

Supplementary Material

Paper ID: 3365

A Victim Population Settings

Collective learning frameworks where agents train on copies of the same environment fall under three lines of research [Parisotto *et al.*, 2019; Marthi *et al.*, 2005; Dimakopoulou *et al.*, 2018; Lupu and Precup, 2020; Yang *et al.*, 2020]. The first line aims to decrease the complexity of solving large Markov decision processes that model sizeable multi-agent systems [Marthi *et al.*, 2005]. They do so by breaking down the problem into tasks that are executed in parallel. Each task houses one or more agents and coordinates at run-time in order to push the agent population towards optimal behavior. The second line of research enables a single agent to better explore the given environment in a more efficient manner, by interacting with it in parallel using different policies [Dimakopoulou *et al.*, 2018]. In the third line of research, the agent aims to learn how to learn and strives to be able to perform efficiently across a family of tasks (Meta RL) [Parisotto *et al.*, 2019]. Herein, the agent trains on different tasks, in parallel. In all three domains, the agents exchange information in order to learn better strategies more quickly. [Yaman *et al.*, 2022] points out that an optimal balance between individual and social learning is critical for learning optimal behaviors in a population. Herein, individual learning is when each agent explores, interacts, and learns from the environment separately while social learning is when agents copy each other's behaviors. Therefore, individual learning leads to innovations which can then efficiently spread through the population via social learning. However, individual learning incurs high exploration cost. On the other hand, social learning is cost-effective but requires construction of a strategy by which the population chooses which behaviors (successful vs majority) to copy. The optimality of these strategies, akin to meta-control strategies in the human brain [Daw *et al.*, 2005], are sensitive to environment variability. In environments that change frequently, social learning might hamper the population's learning as the exchanged information can become invalid very quickly. This work proposes three learning settings namely; Implicit Collective, Collective, and Swarm Collective; wherein individual learning decreases while social learning increases progressively. In these settings, each agent trains to learn its individual task with the support of a separate reward signal while the attacker trains to push the complete (victim) populations towards a target behavior.

In the Implicit Collective setting, the innovation capabil-

ity of the population is maximized via individual learning (at the cost of exploration), as agents practice Decentralized Training, Decentralized Execution (DTDE). This is achieved as a collection of learning agents individually experience a commonly parameterized environment, and interact only implicitly via observations of environment modifications. The environment observations serve as a medium of implicit interaction as the attacker takes a single attack action to modify/poison the blueprint of the victim environment, conditioned on the behaviors of all agents present in the population. Agents, therefore, have some influence on each other, since a failure or stubbornness of one of them has an effect on the next attack presented to the entire population. Drawing inspiration from [Dimakopoulou *et al.*, 2018], exploration is made more efficient during testing by assigning an independent, random number generation seed to each agent in the victim population. This increases diversity as well as commitment of the victim agents. The Implicit Collective setting can find application in the development of an adaptive single-player computer game engine wherein the adaptive game engine represents an attacker that strives to simultaneously push a complete population of victims (players) towards a target behavior (with maximum player stickiness) and each attack action represents a new release/version of the game.

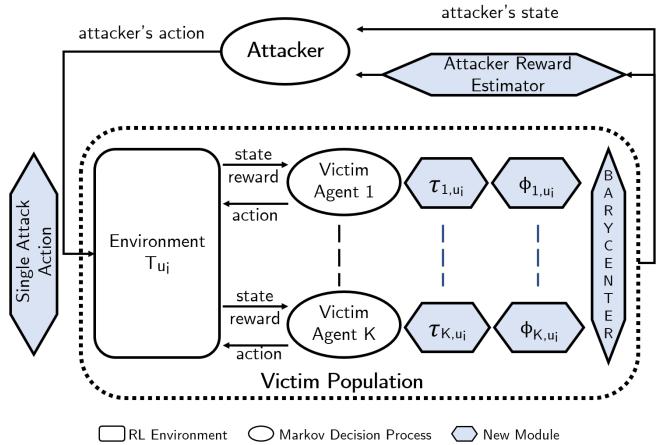


Figure 1: Bi-Level Attack Framework in Collective and Swarm Collective Settings

69 In Collective and Swarm Collective settings the agents occupy
70 the same copy of the victim environment. Herein, each agent observes (anonymized version of) every other agent,
71 and takes actions conditioned on this observation (along
72 with its own position inside the victim environment). The
73 agents practice Centralized Training, Decentralized Execution
74 (CTDE) where a single network is trained using all victims' interactions. As the agents learn from each other's ex-
75 periences, these settings support social learning. In Collective
76 setting, individual learning is enhanced by inculcation of
77 "agent indication" i.e. addition of an indication of the ob-
78 serving agent to its observations [Terry *et al.*, 2020]. This
79 enables the same neural network to represent diverse victim
80 behaviors. In addition, each agent is committed to a sepa-
81 rate, diverse behavior by assigning it an independent, random
82 number generation seed [Dimakopoulou *et al.*, 2018]. Collective
83 setting therefore presents a balanced scenario that houses
84 high support for both individual and social learning. On the
85 other hand, Swarm Collective does not inculcate agent indica-
86 tion and independent random seeds, and therefore strongly
87 promotes social learning. The centrally trained neural net-
88 work in this setting learns a single optimal victim behavior
89 using all victims' interactions. Lastly, a soft competition
90 is injected into both settings by terminating victim-training
91 episodes as soon as at least one victim reaches the goal state
92 (or maximum training time has elapsed); to enable faster
93 adoption of optimal behaviors inside the victim population.
94

96 **B (Expanded) Related Work**

97 This appendix helps to further position the current paper
98 in terms of the proposed methodology as well as the over-
99 all framework through comparison with literature in related
100 domains of Group-Behavior Modelling, Set Representation
101 Learning, and Unsupervised Environment Design respec-
102 tively.

103 **B.1 Group-Behavior Modelling**

104 The existing research works that strive to extract or model
105 the behavior of a group of agents typically follow two ap-
106 proaches. In the first approach, a given agent strives to in-
107 culate the effect of the behavior of other agents in its own
108 learning, by directly feeding the other agents' observations,
109 rewards, and/or behavior trajectories ((state, action, next ac-
110 tion) tuples) into its own quality, value and/or policy net-
111 work(s) [Zhang and Lesser, 2010; Foerster *et al.*, 2018]. In
112 the second approach, a given agent first learns the single-
113 agent behavior models of all the "other" agents [Raileanu
114 *et al.*, 2018; Papoudakis and Albrecht, 2020; Rabinowitz
115 *et al.*, 2018; Grover *et al.*, 2018; Tacchetti *et al.*, 2019;
116 Grover *et al.*, 2018; Shum *et al.*, 2019; He *et al.*, 2016;
117 Papoudakis *et al.*, 2021] and then utilizes the predictions of
118 these models in its own decision making. To the best of our
119 knowledge, this is the first work that creates a third category
120 of group-behavior modeling wherein the complete group's
121 behavior is modeled together to create a single latent repre-
122 sentation that captures the distribution of behaviors present
123 inside the agent population.

