
Adversarial Cheap Talk

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Adversarial attacks in reinforcement learning (RL) often assume highly-privileged
2 access to the learning agent’s parameters, environment or data. Instead, this
3 paper proposes a novel adversarial setting called a *Cheap Talk MDP* in which an
4 Adversary has a *minimal range of influence* over the Victim. Parameterised as a
5 deterministic policy that only conditions on the current state, an Adversary can
6 merely *append* information to a Victim’s observation. To motivate the minimum-
7 viability, we prove that in this setting the Adversary cannot occlude the ground truth,
8 influence the underlying dynamics of the environment, introduce non-stationarity,
9 add stochasticity, see the Victim’s actions, or access their parameters. Additionally,
10 we present a novel meta-learning algorithm to train the Adversary, called adversarial
11 cheap talk (ACT). Using ACT, we demonstrate that the resulting Adversary still
12 manages to influence the Victim’s *training and test* performance despite these
13 restrictive assumptions. Affecting train-time performance reveals a new attack
14 vector and provides insight into the success and failure modes of existing RL
15 algorithms. More specifically, we show that an ACT Adversary is capable of
16 *harming* performance by interfering with the learner’s function approximation
17 and *helping* the Victim’s performance by appending useful features. Finally, we
18 demonstrate that an ACT Adversary can append information during train-time to
19 directly and arbitrarily control the Victim at test-time in a zero-shot manner.

20 1 Introduction

21 Learning agents are often trained in settings where adversaries may have some control over part of the
22 agent’s observations. However, the adversary cannot usually influence the dynamics of the underlying
23 environment or the reward signal (at least not without cost). For example, it is often possible to
24 append arbitrary tags to content later used to train recommender systems. Similarly, an adversary
25 could rent space on interactive bulletin boards near busy traffic intersections to influence data sets
26 used for training self-driving cars. Another instance occurs in financial markets, where an adversary
27 can change the state of the order-book by submitting orders far *out of the money*. While all of these
28 are examples of *useless* features from an information point of view, under the current paradigm of
29 end-to-end deep learning it is common practice to include a superset of useful features as part of the
30 input and to let the model learn which features matter. This paper demonstrates that an actor can
31 heavily influence the behaviour and performance of learning agents by controlling information *only*
32 in these “useless” features.

33 Most past work in adversarial attacks assumes that the adversary can influence the environment
34 dynamics [13, 11]. For example, perturbing images and observations could obscure or alter relevant
35 information, such as the ball’s location in a Pong game [15]. Furthermore, many attacks require

36 access to the trained agent’s weights and parameters to generate the adversarial inputs [27]. Finally,
37 most of these attacks only cause the victim’s policy to fail arbitrarily instead of giving the adversary
38 full control over the victim’s policy at test time [12, 14, 22, 1].

39 In contrast, our work, inspired by recent advancements in the field of opponent shaping [20], in
40 Section 2 proposes a novel, minimum-viable setting to shape a learning agent, called “Cheap Talk
41 MDP”. In this setting, the Adversary can only *append* information to the observation of a Victim
42 as a deterministic function of the current state. The Adversary does not have access to the Victim’s
43 parameters, actions, or even samples from the Victim’s policy. In Section 3, we prove that the
44 Adversary cannot change the dynamics of the underlying environment nor alter the reward functions.
45 Nor can it inject stochasticity into the environment (deterministic) or introduce non-stationarity
46 (function of the current state only). Furthermore, we prove that Adversaries *cannot* influence *tabular*
47 Victims in this setting in Proposition 1; Adversaries can therefore only interfere with a Victim through
48 their function approximator. In this sense, our setting represents a *bare minimum* range of influence,
49 as further justified in Appendix B.3.

50 In Section 3, we also introduce a new meta-learning algorithm to train the Adversary, called adversarial
51 cheap talk (ACT). With an extensive set of experiments in Section 4, we demonstrate that an ACT
52 Adversary can influence a Victim to achieve a number of outcomes:

- 53 1. We show that the ACT Adversary can prevent the Victim from solving a task, resulting
54 in low rewards *during training*. We provide empirical evidence that the Adversary sends
55 messages which induce *catastrophic interference* in the Victim’s neural network.
- 56 2. Conversely, an ACT Adversary can learn to send useful messages that *improve* the Victim’s
57 training process, resulting in higher rewards *during training*.
- 58 3. Finally, we introduce a training scheme that allows the ACT Adversary to arbitrarily control
59 the Victim *at test-time*, in a *zero-shot* manner.

60 **Related Work and Background** For an in-depth discussion on related work and background, we
61 point the reader to Appendix A.

62 2 Problem Setting

63 In this work, we consider two agents that interact in a setting we call a *Cheap Talk MDP*
64 $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mathcal{M}, f, \mathcal{J} \rangle$. Here \mathcal{M} denotes the space of messages. We refer to the agent observing
65 the message as the *Victim* with transition and reward functions \mathcal{P}, \mathcal{R} independent from \mathcal{M} . The
66 agent appending the message is called the *Adversary*, endowed with a deterministic policy (function)
67 $f : \mathcal{S} \rightarrow \mathcal{M}$ to append messages and an objective function \mathcal{J} to optimise (details below).

68 The Victim is a “standard” reinforcement learning agent, selecting actions according to its policy
69 $a_t \sim \pi_\theta(\cdot \mid s, f(s))$, where $a \in \mathcal{A}, s \in \mathcal{S}$. The Victim optimises its policy π_θ with respect to
70 parameters θ , to maximise its expected return J defined in Equation ??.

71 By contrast, the Adversary may only act by appending a message $f_\phi(s)$ to s at every step, where
72 $f_\phi : \mathcal{S} \rightarrow \mathcal{M}$ is a deterministic policy (function) of the current state and ϕ are the Adversary’s
73 parameters. These parameters may only be updated *between* full training/testing episodes of the
74 Victim; the function remains static during episodes to avoid introducing non-stationarity. The
75 Adversary’s objective function \mathcal{J} may be picked arbitrarily, and need not be differentiable if it is
76 optimised using ES.

77 In the train-time setting, we focus on the allied setting, where Adversary and Victim objectives are
78 equal, $\mathcal{J} = J$, and the adversarial setting where objectives are zero-sum, $\mathcal{J} = -J$. In the test-time
79 setting, we use an entirely different objective, such as reaching for an arbitrary circle in Reacher (see
80 Figure 3c). This incentivises the Adversary to manipulate the Victim into maximising \mathcal{J} , even if at
81 the cost of the Victim’s original objective J .

