control," *The Journal of Defense Modeling and Simulation*, vol. 20, no. 4, 2023.
- [19] E. Liang, R. Liaw, R. Nishihara, P. Moritz, R. Fox, K. Goldberg, J. Gonzalez, M. Jordan, and I. Stoica, "Rllib: Abstractions for distributed reinforcement learning," in *International conference on machine learning*. PMLR, 2018.
- [20] V. Mnih, A. P. Badia, M. Mirza, A. Graves, T. Lillicrap, T. Harley, D. Silver, and K. Kavukcuoglu, "Asynchronous methods for deep reinforcement learning," in *International conference on machine learning*. PMLR, 2016.
- [21] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov, "Proximal policy optimization algorithms," *arXiv preprint arXiv:1707.06347*, 2017.
- [22] J. Schulman, S. Levine, P. Abbeel, M. Jordan, and P. Moritz, "Trust region policy optimization," in *International conference on machine learning*. PMLR, 2015.
- [23] S. Huang, N. Papernot, I. Goodfellow, Y. Duan, and P. Abbeel, "Adversarial attacks on neural network policies," *arXiv preprint arXiv:1702.02284*, 2017.
- [24] J. Sun, T. Zhang, X. Xie, L. Ma, Y. Zheng, K. Chen, and Y. Liu, "Stealthy and efficient adversarial attacks against deep reinforcement learning," in *Proceedings of the AAAI Conference on Artificial Intelligence*, 2020.
- [25] X. Wu, W. Guo, H. Wei, and X. Xing, "Adversarial policy training against deep reinforcement learning," in *30th USENIX Security Symposium (USENIX Security* 21), 2021.
- [26] A. Gleave, M. Dennis, C. Wild, N. Kant, S. Levine, and S. Russell, "Adversarial policies: Attacking deep reinforcement learning," *arXiv preprint arXiv:1905.10615*, 2019.
- [27] N. Carlini and D. Wagner, "Towards evaluating the robustness of neural networks," in 2017 IEEE Symposium on Security and Privacy (SP), 2017.
- [28] I. J. Goodfellow, J. Shlens, and C. Szegedy, "Explaining and harnessing adversarial examples," *arXiv preprint arXiv:1412.6572*, 2014.
- [29] A. Madry, A. Makelov, L. Schmidt, D. Tsipras, and A. Vladu, "Towards deep learning models resistant to adversarial attacks," in *International Conference on Learning Representations*, 2018. [Online]. Available: https://openreview.net/forum?id=rJzIBfZAb