124 **B.2 Set Representation Learning**

125 Set-based machine learning tasks have recently come into fo-
126 cus, prompting several research works to propose deep mod-
127eling solutions to permutation-invariant collections of data.
128 PointNet [Qi *et al.*, 2017] and Deep Sets [Zaheer *et al.*, 2017]
129 learn a representation of an input set by first feeding each
130 element of this set individually into a neural network and
131 then aggregating the network outputs by using max-pooling
132 (PointNet) and summation (Deep Sets) respectively. Rep the
133 Set [Skianis *et al.*, 2020] computes the similarity of a given
134 input set, to a pre-defined number (say k) of learnable refer-
135 ence sets and uses these similarity values (vector of length k)
136 as the representation of the input set. Set Transformer [Lee
137 *et al.*, 2019] is a permutation-invariant attention-based neu-
138 ral network developed by removing the positional encoding
139 present in the original transformer and inculcating latent vari-
140 ables to reduce computational complexity. All these models
141 are completely flexible wrt the size of the input set except
142 Set Transformer which works with all set sizes smaller or
143 equal to a pre-decided maximum. However, unlike the ap-
144 proach proposed in this paper, the aforementioned models do
145 not support unsupervised learning. Optimal Transport Ker-
146 nel Embedding (OTKE) [Mialon *et al.*, 2020] is the first set-
147 based representation learning model that supports unsuper-
148 vised learning. OTKE constructs a reference set, of a pre-
149 defined length (say k), consisting of "k" cluster centroids or
150 "k" Wasserstein barycenters of the input data. Thereafter,
151 OTKE uses the Gaussian kernel to map all elements of the
152 input and reference set to a non-linear space. Finally, the in-
153 put set's representation constitutes a set of k-weighted sums
154 of the input set elements wherein the j^{th} sum's weights are
155 the similarity scores between all input elements and reference
156 element j. The similarity score between the two elements is
157 given by the transport map between them. OTKE is therefore
158 a soft clustering technique based on the Wasserstein metric.
159 The approach proposed in this work can be seen as a sim-
160 plification of OTKE as it treats the individual behaviors of
161 victims as a single behavior cluster. This simplified approach
162 is justified as the victim populations considered in this work
163 consist of homogeneous agents. Future works on attacks on
164 heterogeneous victim populations can employ OTKE for con-
165 structing population behavior representations.

166 **B.3 Unsupervised Environment Design**

167 Design of RL environments takes a lot of time and effort, is
168 error-prone, and is infused with designer bias. Unsupervised
169 Environment Design (UED) is a recent paradigm that aims
170 to not only automate this step but generate environment dis-
171 tributions that are conducive to emergent complexity, robust-
172 ness, and efficient transfer learning in RL agents. Domain
173 Randomization [Tobin *et al.*, 2017] creates a distribution of
174 environments by assigning random values to the free param-
175 eters of the environment. Even though it exposes agents to
176 a wide variety of environments, it does not create environ-
177 ments with complex structures and is not reactive to the capa-
178 bilities of the learning agent. Minimax Adversarial Training
179 [Morimoto and Doya, 2005] makes environment distribution
180 generation reactive to the learning agent's capabilities by in-
181 troducing an adversary that sequentially creates environments

182 that minimize the learning agent’s rewards. This worst-case
 183 tendency has however shown to create unsolvable environments
 184 that hinder the learning agent’s progress. PAIRED
 185 [Dennis *et al.*, 2020] does away with this worst-case ten-
 186 dency by introducing an additional ”antagonist” agent which
 187 is smarter than the learning agent (protagonist). The adver-
 188 sary’s objective is to maximize regret i.e. the difference be-
 189 tween these two agents’ rewards. This enables the generation
 190 of challenging but solvable environments. However, training
 191 the environment-creating adversary is challenging and is
 192 shown to produce weaker results than when main agents are
 193 trained on randomly-selected high-regret environments [Jiang
 194 *et al.*, 2021a]. This difficulty in training the adversary is due
 195 to sparse rewards and long-horizon credit assignment prob-
 196 lems. Inspired by the Teacher-Student Curriculum Learning
 197 paradigm [Matiisen *et al.*, 2019], Prioritized Level Re-
 198 play (PLR) [Jiang *et al.*, 2021b], and PLR+ [Jiang *et al.*,
 199 2021a] do away with adversary-based mechanism and instead
 200 strategically sample the next environment by prioritizing en-
 201 vironments that have larger future learning potential. How-
 202 ever, they cannot create new challenging environments with
 203 complex structures. This limits the frequency with which
 204 the learning agent can be exposed to complex structures and
 205 can thus hinder its progress, especially in high-dimensional
 206 complex environments. ACCEL [Parker-Holder *et al.*, 2022]
 207 does away with both, adversary-training-based and sampling-
 208 based methodologies’ drawbacks by utilizing an evolutionary
 209 environment generator and a regret-based curator. The gen-
 210 erator makes small modifications (mutations) to high-regret
 211 environments present inside the replay buffer. These modi-
 212 fied environments are added to the replay buffer if they result
 213 in high regret. Regret thus functions as the fitness function
 214 for evolution and helps develop challenging environments for
 215 the learning agent.

216 Training-time environment-poisoning attacks can be com-
 217 pared to hypothetical negative-UED methodologies whose
 218 aim could be automatic design of environment distributions
 219 that result in minimization of learning agent’s (victim’s) per-
 220 formance. However, a straightforward negative of existing
 221 UED methodologies is not equivalent to environment-
 222 poisoning attacks. First of all, adversary-based and evolution-
 223 based UED create challenging environments by aiming to
 224 minimize learning agent’s rewards, while sample-based UED
 225 sample environments with high learning potential. All these
 226 mechanisms serve to improve the learning agent’s perfor-
 227 mance. A straightforward opposite of these approaches will
 228 only lead to reduced victim rewards. Environment-poisoning
 229 attacks on the other hand strive to push the victim to adopt an
 230 attacker-desired target behavior employing minimal attacker
 231 effort, while the victim trains to optimize its own objective.
 232 Moreover, unlike UED, in environment-poisoning attacks,
 233 the sequence of modifications is critical. While the victim
 234 begins training in its environment, the attacker begins to peri-
 235 odically modify the victim’s environment. Since the attacker
 236 intends to push the victim towards a target behavior that the
 237 victim has a low tendency of adopting by itself, in the orig-
 238 inal environment; the attacker must strategize modifications
 239 by carefully considering the current behaviors being learned
 240 by the victim. Moreover, the attack paradigm imposes con-

straints on the attacker in terms of the permissible magnitude
 241 of environment changes allowed at a single step. Any sizeable
 242 modification to the environment can therefore only be imple-
 243 mented through a sequence of small changes. The attack’s
 244 sensitivity to current victim behaviors as well as magnitude
 245 constraints render the sequence of environment changes crit-
 246 ical in the attack domain.