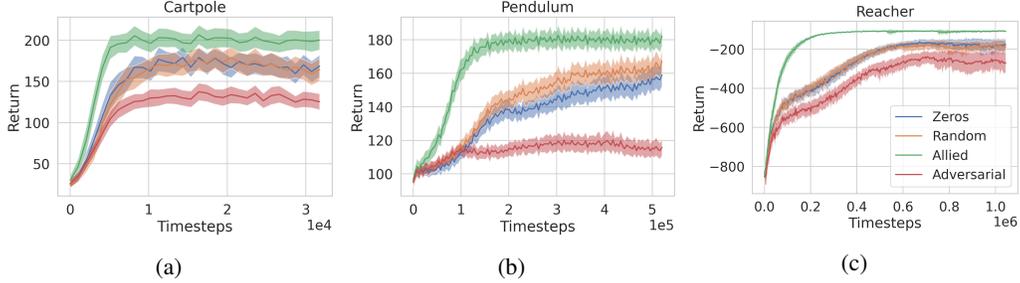


Figure 1: Visualizations of the training curves of the Victim across different number of dimensions of the messages for (a) Cartpole, (b) Pendulum, and (c) Reacher. Error bars denote the standard error across 10 seeds of Victims trained against a single trained Adversary.

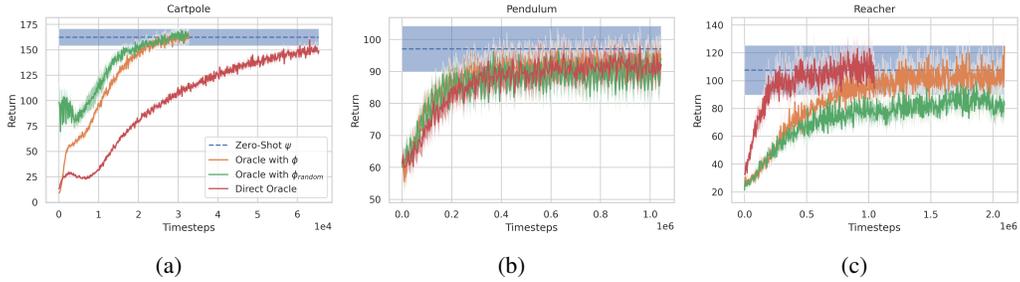


Figure 2: Training curves of the different agents in (a) Goal-Conditioned Cartpole (b) Goal-Conditioned Pendulum (c) Goal-Conditioned Reacher. The ablations show that the train- and test-time Adversaries learn near-optimal performance in comparison to the oracles. Error bars denote the standard error across 10 seeds of Victim trained against a single trained Adversary.

82 3 Method

83 **Meta-Training Procedure** Our method treats the problem setting as a meta-learning problem. The
 84 Adversary’s parameters ϕ are only updated after a *full training (and testing) run* of the Victim’s
 85 parameters θ . Note that ϕ is *static* during the whole training run (inner loop) of θ and only gets
 86 updated once the inner loop completes. In the outer loop, we optimise the Adversary’s objective \mathcal{J}
 87 with respect to ϕ using ES as a black-box optimisation technique.

88 **Train-Time Manipulation** When influencing the agent’s performance during train-time, we consider
 89 the allied and the adversarial settings. Pseudocode is provided in Algorithm 1 (see Appendix E,
 90 where E is the number of Victim training episodes and N is the ES population size. Letting $c = 1$
 91 for allied and $c = -1$ for adversarial, the Adversary’s objective is c times the Victim’s mean reward
 92 accumulated over training.

93 **Zero-Shot Test-Time Manipulation** In zero-shot test-time manipulation, the Adversary attempts to
 94 maximise its objective \mathcal{J} during some notion of test-time starting at time I . In practice, we introduce
 95 a separate Adversary for test-time, parameterized by its own set of parameters ψ . However, both the
 96 train-time Adversary ϕ and test-time Adversary ψ have identical objective function \mathcal{J} . The train-time
 97 Adversary wants to *create* a backdoor to make the Victim susceptible to manipulation at test-time.
 98 The test-time Adversary wants to *use* this backdoor to control the Victim. The test-time Adversary ψ
 99 operates zero-shot because it has not seen the specific, trained test-time parameters of the Victim θ'
 100 of the current meta-episode before interacting with it.

101 **Theoretical Results** For the theoretical results, we point the reader to Appendix B.1 and B.2.
 102 Moreover, in Appendix B.3, we informally show that removing any component from a Cheap
 103 Talk MDP would either nullify all possibility of influence or make the setting so limited as to be
 104 uninteresting.

105 4 Experiments and Results

106 We evaluate ACT on three different environments: Cartpole, Pendulum, and Reacher [4]. The Victim
107 is trained with Proximal Policy Optimisation [24, PPO]. The Adversary is trained using ES [23].

108 4.1 Train-Time Influence

109 Figure 1 show the results of training Victims alongside different Adversaries. It is evaluated on four
110 different Adversaries:

- 111 1. **Ally**: meta-trained to *maximize* the Victim’s mean reward throughout training.
- 112 2. **Adversary**: meta-trained to *minimize* the Victim’s mean reward throughout training.
- 113 3. **Random Adversary**: randomly initialise and fix the Adversary’s parameters ϕ .
- 114 4. **Zeroes Adversary**: appends only zeroes as messages.

115 The ally clearly improves performance while the Adversary significantly harms it. We hypothesise
116 that the Adversary may be inducing catastrophic interference within the environment, which was
117 observed by Fedus et al. [8] in Atari 2600 games. In Figure 6 Appendix D, we demonstrate that
118 the Adversary induces catastrophic interference in both the Adversarial setting by influencing the
119 correlation between gradient updates between different parts of a single inner loop episode. We
120 also study how the cheap-talk channel size affects the performance of the Adversary in Figure 5a
121 Appendix D.

122 4.2 Zero-Shot Test-Time Manipulation

123 In the setting of zero-shot test-time manipulation, the Adversary’s objective is to maximize the score
124 of a *goal-conditioned objective*. Consequently, the Adversary needs to learn to introduce a backdoor
125 during train-time and use the backdoor during test-time to fully control the Victim. We describe the
126 environment-specific rewards and how these goals are parameterized in Figure 3, Appendix D.

127 All results are shown in Figure 2. For an explanation of each ablation, see the Appendix C We
128 can use the Direct Oracle as a baseline to measure how effective the train-time Adversary ϕ and
129 test-time Adversary ψ are at achieving the maximal possible return jointly. As Figure 2 shows, the
130 ES optimized train- and test-time Adversaries perform near-optimally. We investigate this further in
131 Figure 4, Appendix D, where we compare the range and variance of Victims trained with ES optimized
132 Adversaries ϕ and Victims trained with random Adversaries ϕ_{random} across different message values.