C Attacker’s Reward

The aforementioned Implicit Collective, Collective and
 249 Swarm Collective settings enable the victim populations to
 250 practice different levels of individual vs social learning. To
 251 show that environment-poisoning attacks can be effective un-
 252 der different levels of uncertainty, Implicit Collective set-
 253 ting is set as a Blackbox setting for the adversary while
 254 Collective and Swarm Collective settings are set as Ultra-
 255 Blackbox settings. In both Blackbox and Ultra-Blackbox, the
 256 attacker must extract individual victim behaviors from across-
 257 policy, on-process behavior traces, and additionally, in Ultra-
 258 Blackbox setting, the attacker is also not supported with any
 259 extrinsic reward signal, during its training. The attacker must
 260 therefore construct an intrinsic reward and estimate the effi-
 261 cacy of its own attack strategy by observing the change in the
 262 victim population’s behavior after each attack step, i.e. the
 263 victim population’s behavior adaptation in response to the at-
 264 tacker’s sequential attack.

This work extends the reward space design of [Xu *et al.*,
 266 2021] to the multi-victim setting by treating the underlying
 267 environment dynamics coupled with each population mem-
 268 ber’s behavior as a stochastic Markov process over state-
 269 action pairs $P_{k,u_i}(s_{j+1}, a_{j+1} | s_j, a_j)$. Herein, j is used to
 270 denote victim-level time step, just as i has been used to de-
 271 note attacker-level time step. Construction of this victim pro-
 272 cess enables computation of the joint probability distribution
 273 of current environment dynamics and the current victim be-
 274 havior which is critical as environment dynamics can impede
 275 certain victim tendencies while facilitating others. Kullback
 276 Leibler Divergence Rate, D_{klr} [Rached *et al.*, 2004]) between
 277 the ideal (in accordance with the attacker’s objectives) and
 278 the current process of each member of the victim popula-
 279 tion is computed and their negative mean is taken as the at-
 280 tacker’s reward for the given time-step. Herein the ideal pro-
 281 cess $P_{k,u_0}^*(s_{j+1}, a_{j+1} | s_j, a_j)$ is one where victim k adopts
 282 the attacker-desired behavior, τ_{k,u_i}^* aka the target behavior,
 283 without requiring any modifications to the underlying envi-
 284 ronment dynamics T_{u_0} . In this work, the target behavior is
 285 partial in nature:

$$\tau_{k,u_i}^*(s_n) = \begin{cases} a_n^* & s_n \in S^* \\ \tau_{k,u_i}(s_n) & s_n \notin S^* \end{cases} \quad (1)$$

Here S^* is the target state set, a_n^* is the target action for
 287 target state s_n , and $\tau_{k,u_i}(s_n)$ is the observed behavior of vic-
 288 tim k as it trains in environment with transition dynamics T_{u_i} .
 289 τ_{k,u_i}^* is the partial target behavior of victim k . The partiality
 290 of this design allows the attacker to specify target behavior for
 291 only those states that it cares about. This design, unlike the
 292 complete target behavior design proposed in [Rabinovich *et*
 293

294 *al.*, 2010], is efficient in high-dimensional discrete and
295 continuous state spaces.

296 The current process for victim k is defined using the cur-
297 rent underlying (modified) dynamics T_{u_i} as well as victim k 's
298 behavior after training in this modified environment τ_{k,u_i} . In
299 the Blackbox setting, the attacker is provided with an extrin-
300 sic reward signal that can access a victim k 's Q function to
301 compute its policy π_{k,u_i} ; while in the Ultra-Blackbox setting,
302 the attacker constructs its own intrinsic reward signal using
303 the behavior traces of victim k . Herein, the attacker approxi-
304 mates π_{k,u_i} as a distribution using the last h actions taken by
305 victim k in each state. The current process P_{k,u_i} of victim
306 k in terms of current environment dynamics and victim k 's
307 current policy is defined a:

$$P_{k,u_i}(s_{j+1}, a_{j+1}|s_j, a_j) = T_{u_i}(s_{j+1}|s_j, a_j) \pi_{k,u_i}(a_{j+1}|s_{j+1}) \quad (2)$$

308 The target policy π_{k,u_i}^* is constructed by taking a copy of
309 π_{k,u_i} and modifying it by assigning probability 1.0 to all tar-
310 get actions (and 0.0 to non-target actions) w.r.t. each target
311 state. The ideal process P_{k,u_0}^* of victim k in terms of default
312 dynamics and target policy is defined as:

$$P_{k,u_0}^*(s_{j+1}, a_{j+1}|s_j, a_j) = T_{u_0}(s_{j+1}|s_j, a_j) \pi_{k,u_i}^*(a_{j+1}|s_{j+1}) \quad (3)$$

313 The reward measuring the current performance of the at-
314 tack strategy is then defined as negative of the average Kull-
315 back Leibler Divergence Rate between the ideal process and
316 the current processes of all agents in the victim population. In
317 equation 4 given below q_i is the stationary state distribution
318 of the victim environment with dynamics T_{u_i} .

$$\begin{aligned} r_{\alpha,i} &= -(1/K) \sum_k D_{klr}(P_{k,u_0}^* || P_{k,u_i}) \\ D_{klr}(P_{k,u_0}^* || P_{k,u_i}) &= \sum_j q_i(s_j) \pi_{k,u_i}(a_j|s_j) D_i^{kl}(s_j, a_j) \\ D_i^{kl}(s_j, a_j) &= \sum_{s_{j+1}, a_{j+1}} \left(P_{k,u_i}(s_{j+1}, a_{j+1}|s_j, a_j) \right. \\ &\quad \left. \log \frac{P_{k,u_i}(s_{j+1}, a_{j+1}|s_j, a_j)}{P_{k,u_0}^*(s_{j+1}, a_{j+1}|s_j, a_j)} \right) \\ q_i(s_{j+1}) &= \sum_j q_i(s_j) \pi_{k,u_i}(a_j|s_j) T_{u_i}(s_{j+1}|s_j, a_j) \end{aligned} \quad (4)$$

319 D Experimental Setting - Victim Population

320 Each member of the victim population utilises independent Q
321 Learning under Implicit Collective setting and DQN (Deep-
322 Q Networks) under Collective and Swarm Collective settings.
323 Under Implicit Collective setting, in each attack step, the vic-
324 tim population trains for 80 episodes using Q learning with
325 discount factor $\gamma = 0.90$, and learning rate $\alpha = 0.100$. On
326 the other hand, under Collective and Swarm Collective set-
327 tings, the victim population trains for 40 episodes in each

328 attack step, using DQN with discount factor $\gamma = 0.99$,
329 and learning rate $\alpha = 0.001$. These values were chosen
330 such that the victim population converges to optimal behav-
331 ior (w.r.t victim population's objectives) in the default (un-
332 attacked) environment. Furthermore, it is important to note
333 that the state-of-the-art single-agent environment poisoning
334 work, TEPA [Xu *et al.*, 2021] whose multi-victim extension is
335 constructed and utilized as baseline in this paper, tested their
336 attack against an ϵ -Greedy victim that is heavily biased to-
337 wards exploitation from the beginning of the training period.
338 However, we test the barycenter-based attack against more-
339 realistic SoftMax victims that spend substantial time explor-
340 ing the environment when they begin their training.