133 5 Conclusion & Future Work

134 In this paper, we propose a novel, minimum-viable, adversarial setting for RL agents, where the
135 Adversary can only influence the Victim over messages, and can only do so with deterministic
136 function that only depends on the current state. By training a Adversary with adversarial cheap talk
137 (ACT), we show that appending to the observations of a learning agent, even with strong constraints,
138 is sufficient to drastically improve or *decrease* a learning agent’s train-time performance or introduce
139 a backdoor to control the learning agent at test time completely. Our test-time ablation studies
140 demonstrate that the train- and test-time Adversaries achieve near-optimal performance individually
141 as well as jointly, when compared against strong oracle baselines. We also provide in-depth analysis
142 on how the Adversaries work. As RL models become more widespread, we believe practitioners
143 should consider this new class of minimum viable attacks. Therefore, we propose identifying and
144 filtering out seemingly-superfluous information as the first defence measure. In future work we will
145 investigate different defence strategies, such as the identification of messages, and larger-scale input
146 settings.

References

- [1] C. Ashcraft and K. Karra. Poisoning deep reinforcement learning agents with in-distribution triggers. *arXiv preprint arXiv:2106.07798*, 2021.
- [2] E. Bengio, J. Pineau, and D. Precup. Interference and generalization in temporal difference learning. In *Proceedings of the 37th International Conference on Machine Learning*, volume 119 of *Proceedings of Machine Learning Research*, pages 767–777, 2020.
- [3] J. Bradbury, R. Frostig, P. Hawkins, M. J. Johnson, C. Leary, D. Maclaurin, G. Necula, A. Paszke, J. VanderPlas, S. Wanderman-Milne, and Q. Zhang. JAX: composable transformations of Python+NumPy programs, 2018.
- [4] G. Brockman, V. Cheung, L. Pettersson, J. Schneider, J. Schulman, J. Tang, and W. Zaremba. Openai gym, 2016.
- [5] K. Cao, A. Lazaridou, M. Lanctot, J. Z. Leibo, K. Tuyls, and S. Clark. Emergent communication through negotiation. In *6th International Conference on Learning Representations*, 2018.
- [6] V. P. Crawford and J. Sobel. Strategic information transmission. *Econometrica*, 50(6):1431–1451, 1982.
- [7] J. Farrell. Cheap talk, coordination, and entry. *The RAND Journal of Economics*, 18(1):34–39, 1987.
- [8] W. Fedus, D. Ghosh, J. D. Martin, M. G. Bellemare, Y. Bengio, and H. Larochelle. On catastrophic interference in atari 2600 games. *arXiv preprint arXiv:2002.12499*, 2020.
- [9] J. Foerster, R. Y. Chen, M. Al-Shedivat, S. Whiteson, P. Abbeel, and I. Mordatch. Learning with opponent-learning awareness. In *Proceedings of the 17th International Conference on Autonomous Agents and MultiAgent Systems*, pages 122–130, 2018.
- [10] J. N. Foerster, Y. M. Assael, N. de Freitas, and S. Whiteson. Learning to communicate with deep multi-agent reinforcement learning. In *Advances in Neural Information Processing Systems*, volume 29, pages 2137–2145, 2016.
- [11] A. Gleave, M. Dennis, C. Wild, N. Kant, S. Levine, and S. Russell. Adversarial policies: Attacking deep reinforcement learning. In *8th International Conference on Learning Representations*, 2020.
- [12] T. Gu, B. Dolan-Gavitt, and S. Garg. Badnets: Identifying vulnerabilities in the machine learning model supply chain. *arXiv preprint arXiv:1708.06733*, 2017.
- [13] S. H. Huang, N. Papernot, I. J. Goodfellow, Y. Duan, and P. Abbeel. Adversarial attacks on neural network policies. In *5th International Conference on Learning Representations, Workshop Track Proceedings*, 2017.
- [14] P. Kiourti, K. Wardega, S. Jha, and W. Li. Trojdr!l: Evaluation of backdoor attacks on deep reinforcement learning. In *57th ACM/IEEE Design Automation Conference*, pages 1–6, 2020.
- [15] J. Kos and D. Song. Delving into adversarial attacks on deep policies. In *5th International Conference on Learning Representations, Workshop Track Proceedings*, 2017.
- [16] R. T. Lange. evosax: Jax-based evolution strategies, 2022.
- [17] D. Lenton, F. Pardo, F. Falck, S. James, and R. Clark. Ivy: Templated deep learning for inter-framework portability. *arXiv preprint arXiv:2102.02886*, 2021.
- [18] A. Letcher, D. Balduzzi, S. Racanière, J. Martens, J. N. Foerster, K. Tuyls, and T. Graepel. Differentiable game mechanics. *J. Mach. Learn. Res.*, 20:84:1–84:40, 2019.

- 189 [19] A. Letcher, J. N. Foerster, D. Balduzzi, T. Rocktäschel, and S. Whiteson. Stable opponent
190 shaping in differentiable games. In *7th International Conference on Learning Representations*,
191 2019.
- 192 [20] C. Lu, T. Willi, C. Schroeder de Witt, and J. Foerster. Model-free opponent shaping. *arXiv*
193 *preprint arXiv:2205.01447*, 2022.
- 194 [21] C. Lyle, M. Rowland, and W. Dabney. Understanding and preventing capacity loss in reinforce-
195 ment learning. *arXiv preprint arXiv:2204.09560*, 2022.
- 196 [22] A. Salem, R. Wen, M. Backes, S. Ma, and Y. Zhang. Dynamic backdoor attacks against machine
197 learning models. *arXiv preprint arXiv:2003.03675*, 2020.
- 198 [23] T. Salimans, J. Ho, X. Chen, and I. Sutskever. Evolution strategies as a scalable alternative to
199 reinforcement learning. *arXiv preprint arXiv:1703.03864*, 2017.
- 200 [24] J. Schulman, F. Wolski, P. Dhariwal, A. Radford, and O. Klimov. Proximal policy optimization
201 algorithms. *arXiv preprint arXiv:1707.06347*, 2017.
- 202 [25] X. Song, Y. Jiang, S. Tu, Y. Du, and B. Neyshabur. Observational overfitting in reinforcement
203 learning. In *8th International Conference on Learning Representations*, 2020.
- 204 [26] H. van Hasselt, Y. Doron, F. Strub, M. Hessel, N. Sonnerat, and J. Modayil. Deep reinforcement
205 learning and the deadly triad. *arXiv preprint arXiv:1812.02648*, 2018.
- 206 [27] L. Wang, Z. Javed, X. Wu, W. Guo, X. Xing, and D. Song. BACKDOORL: backdoor attack
207 against competitive reinforcement learning. In *Proceedings of the Thirtieth International Joint*
208 *Conference on Artificial Intelligence*, pages 3699–3705, 2021.
- 209 [28] T. Willi, J. Treutlein, A. Letcher, and J. Foerster. COLA: consistent learning with opponent-
210 learning awareness. *arXiv preprint arXiv:2203.04098*, 2022.

211 **A Related Work**

212 **A.1 Test-Time Adversarial Attacks**

213 Most work investigating adversarial attacks on deep RL systems focuses on attacks at test-time, i.e.
214 those that assume a fully trained, static policy. Gleave et al. [11] learn adversarial policies to attack a
215 pretrained agent at test-time. In contrast to our method, the adversarial agent can directly interact
216 with the environment and the victim agent, thus introducing non-stationarity and assuming sampling
217 access to the (static) victim. Also, they do not investigate adversarial agents that affect training
218 performance. Huang et al. [13] also investigate adversarial attacks to influence test-time performance.
219 In contrast to our work, they directly perturb the observation space and do not simply append to it,
220 thus assuming access to the observation space of the victim. Kos and Song [15] also attack test-time
221 performance by directly perturbing observations.

222 **A.2 Backdoor Attacks**

223 Backdoor attacks in reinforcement learning aim to introduce a vulnerability during train-time, which
224 can be used at test-time. Backdoors can be static, meaning they get activated with fixed patterns,
225 or dynamic, which is when the backdoor gets activated by context-dependent patterns [22]. For
226 static backdoors, the adversary often directly perturbs the observation space [12, 14, 1]. To introduce
227 dynamic backdoors, the threat model assumes that the adversary has full control over the training
228 process of the agent, giving the adversary the ability to introduce backdoors at train-time [27]. In
229 contrast, in our threat model, we assume a minimal range of influence by only appending to the
230 observations. Furthermore, instead of perturbing the observations directly, Wang et al. [27] deploy
231 the adversarial agent directly in the environment. Directly interacting with the environment allows
232 the adversary to introduce non-stationarity and stochasticity. In contrast, our setting does not allow
233 the Adversary to introduce either.

234 **A.3 Failure Modes in Deep Reinforcement Learning**

235 Previous works have shown that using neural networks as function approximators in reinforcement
236 learning often results in multiple failure modes due to the non-stationarity of value function boot-
237 strapping [26]. In particular, works have shown that catastrophic interference [2] and capacity loss
238 [21] often occur, even within a single episode of an environment [8]. Song et al. [25] shows that deep
239 reinforcement learning algorithms can often overfit to spurious correlations in the observation space.
240 By appending to the observation space, we learn to induce the observational failure modes described
241 in these works.

242 **A.4 Opponent Shaping / Cheap Talk**

243 Our method is closely related to the field of opponent shaping. Originally, most opponent shaping
244 algorithms assumed white-box access to their opponents to shape the flow of the opponent’s gradient
245 [9, 18, 19, 28]. Instead, Lu et al. [20] introduce a method to shape opponents without white-box
246 access. However, they still deploy an agent to interact directly in the environment. In contrast, we
247 propose a method to shape other agents without having to interact in the environment at all, solely
248 by appending messages through a cheap talk channel. Cheap talk is communication that incurs no
249 cost, is non-binding (it can be ignored and does not limit the agent’s action space), and is unverifiable
250 (meaning any information, true or false, can be communicated) [7]. In RL terms, a cheap talk
251 channel is a part of the state space which can be observed by other agents but does not alter transition
252 dynamics or reward functions. Cheap talk channels [6] in deep reinforcement learning have been
253 used to learn emergent communication [10] and to solve coordination problems [5]. To the best of
254 our knowledge, this paper is the first to use a cheap talk channel (and only a cheap talk channel) to
255 shape learning agents.

256 **B Proofs**

257 **B.1 Proof of Proposition 1**

258 In this section, we further justify the claim that our setting represents the bare minimum range of
 259 influence. To begin, we prove that Adversaries *cannot* influence *tabular* Victims in Cheap Talk
 260 MDPs; Adversaries can therefore only interfere with a Victim through their function approximator.

261 **Proposition 1.** *For any deterministic Adversary $f : \mathcal{S} \rightarrow \mathcal{M}$, the return of a tabular Victim initialised*
 262 *uniformly along the \mathcal{M} axis is independent from f . Moreover, any Victim which is guaranteed to*
 263 *converge to optimal policies in MDPs will, for any Cheap Talk MDP, converge to a policy whose*
 264 *expected return is the optimal return for the original no-channel MDP – even in non-tabular settings*
 265 *and regardless of initialisation.*

266 *Proof.* We begin with the tabular case.

267 **Tabular Victims.** In a Cheap Talk MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mathcal{M}, f, \mathcal{J} \rangle$, a tabular Victim arbitrarily orders
 268 states as $\{s_1, \dots, s_d\}$ and messages as $\{m_1, \dots, m_k\}$, where $d = |\mathcal{S}|$ and $k = |\mathcal{M}|$, and stores
 269 policies $\pi_t(\cdot | s_i, m_j)$ at time t of the learning process for all $i \in [d], j \in [k]$. The argument follows
 270 identically for value functions. Assuming uniform initialisation along the \mathcal{M} axis means that

$$\pi_0(\cdot | s_i, m_j) = \pi_0(\cdot | s_i, m_{j'})$$

271 for all $j, j' \in [k]$. Now consider any two Adversaries f, g and their influence on two copies of the
 272 same Victim V, W with respective policies π, χ . The only states encountered in the environment are
 273 of the form $(s, f(s))$ and $(s, g(s))$ respectively, so Victims only update the corresponding policies

$$\pi_t(\cdot | s_i, f(s_i)) \quad \text{and} \quad \chi_t(\cdot | s_i, g(s_i)).$$