341 E Experimental Setting - Attacker

342 The attacker in the baseline as well as proposed methodology
343 is trained using Deep Deterministic Policy Gradient (DDPG)
344 with discount factor $\gamma = 0.95$, and target network update
345 rate $\tau = 0.005$. TEPA attacker collects data for 100 attack
346 episodes before beginning to train while we reduce the initial
347 data collection period to 30 attack episodes in the barycenter-
348 based attacker. The policy function of the attacker is a fully
349 connected, feedforward neural network with specification:
350 INPUT(21)-FC(400)-ReLU-FC(300)-ReLU-FC(16)-Tanh.

351 TEPA attacker uses an auto-encoder to capture latest vic-
352 tim behaviors. Its encoder is a fully connected, feedfor-
353 ward neural network with specification: INPUT(12)-FC(36)-
354 ReLU-FC(36)-ReLU-FC(5). The proposed barycenter-based
355 attacker on the other hand utilises a variational auto-encoder
356 (VAE) whose encoder is a fully connected, feedforward neu-
357 ral network with specification: INPUT(32)-FC(36)-ReLU-
358 FC(36)-ReLU-FC(5). The learning rates of the two models
359 are 0.001 and 0.00001 respectively. It is important to note
360 here that unlike the pre-trained VAE employed in the pro-
361 posed methodology, the auto-encoder in TEPA trains along
362 with the attacker in order to ensure that the latest victim be-
363 haviors are mapped to the latent space with high accuracy.
364 To enable evaluation of TEPA-based concatenation strategies,
365 the auto-encoder parameters are saved during training, after
366 every 20 attack episodes. At the time of evaluation of a par-
367 ticular strategy, the auto-encoder parameters saved closest to
368 the given strategy are used by the attacker.

369 F Performance Metrics

370 The performance of the attacker is measured in terms of the
371 accuracy of the victim population's adoption of the target be-
372 havior (Attack Accuracy and Attack SoftMax Accuracy) as
373 well as the cumulative changes brought about in the victim
374 environment by the attacker (Attacker Effort).

375 Attack Accuracy computes the level of adoption of the tar-
376 get behavior by the victim population in the given environ-
377 ment. Given that π_{k,u_i} is victim k 's probabilistic policy in
378 environment with dynamics T_{u_i} , K is the number of agents in
379 the victim population and S^* is the set of target states, let
380 N^* be the number of target states, a_n^* be the target action in
381 target state s_n , and $f_a(s, \pi_{k,u_i})$ be a function that takes a tar-
382 get state and a victim policy as input and outputs 1 if the given

383 policy assigns highest probability to the target action in that
 384 state and 0 otherwise.

$$f_a(s_n, \pi_{k,u_i}) = \begin{cases} 1 & \pi_{k,u_i}(a_n^*|s_n) > \pi_{k,u_i}(a_n|s_n) \forall a_n \\ 0 & \text{otherwise} \end{cases} \quad (5)$$

$$\text{Attack Accuracy} = \frac{1}{K} \sum_k \left(\frac{1}{N^*} \sum_{s_n \in S^*} f_a(s_n, \pi_{k,u_i}) \right) \quad (6)$$

385 Attack SoftMax Accuracy computes the strength with
 386 which the target behavior is adopted by the victim popula-
 387 tion by calculating the probability assigned to the target path
 388 by the victim population.

Attack SoftMax Accuracy

$$= \frac{1}{K} \sum_k \left(\frac{1}{N^*} \sum_{s_n \in S^*} \pi_{k,u_i}(a_n^*|s_n) \right) \quad (7)$$

389 Attacker Effort computes the degree to which the at-
 390 tacker modifies the victim environment. Let $G =$
 391 $[g_{u_i}^1, g_{u_i}^2, g_{u_i}^3, \dots, g_{u_i}^M]$ be the altitudes of the M grid cells in
 392 victim environment with dynamics T_{u_i} .

$$\text{Attacker Effort} = \frac{1}{M} \sum_{m=1}^M |g_{u_i}^m - g_{u_{i-1}}^m| \quad (8)$$

393 G Attack Strategy Selection

394 The attacker training episodes are 15-step sequential attacks
 395 on freshly initialized victim populations. Each episode begins
 396 with a slightly perturbed environment so that the attacker's re-
 397 ward that aims to minimize distance between current and de-
 398 fault environment (along with minimizing distance between
 399 current and target population behavior) does not get stuck on
 400 the default environment without ensuring victim population's
 401 adoption of target behavior. The new population begins train-
 402 ing to maximize its objectives on this slightly perturbed en-
 403 vironment. At attack time step 1, the attacker makes its first
 404 modification of the victim environment. After each episode,
 405 the attack strategy employed in that episode is saved if it is
 406 better or equal to the best attack strategy found so far, with
 407 respect to last-timestep, mean or cumulative value of at least
 408 one strategy quality criterion. A given strategy's quality is
 409 approximated using 3 internal and 5 external quality criteria.
 410 Herein criteria that are approximated by the attacker are re-
 411 ferred to as internal while criteria computed by the external
 412 system for the purpose of training the attacker are termed ex-
 413 ternal. The 8 quality criteria are:

- 414 1. KLR-Full (Internal): KLR between current (current env
 * current victim behavior) and perfect process (default
 env * target victim behavior)
- 415 2. KLR-Environment (Internal): KLR between (current
 env * target victim behavior) and perfect process (de-
 fault env * target victim behavior)

3. KLR-Behavior (Internal): KLR between (default env *
 current victim behavior) and perfect process (default env
 * target victim behavior)
4. Attack Accuracy (External): same as Attack Accuracy
5. Attack Softmax Accuracy (External): same as Attack
 SoftMax Accuracy performance metric
6. Attack Partial-Softmax Accuracy (External): unlike At-
 tack SoftMax Accuracy where probability of attacker-
 desired actions in all attacker-desired states are added,
 here probability of attacker-desired actions in only those
 attacker-desired states are added where Attack Accuracy
 is 1.0 i.e. only probability of those attacker-desired ac-
 tions are taken into account which are already assigned
 maximum probability by the victim. This quality crite-
 rion enables the attacker to identify (and thereby save)
 strategies that are capable of inducing strong adoption
 of target behavior but were not able to achieve this in all
 target states in the given trial.
7. Attacker Effort (External): same as Attacker Effort per-
 formance metric
8. Attack Time (External): computational time correspond-
 ing to each attack step (inclusive of time taken by the
 victims to train in the attacked/poisoned/modified envi-
 ronment).

The experiment graphs present in the main body of the paper demonstrate performance of the best attack strategies found by the different models. These best strategies are selected by prioritising Attack Accuracy, as the main goal of this work is to find strategies that push victim populations to adopt the target behavior. The strength of target behavior adoption (Attack Softmax Accuracy) and amount of changes made to the environment (Attacker Effort) in order to achieve this, demonstrate additional capabilities of the attack strategies. The selection procedure begins with a filtering process that shortlists only those strategies that acquire Attack Accuracy of 1.0 by the end of the attack. The shortlisted strategies are then arranged in decreasing order of mean Attack Accuracy. Thereafter, the top $20 + j$ strategies are selected where all strategies between ranks 20 and $20+j$ have the same mean Attack Accuracy. In case less than 20 strategies acquire Attack Accuracy 1.0 by the end of the attack, then top $20+j$ strategies are selected solely on the basis of last-step Attack Accuracy. Lastly, each of these $20+j$ strategies is utilised to attack 10 same-sized populations that utilize the same random seed (Implicit Collective) or utilize neural networks initialized with random numbers from the same range (Collective and Swarm Collective); as used during training.