274 We prove by induction that these quantities are equal for all t . The base case holds by uniform
 275 initialisation along \mathcal{M} ; assume the claim holds for all fixed $0 \leq t \leq T$. The Victims update their
 276 policies at time $T + 1$ according to the same learning rule, as a function of the transitions and returns
 277 under current and past policies π_t and χ_t respectively. Transitions take the form $(s, f(s), a, s', f(s'))$
 278 for V and $(s, g(s), a, s', g(s))$ for W , which have identical probabilities and returns because

$$\begin{aligned} \pi_t(a | s_i, f(s_i)) &= \chi_t(a | s_i, g(s_i)); \\ \mathcal{P}(s', f(s') | s, f(s), a) &= \mathcal{P}(s', g(s') | s, g(s), a); \\ \mathcal{R}(s, f(s), a) &= \mathcal{R}(s, g(s), a) \end{aligned}$$

279 by inductive assumption and independence of \mathcal{P}, \mathcal{R} from \mathcal{M} . This probability- and return-preserving
 280 bijection between transitions, as well as being copies with identical initialisation in the environment,
 281 implies that policies $\pi_T(\cdot | s_i, f(s_i)) = \chi_T(\cdot | s_i, g(s_i))$ are updated identically to

$$\pi_{T+1}(\cdot | s_i, f(s_i)) = \chi_{T+1}(\cdot | s_i, g(s_i))$$

282 as required to complete induction. Note that this could not necessarily be accomplished in non-tabular
 283 settings, where updating parameters θ of the function approximator for some state s_i may alter
 284 the policy on some other state s_j . It now follows that trajectories $\tau = (s^k, f(s^k), a^k)_k$ for V and
 285 $\omega = (s^k, g(s^k), a^k)_k$ for W have identical probabilities and hence produce identical returns

$$\mathbb{E}_{\tau \sim \pi_t} [\mathcal{R}(\tau)] = \mathbb{E}_{\omega \sim \chi_t} [\mathcal{R}(\omega)]$$

286 at any timestep t of the learning process, concluding independence from Adversaries.

287 **Optimally Convergent Victims.** By assumption, the Victim is guaranteed to converge to an optimal
 288 policy $\bar{\pi}$ in the Cheap Talk MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \mathcal{M}, f, \mathcal{J}, \gamma \rangle$, since a Cheap Talk MDP is itself an
 289 MDP with an augmented state space $\mathcal{S} \times \mathcal{M}$ and augmented transition/reward functions that are
 290 defined to be independent from \mathcal{M} . Now $\bar{\pi}$ naturally induces a policy π on the no-channel MDP,

291 given by $\pi(\cdot | s) := \bar{\pi}(\cdot | s, f(s))$, and in particular $Q(s, a) = \bar{Q}(s, f(s), a)$ by independence of
 292 transitions and rewards from \mathcal{M} . Optimality of π follows directly from the Bellman equation

$$\begin{aligned} Q(s, a) &= \bar{Q}(s, f(s), a) = \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a), r \sim \mathcal{R}(\cdot | s, a)} \left[r + \gamma \max_{a' \in \mathcal{A}} \bar{Q}(s', f(s'), a') \right] \\ &= \mathbb{E}_{s' \sim \mathcal{P}(\cdot | s, a), r \sim \mathcal{R}(\cdot | s, a)} \left[r + \gamma \max_{a' \in \mathcal{A}} Q(s', a') \right]. \end{aligned}$$

293 Now trajectories $\bar{\tau} = (s^k, f(s^k), a^k)_k$ and $\tau = (s^k, a^k)_k$ have identical probability and return under
 294 π and $\bar{\pi}$ respectively, so the Victim has expected return

$$\mathbb{E}_{\bar{\tau} \sim \bar{\pi}} [\mathcal{R}(\bar{\tau})] = \mathbb{E}_{\tau \sim \pi} [\mathcal{R}(\tau)]$$

295 which is the optimal expected return of the original no-channel MDP. \square

296 B.2 Proof of Proposition 2

297 For completeness, we also formally prove our claims from the introduction regarding what the
 298 Adversary *cannot* do in Cheap Talk MDPs.

299 **Proposition 2.** *In a Cheap Talk MDP, the Adversary cannot (1) occlude the ground truth, (2) influence*
 300 *the environment dynamics / reward functions, (3) see the Victim's actions or parameters, (4) inject*
 301 *stochasticity, or (5) introduce non-stationarity.*

302 *Proof.* Mostly by definition. Formally, consider a Cheap Talk MDP $\langle \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \mathcal{M}, f, \mathcal{J} \rangle$ as
 303 defined in Section 2. For a fixed training / testing run of the Victim on the MDP, the Adversary
 304 outputs a message $f(s)$ at each step according to a fixed deterministic function $f : \mathcal{S} \rightarrow \mathcal{M}$.

- 305 (1) The message is *appended* to the state s and the Victim acts with full visibility of the ground
 306 truth (state) s according to its policy: $a \sim \pi(\cdot | s, f(s))$.
- 307 (2) The transition and reward functions \mathcal{P}, \mathcal{R} are defined to be independent from \mathcal{M} . Formally we
 308 have $\mathcal{P}(\cdot | s, m, a) = \mathcal{P}(\cdot | s, m', a)$ for all $m, m' \in \mathcal{M}$ (similarly for \mathcal{R}), so the Adversary's
 309 choice of message $m = f(s)$ cannot influence \mathcal{P} or \mathcal{R} .
- 310 (3) $f : \mathcal{S} \rightarrow \mathcal{M}$ is defined as a function of \mathcal{S} only, so the Adversary cannot condition its policy
 311 based on the Victim's actions or parameters (i.e. it cannot see them for all practical purposes).
- 312 (4) f is a deterministic function, so $\pi(\cdot | s, f(s))$ is a distribution only on actions \mathcal{A} . The transition
 313 and reward functions are independent from f , so they are distributions only on state-action pairs
 314 $\mathcal{S} \times \mathcal{A}$. It follows that the Adversary injects no further stochasticity into the MDP.
- 315 (5) f is static for a fixed training / testing run, so $s_t = s_{t'}$ implies $f(s_t) = f(s_{t'})$ for all timesteps
 316 t, t' in the run. It follows that any given Victim policy π is stationary, namely $\pi(\cdot | s_t, f(s_t)) =$
 317 $\pi(\cdot | s_{t'}, f(s_{t'}))$ for all $s_t = s_{t'}$. Since \mathcal{P} and \mathcal{R} are stationary (as defined by a standard MDP)
 318 and independent from \mathcal{M} , their stationarity is also preserved. \square

319 B.3 Informal Justification of Minimality

320 Finally, let us informally show that removing any component from a Cheap Talk MDP would either
 321 nullify all possibility of influence or make the setting so limited as to be uninteresting.