468 H Experimental Data (Expanded)

This appendix includes results corresponding to experiments under the Swarm Collective setting as well as additional results corresponding to experiments under Implicit Collective and Collective settings.

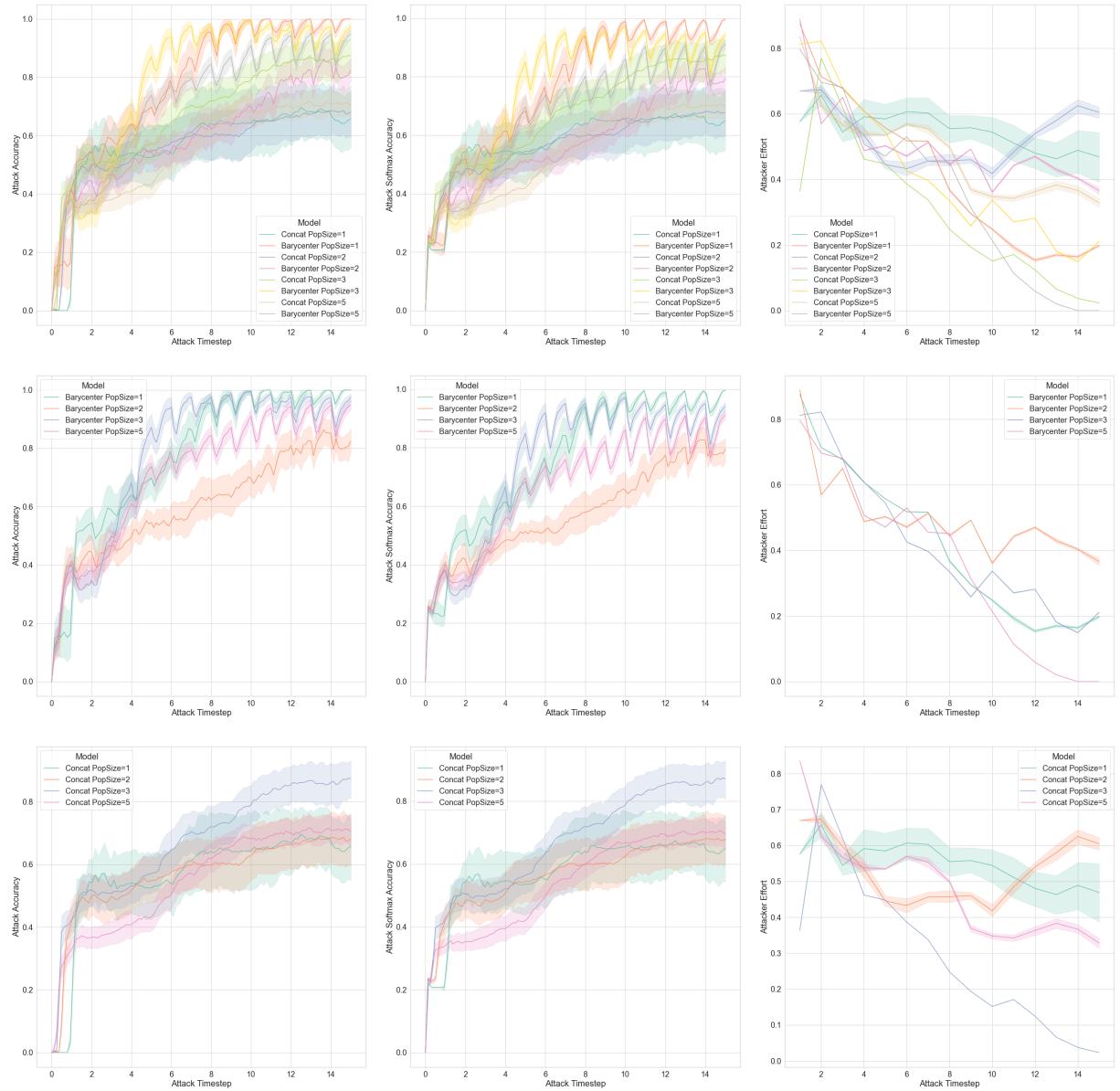


Figure 2: Accuracy, Softmax Accuracy and Effort of attack, trained and tested on same sized Swarm Collective victim

473 H.1 Swarm Collective Setting

474 Experiment A (see Figure 2) tests the capabilities of the pro-
475 posed and baseline methods to attack Swarm Collectives of
476 sizes 1, 2, 3 and 5. In this setting, majority of both concatena-
477 tion and barycenter based strategies achieve better accuracies
478 while exerting higher effort, compared to their Collective set-
479 ting counterparts. All baseline strategies except the strategy
480 trained on size-3 populations, achieve Attack Accuracy and
481 Attack SoftMax Accuracy between 0.6 and 0.8, in contrast to
482 the Collective setting, wherein they achieve both accuracies
483 between 0.4 and 0.6 by the end of the attack. Size-3 baseline
484 strategy achieves 0.9 Attack *and* Attack-SoftMax accuracy
485 by the end of the attack in both Collective and Swarm Col-
486 lective settings. Higher accuracies for Swarm Collectives are
487 accompanied with higher Attacker Effort as all strategies ex-
488 cept size-3 strategy exert effort between 0.3 and 0.6, in con-
489 trast to Collective setting wherein they reduce Attacker Ef-
490 fort to below 0.2 by the end of the attack. Size-3 strategy in
491 both Collective and Swarm Collective settings converge to-
492 wards 0.0 Attacker Effort by the end of the attack. This shows
493 that the baseline method that uses concatenation of trajectory-
494 based individual victims’ behaviors to capture the population
495 behavior is less capable of accurately understanding popula-
496 tion behavior of populations of sizes $\neq 3$. Barycenter-based
497 strategies that achieve high Attack Accuracy of above 0.8 and
498 high Attack Softmax Accuracy of above 0.7 in Collective set-
499 ting achieve even higher accuracies for Swarm Collectives
500 with majority reaching above 0.9 by the end of the attack.
501 This performance boost comes with slightly higher Attacker
502 Effort which is between 0.2 and 0.4 for Swarm Collectives
503 and between 0.0 and 0.4 in the Collective setting. Lastly, like
504 Collective setting, barycenter-based attacks exhibit the climb-
505 ing zig-zag behavior in terms of Attack and Attack-Softmax
506 Accuracies for Swarm Collectives.