- 322 (1) Removing the set \mathcal{M} or the policy $f : \mathcal{S} \rightarrow \mathcal{M}$ entirely would result in the Victim being
 323 completely independent from the Adversary, since nothing would be appended to its observation.
- 324 (2) Restricting the capacity of \mathcal{M} to a certain number of bits would further restrict an Adversary's
 325 range of influence, so one could say that the *truly* minimum-viable setting is to impose a set of
 326 size $|\mathcal{M}| = 1$. However, cheap talk is still cheap talk when varying capacity, and there is no
 327 reason to arbitrarily restrict the size to 1 if we are to apply our setting to complex environments
 328 likely requiring more than a single bit of communication to witness interesting results.

329 (3) Not allowing Adversaries to see states, namely removing \mathcal{S} as inputs to f , yields a function
330 $f : \{0\} \rightarrow \mathcal{M}$ which always outputs the same message $f(0) = m \in \mathcal{M}$. This is equivalent to the
331 previous restriction of imposing a set \mathcal{M} of size 1, since in this case any function $f : \mathcal{S} \rightarrow \mathcal{M}$
332 would have to output the unique element $f(s) = m$ for all input states s .

333 (4) The Adversary must have some objective function \mathcal{J} in order for an adversarial setting to make
334 sense – removing it would remove the Adversary’s rationale for existence, since it would have
335 no incentive to learn parameters that influence the Victim according to some goal.

336 (5) Restricting the function class of objectives \mathcal{J} is a valid minimisation of the setting, but simply
337 restricts our interesting the setting itself. The setting should at the very least allow for adversarial
338 objectives of the form $\mathcal{J} = -J$ as we consider in the train-time setting. In test-time, our aim is
339 to show how Adversaries can exert arbitrary control over Victims despite cheap talk restrictions,
340 and we therefore consider more general objective functions.

341 **C Zero-Shot Ablations**

342 1. **Direct Oracle:** In this baseline, there is no cheap talk. We simply train a PPO agent to
343 maximize the *goal-conditioned return*. It can observe the full state and directly output
344 actions in the environment.

345 2. **Zero-Shot Adversary:** First, we train a Victim θ alongside a train-time Adversary ϕ . We
346 then evaluate the return of the test-time Adversary ψ according to the goal-conditioned
347 return (as described in Algorithm 2). The test-time Adversary ψ operates zero-shot because
348 it was not trained with the specific, trained instance of the Victim θ before interacting with
349 it. It is thus represented by a horizontal line in Figure 2.

350 3. **Oracle with Learned Adversary:** First, we optimize the Victim θ by training it alongside
351 our train-time Adversary ϕ . Then, instead of ES, we use PPO to train the test-time Adversary
352 ψ^* against the Victim θ . Unlike the zero-shot Adversary, the oracle ψ^* is allowed to train
353 against the pretrained and fixed Victim θ to maximize its returns, as described in Algorithm
354 4 in Appendix E.

355 4. **Oracle with Random Adversary:** First, we obtain a Victim θ by training it alongside a
356 *random* train-time Adversary, ϕ_{random} , with randomly initialized and fixed parameters. Next,
357 we use PPO to train the test-time Adversary ψ^* to maximize the goal-conditioned return.

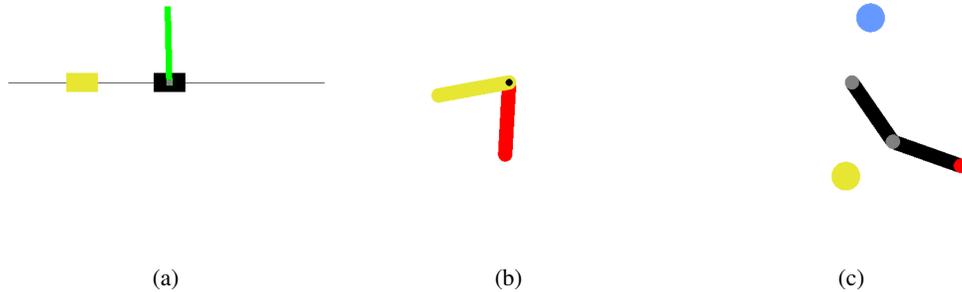


Figure 3: Visualizations of our goal-conditioned environments (a) In Cartpole, the Adversary’s target is a randomly selected point on the x-axis, indicated by the yellow box. (b) In Pendulum, the Adversary’s goal is a randomly selected angle indicated by the yellow pole. (c) In Goal-Conditioned Reacher, the Adversary’s goal is a specific point, denoted by the yellow circle, while the Victim’s goal is the blue circle.

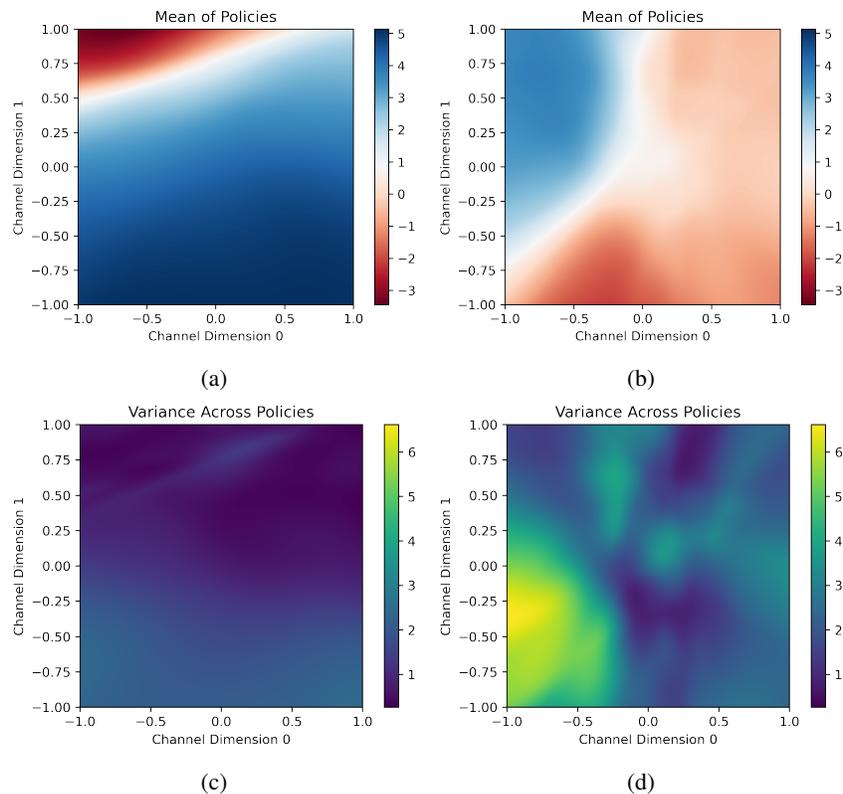


Figure 4: We train 10 different Victims alongside the Learned ϕ (a & c), as well as 10 different Victims alongside a randomly generated ϕ (b & d) in the Pendulum environment. (a) and (b) show the mean of the policy output across the 10 Victims as we vary the value of the message in a fixed randomly selected state. Notably, the policies trained with the learned ϕ achieve a much wider range of outputs. (c) and (d) show the variance of the policy output across the 10 Victims. Notably, the policies trained with the learned ϕ display very little variance, implying that the learned ϕ shapes the Victim in a consistent way.