507 Experiment B, against Swarm Collectives, is presented in
508 Figures 3 and 6. Figure 3 presents strategies trained on pop-
509 ulations of sizes 3, 5, 10, and 20, and tested on populations
510 of sizes 3, 5, 10, and 20 in the same graph in order to demon-
511 strate that the proposed barycenter-based method showcases
512 better size-agnosticity verus Swarm Collectives, than the Col-
513 lective setting. On the other hand, figure 6 presents separate
514 graphs for strategies trained on populations of sizes 3, 5, 10,
515 and 20 respectively. Each strategy is tested on populations
516 of sizes 1, 2, 3, 4, 5, 10, and 20 to facilitate understanding of
517 how a particular strategy’s performance changes with the test-
518 population size. Last-step Attack Accuracy of size-3 strat-
519 egy (strategy trained on size-3 populations) is highest when
520 tested on size-3 populations, slightly lower for size 1 and 2
521 populations, and progressively lower for larger populations.
522 Last-step Attack Accuracy of all tests carried using size-5
523 strategy is higher than their size-3 counterparts and show the
524 same trend of decreasing with increasing population size. The
525 climbing zig-zag behavior of size-10 strategy degrades with
526 decreasing test-population size starting with size-5 test pop-
527 ulations. Moreover, last-step Attack Accuracy of almost all
528 tests carried using size-10 strategy is lower than their size-
529 3 and size-5 counterparts. The climbing zig-zag behavior of
530 tests carried out using size-20 strategy degrades further than
531 their size-10 counterparts, on small populations (1-5). How-

ever, performance of size-20 strategy on tests carried out on
532 large populations (10,20) is better than the performance of all
533 other strategies. These results imply that strategies trained
534 on smaller populations and tested on larger populations per-
535 fectly retain the climbing zig-zag behavior but this behavior
536 slightly degrades for strategies trained on larger and tested on
537 smaller populations. Attack SoftMax Accuracy shows simi-
538 lar trends to Attack Accuracy while Attacker Effort remains
539 largely unchanged across tests on different sized test popula-
540 tions. Lastly, even while taking into account all the variations
541 discussed above, all strategies showcase high level of size-
542 agnosticity as all performance plots follow the same trend
543 and remain within a certain range to the plot corresponding
544 to attacks trained and tested on the same-sized populations.
545

546 H.2 Implicit Collective and Collective Settings

547 Figures 5 and 4 present Experiment B against Collectives and
548 Implicit Collectives, respectively; for all strategies tested on
549 populations of sizes 1-5, 10, and 20 instead of testing only
550 on 3, 5, 10 and 20 as done in the main paper. Lesser pop-
551 ulations were shown in the main paper to enable readers to
552 catch overall trends quickly while more extensive results are
553 presented here for completeness. Figure 4 shows that in con-
554 trast to other barycenter-based strategies that converge within
555 2 attack steps, size-1 strategy only converges by the end of
556 the attack in Implicit Collective setting. However, even then,
557 size-1 strategy is size-agnostic and achieves last-step Attack
558 Accuracy of 1.0 when used to attack populations of sizes 1-5,
559 10, and 20. Other than this result, figures 4 and 5 are similar
560 to their “clearer” counterparts presented in the main paper.

561 Figure 7 presents separate graphs for size 3, 5, 10, and 20
562 strategies, respectively. Each strategy is tested on popula-
563 tions of sizes 1-5, 10, and 20 to facilitate understanding of
564 how a particular strategy’s performance changes with test-
565 population size. Last-step Attack Accuracy of size-3 strat-
566 egy mostly decreases with increasing test-population size but
567 stays above 0.6 until size-10 test populations and drops to
568 0.3 for size-20 test-populations. Last-step Attack Accuracy
569 of size-5 strategy reaches 1.0 for small populations (1-5);
570 above 0.8 for size-10 test-populations; and above 0.4 for
571 size-20 test-populations. For size-10 strategy, last-step At-
572 tack Accuracy of small populations (1-5) drops to between
573 0.8 and 0.9; stays the same for size-10 test-populations and
574 climbs to 0.6 for size-20 test-populations. Lastly, last-step
575 Attack Accuracy of size-20 strategy climbs to above 0.9 for
576 small and medium populations (1-10) and to 0.8 for size-
577 20 large test-populations. Also, Attack and Attack SoftMax
578 Accuracies show similar trends. These results imply that
579 barycenter-based method against Collectives achieves better
580 size-agnosticity w.r.t. accuracies when trained on larger pop-
581 ulations. On the other hand, Attacker Effort reduces to be-
582 low 0.2 for size 3 and 5 strategies; between 0.2 and 0.4 for
583 size 10 strategy; and between 0.3 and 0.5 for size-20 strat-
584 egy by the end of the attack. In addition, Attacker Effort
585 shows higher variability across test-populations with increas-
586 ing training-population size. Therefore, barycenter-based
587 method achieves better size-agnosticity w.r.t. accuracies at
588 the cost of higher and lesser size-agnostic effort with increas-
589 ing training-population size.

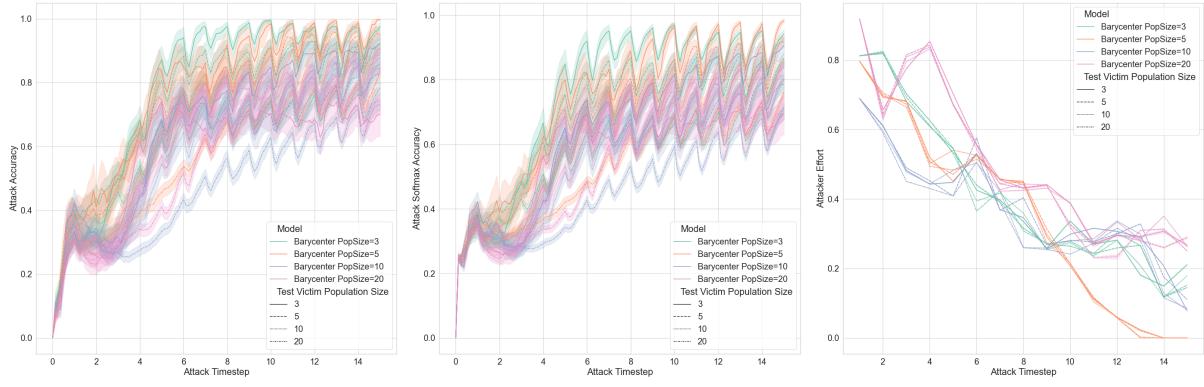


Figure 3: Accuracy, Softmax Accuracy and Effort of Barycenter attacks, tested on Swarm Collectives of sizes 3,5,10 and 20

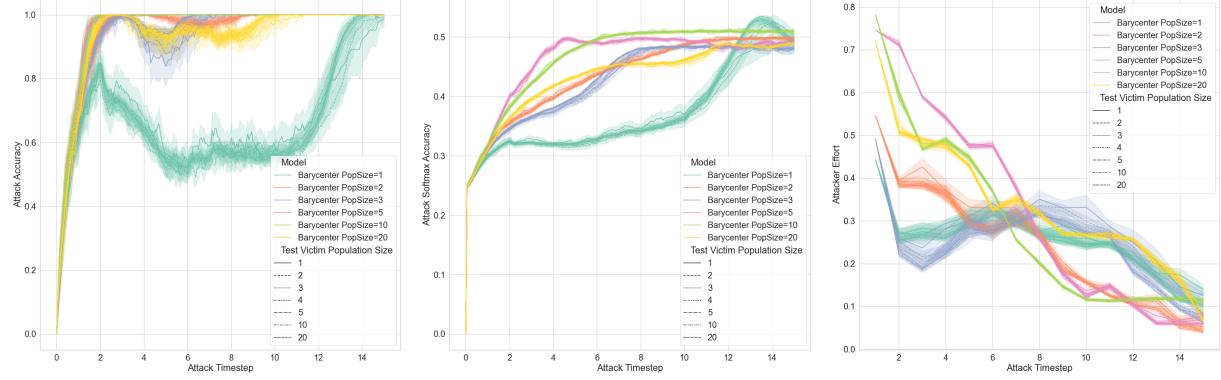


Figure 4: Accuracy, Softmax Accuracy and Effort of Barycenter attacks, tested on Implicit Collectives of sizes 1,2,3,4,5,10 and 20

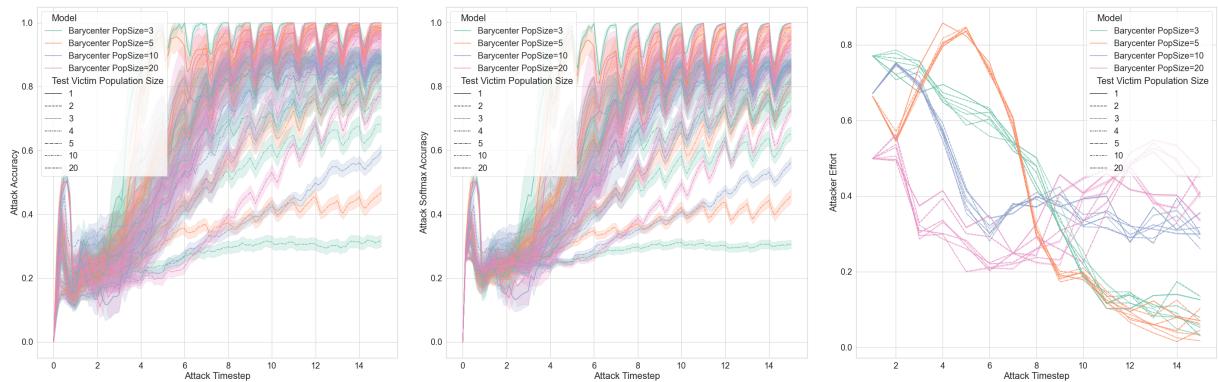


Figure 5: Accuracy, Softmax Accuracy and Effort of Barycenter attacks, tested on Collective victim populations of sizes 1,2,3,4,5,10 and 20

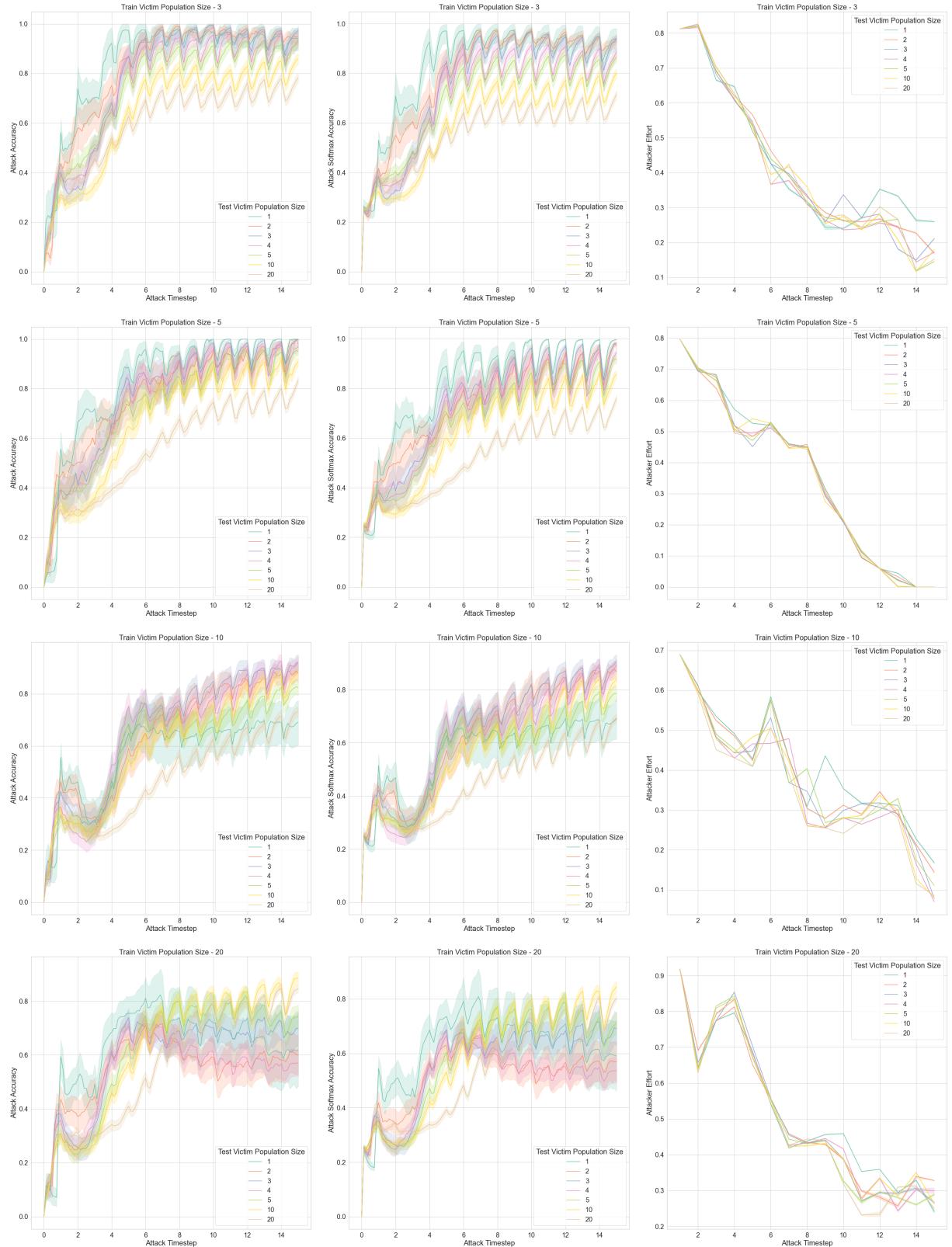


Figure 6: Accuracy, Softmax Accuracy and Effort of Barycenter attacks, tested on Swarm Collective victim populations of sizes 1,2,3,4,5,10 and 20, presented in separate graphs