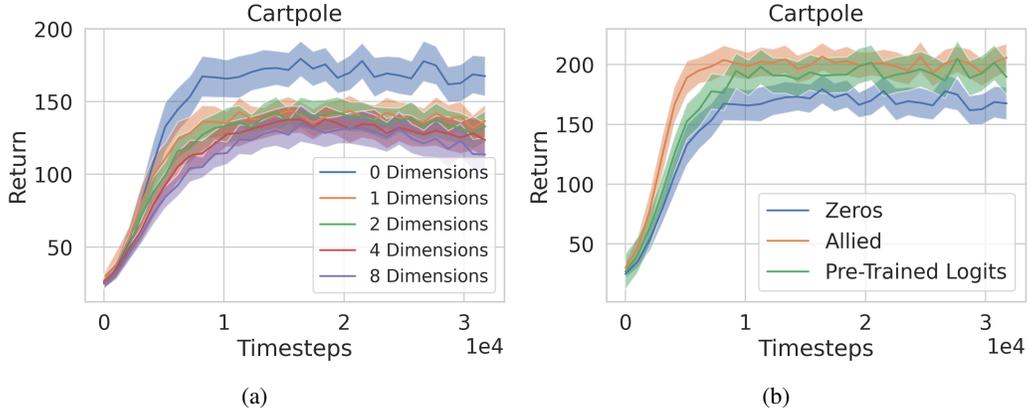


Figure 5: (a) Ablations on the different number of cheap talk dimensions for the Adversary in Cartpole (b) Comparing the ally with an Adversary that outputs the optimal logits in Cartpole. Error bars denote the standard error across 10 seeds of a Victim trained against a single meta-trained Adversary.

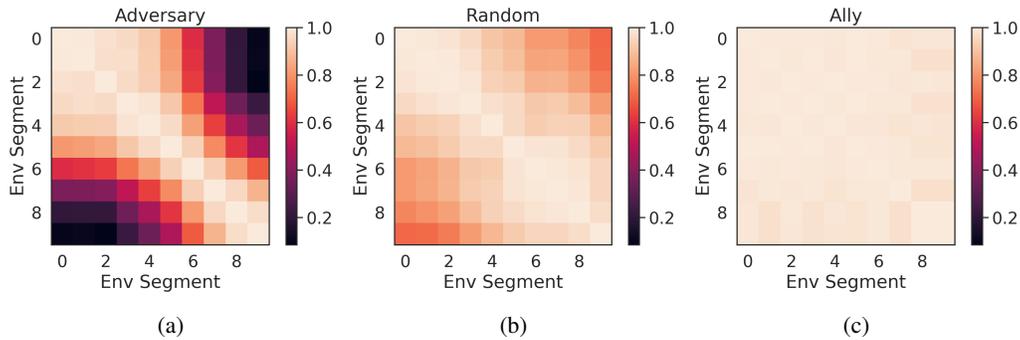


Figure 6: To perform this analysis, we collect each Victim’s experience buffer, before the agents have converged in training, and split each one into 10 bins, ordered by the time-step within the environment. We then calculate the gradient update the agents would perform on each of these bins. In the Adversarial setting (a), the gradient updates performed for transitions sampled early in an episode can *interfere* with the gradient updates performed for transitions later in an episode. Meanwhile, in the Allied setting (c), those gradient updates are positively correlated, suggesting that the gradient updates aid each other.

Algorithm 1 Train-time ACT

```
1: Set  $c = \pm 1$  for allied / adversarial
2: Initialize Adversary parameters  $\phi$ 
3: for  $m = 0$  to  $M$  do
4:   Sample  $\phi_n \sim \phi + \sigma\epsilon_n$  where  $\epsilon_1, \dots, \epsilon_N \sim \mathcal{N}(0, I)$ 
5:   for  $n = 0$  to  $N$  do
6:     Initialize Victim parameters  $\theta$ 
7:     rewards = []
8:     for  $e = 0$  to  $E$  do
9:       s = env.reset()
10:      while not done do
11:         $\bar{s} = [s, f_{\phi_n}(s)]$ 
12:         $a \sim \pi_{\theta}(\cdot | \bar{s})$ 
13:         $r, s, \text{done} = \text{env.step}(a)$ 
14:        rewards.append( $r$ )
15:      end while
16:      Update  $\theta$  with PPO to maximise  $J$ 
17:    end for
18:     $\mathcal{J}_n = c \cdot \text{sum}(\text{rewards})/\text{len}(\text{rewards})$ 
19:  end for
20:  Update  $\phi$  using ES to maximise  $\mathcal{J}$ 
21: end for
```

Algorithm 2 Test-time ACT

```
1: Initialize train-time ACT parameters  $\phi$ 
2: Initialize test-time ACT parameters  $\psi$ 
3: for  $m = 0$  to  $M$  do
4:   Sample  $\phi_n \sim \phi + \sigma\epsilon_n$  where  $\epsilon_1, \dots, \epsilon_N \sim \mathcal{N}(0, I)$ 
5:   Sample  $\psi_n \sim \psi + \sigma\epsilon_n$  where  $\epsilon_1, \dots, \epsilon_N \sim \mathcal{N}(0, I)$ 
6:   for  $n = 0$  to  $N$  do
7:     Initialize policy params  $\theta$ 
8:     rewards = []
9:     for  $e = 0$  to  $E$  do
10:      s = env.reset()
11:      while not done do
12:         $m = f_{\phi_n}(s)$ 
13:         $\bar{s} = [s, m]$ 
14:         $a \sim \pi_{\theta}(\cdot | \bar{s})$ 
15:         $r, s = \text{env.step}(a)$ 
16:      end while
17:      Update  $\theta$  using PPO to maximise  $J$ 
18:    end for
19:    for  $i = 0$  to  $I$  do
20:      s = env.reset()
21:      while not done do
22:         $m = f_{\psi_n}(s)$ 
23:         $\bar{s} = [s, m]$ 
24:         $a \sim \pi_{\theta}(\cdot | \bar{s})$ 
25:         $r, s, \text{done} = \text{env.step}(a)$ 
26:         $r_t^S = R^S(s, a)$ 
27:        rewards.append( $r_t^S$ )
28:      end while
29:    end for
30:  end for
31:  Update  $\phi$  using ES to maximise  $\mathcal{J}$ 
32:  Update  $\psi$  using ES to maximise  $\mathcal{J}$ 
33: end for
```

Algorithm 3 Test-time Oracle PPO ACT

```
1: Initialize train-time ACT parameters  $\phi$ 
2: Obtain trained  $\phi, \theta$  from Algorithm 2
3: Initialize test-time ACT parameters  $\psi^*$ 
4: for  $i = 0$  to  $I$  do
5:    $s = \text{env.reset}()$ 
6:   while not done do
7:      $m \sim \pi_{\psi^*}(\cdot | s)$ 
8:      $\bar{s} = [s, m]$ 
9:      $a \sim \pi_{\theta}(\cdot | \bar{s})$ 
10:     $r, s, \text{done} = \text{env.step}(a)$ 
11:     $r_t^S = R^S(s, a)$ 
12:    rewards.append( $r_t^S$ )
13:  end while
14:  Update  $\psi^*$  using PPO to maximise  $\mathcal{J}$ 
15: end for
```

Algorithm 4 Test-time Random Shaper

```
1: Initialize train-time ACT parameters  $\phi_{\text{random}}$ 
2: Initialize policy params  $\theta$ 
3: rewards = []
4: for  $e = 0$  to  $E$  do
5:    $s = \text{env.reset}()$ 
6:   while not done do
7:      $m = f_{\phi_{\text{random}}}(s)$ 
8:      $\bar{s} = [s, m]$ 
9:      $a \sim \pi_{\theta}(\cdot | \bar{s})$ 
10:     $r, s = \text{env.step}(a)$ 
11:   end while
12:   Update  $\theta$  using PPO to maximise  $J$ 
13: end for
14: Initialize test-time ACT parameters  $\psi^*$ 
15: for  $i = 0$  to  $I$  do
16:    $s = \text{env.reset}()$ 
17:   while not done do
18:      $m \sim \pi_{\psi^*}(\cdot | s)$ 
19:      $\bar{s} = [s, m]$ 
20:      $a \sim \pi_{\theta}(\cdot | \bar{s})$ 
21:      $r, s = \text{env.step}(a)$ 
22:      $r_t^S = R^S(s, a)$ 
23:     rewards.append( $r_t^S$ )
24:   end while
25:   Update  $\psi^*$  using PPO to maximise  $\mathcal{J}$ 
26: end for
```

360 **F Hyperparameter Details**

361 We train thousands of agents per minute on a single V100 GPU by vectorising both the *PPO*
 362 *algorithm itself* and the environments using Jax [3]. This allows us to JIT-compile the *full training*
 363 *pipeline* and perform end-to-end deep RL training completely on GPUs. We adapt the environment
 364 implementations from Brockman et al. [4] and Lenton et al. [17] and use the ES implementation
 365 from Lange [16]. This compute setup allows us to efficiently perform outer-loop ES on the full
 366 training trajectories of inner-loop PPO agents. For example, in Cartpole, we run 8192 PPO Victims
 367 alongside 8192 train-time Adversaries and 8192 test-time Adversaries, each over four instances of the
 368 environment on a single V100 GPU. Over 1024 generations of ES, this results in training 8,388,608
 369 PPO agents from scratch in 2 hours on 4 V100 GPUs.

370 We report the hyperparameter values used for each environment in our experiments.

Table 1: Important parameters for the Cartpole environment

Parameter	Value
State Size	4
message Size	2
Number of Environments	4
Maximum Grad Norm	0.5
Number of Updates	32
Update Period	256
Outer Discount Factor γ	0.99
Number of Epochs per Update	16
PPO Clipping ϵ	0.2
General Advantage Estimation λ	0.95
Critic Coefficient	0.5
Entropy Coefficient	0.01
Learning Rate	0.005
Population Size	1024
Number of Generations	2049
Outer Agent (OA) Hidden Layers	2
OA Size of Hidden Layers	64
OA Hidden Activation Function	ReLU
OA Output Activation Function	Tanh
Inner Agent (IA) Actor Hidden Layers	2
IA Size of Actor Hidden Layers	32
IA Number of Critic Hidden Layers	2
IA Size of Critic Hidden Layers	32
IA Activation Function	Tanh
Number of Rollouts	4

Table 2: Important parameters for the Pendulum environment

Parameter	Value
State Size	3
message Size	2
Number of Environments	16
Maximum Grad Norm	0.5
Number of Updates	128
Update Period	256
Outer Discount Factor γ	0.95
Number of Epochs per Update	16
PPO Clipping ϵ	0.2
General Advantage Estimation λ	0.95
Critic Coefficient	0.5
Entropy Coefficient	0.005
Learning Rate	0.02
Population Size	768
Number of Generations	2049
Outer Agent (OA) Hidden Layers	2
OA Size of Hidden Layers	64
OA Hidden Activation Function	ReLU
OA Output Activation Function	Tanh
Inner Agent (IA) Actor Hidden Layers	1
IA Size of Actor Hidden Layers	32
IA Number of Critic Hidden Layers	1
IA Size of Critic Hidden Layers	32
IA Activation Function	Tanh
Number of Rollouts	4

Table 3: Important parameters for the Reacher environment

Parameter	Value
State Size	10
message Size	4
Number of Environments	32
Maximum Grad Norm	0.5
Number of Updates	256
Update Period	128
Outer Discount Factor γ	0.99
Number of Epochs per Update	10
PPO Clipping ϵ	0.2
General Advantage Estimation λ	0.95
Critic Coefficient	0.5
Entropy Coefficient	0.0005
Learning Rate	0.004
Population Size	128
Number of Generations	2049
Outer Agent (OA) Hidden Layers	2
OA Size of Hidden Layers	64
OA Hidden Activation Function	ReLU
OA Output Activation Function	Tanh
Inner Agent (IA) Actor Hidden Layers	2
IA Size of Actor Hidden Layers	128
IA Number of Critic Hidden Layers	2
IA Size of Critic Hidden Layers	128
IA Activation Function	ReLU
Number of Rollouts	4