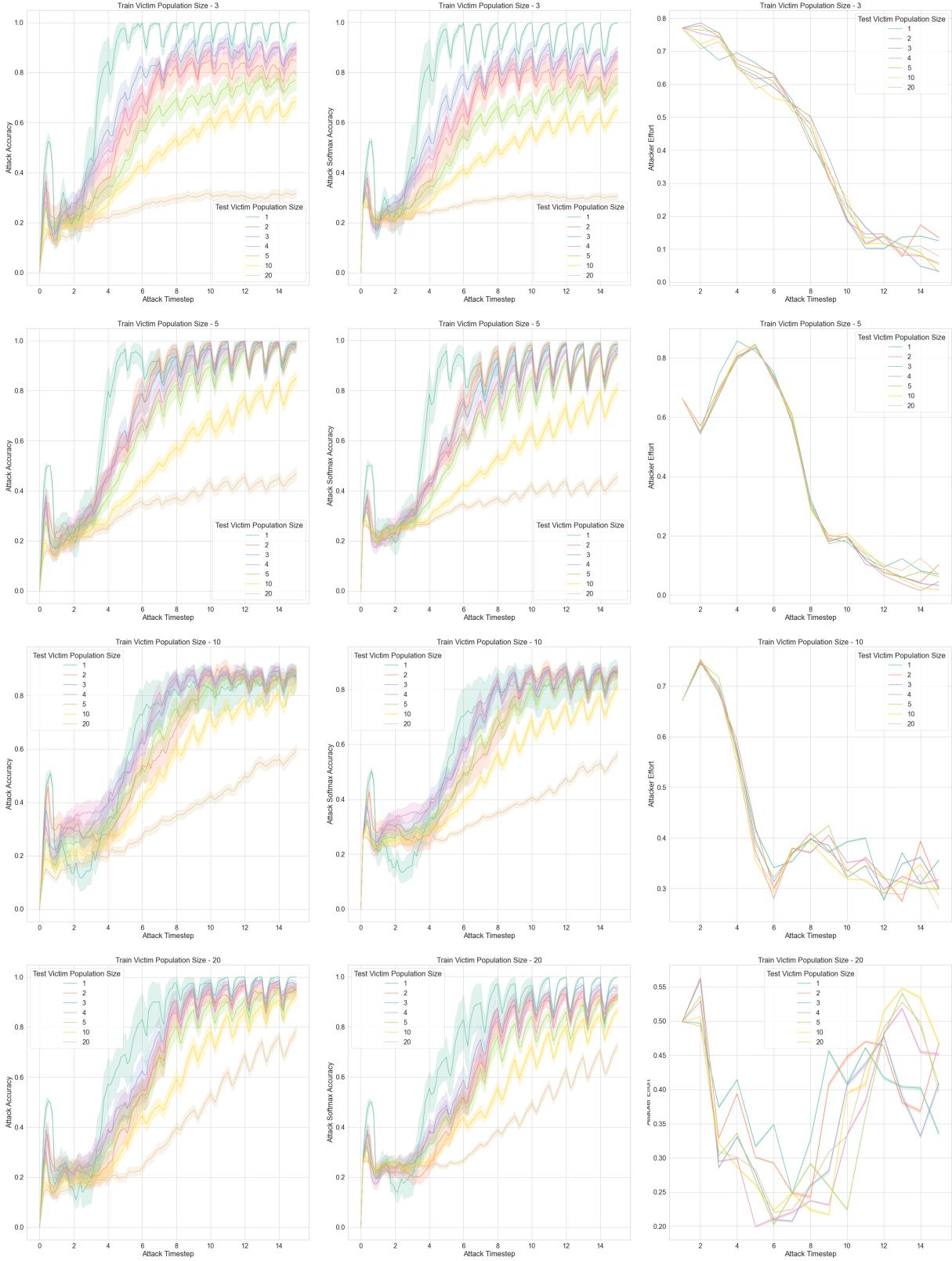


Figure 7: Accuracy, Softmax Accuracy and Effort of Barycenter attacks, tested on Collective victim populations of sizes 1,2,3,4,5,10 and 20, presented in separate graphs

References

- [Daw *et al.*, 2005] Nathaniel D Daw, Yael Niv, and Peter Dayan. Uncertainty-based competition between prefrontal and dorsolateral striatal systems for behavioral control. *Nature neuroscience*, 8(12):1704–1711, 2005.
- [Dennis *et al.*, 2020] M. Dennis, N. Jaques, E. Vinitsky, A. Bayen, S. Russell, A. Critch, and S. Levine. Emergent complexity and zero-shot transfer via unsupervised environment design. In *NeurIPS*, pages 13049–13061, 2020.
- [Dimakopoulou *et al.*, 2018] M. Dimakopoulou, I. Osband, and B. Van Roy. Scalable coordinated exploration in concurrent reinforcement learning. In *NeurIPS*, 2018.
- [Foerster *et al.*, 2018] J. Foerster, R. Y. Chen, M. Al-Shedivat, S. Whiteson, P. Abbeel, and I. Mordatch. Learning with opponent-learning awareness. In *AAMAS*, pages 122–130, 2018.
- [Grover *et al.*, 2018] A. Grover, M. Al-Shedivat, J. Gupta, Y. Burda, and H. Edwards. Learning policy representations in multiagent systems. In *ICML*, pages 1802–1811, 2018.
- [He *et al.*, 2016] H. He, J. Boyd-Graber, K. Kwok, and H. Daumé III. Opponent modeling in deep reinforcement learning. In *ICML*, pages 1804–1813, 2016.
- [Jiang *et al.*, 2021a] M. Jiang, M. Dennis, J. Parker-Holder, J. Foerster, E. Grefenstette, and T. Rocktäschel. Replay-guided adversarial environment design. In *NeurIPS*, pages 1884–1897, 2021.
- [Jiang *et al.*, 2021b] M. Jiang, E. Grefenstette, and T. Rocktäschel. Prioritized level replay. In *ICML*, pages 4940–4950, 2021.
- [Lee *et al.*, 2019] J. Lee, Y. Lee, J. Kim, A. Kosiorek, S. Choi, and Y. W. Teh. Set transformer: A framework for attention-based permutation-invariant neural networks. In *ICML*, pages 3744–3753, 2019.
- [Lupu and Precup, 2020] A. Lupu and D. Precup. Gifting in multi-agent reinforcement learning. In *AAMAS*, pages 789–797, 2020.
- [Marthi *et al.*, 2005] B. Marthi, S. Russell, D. Latham, and C. Guestrin. Concurrent hierarchical reinforcement learning. In *IJCAI*, pages 779–785, 2005.
- [Matiisen *et al.*, 2019] T. Matiisen, A. Oliver, T. Cohen, and J. Schulman. Teacher–student curriculum learning. *IEEE transactions on neural networks and learning systems*, 31(9):3732–3740, 2019.
- [Mialon *et al.*, 2020] G. Mialon, D. Chen, A. d’Aspremont, and J. Mairal. A trainable optimal transport embedding for feature aggregation. In *ICLR*, 2020.
- [Morimoto and Doya, 2005] J. Morimoto and K. Doya. Robust reinforcement learning. *Neural computation*, 17(2):335–359, 2005.
- [Papoudakis and Albrecht, 2020] G. Papoudakis and S. V. Albrecht. Variational autoencoders for opponent modeling in multi-agent systems. In *AAAI WS on RL in Games*, 2020.
- [Papoudakis *et al.*, 2021] G. Papoudakis, F. Christianos, and S. Albrecht. Agent modelling under partial observability for deep reinforcement learning. pages 19210–19222, 2021.
- [Parisotto *et al.*, 2019] E. Parisotto, S. Ghosh, S. B. Yalamanchi, V. Chinnaobireddy, Y. Wu, and R. Salakhutdinov. Concurrent meta reinforcement learning. *arXiv preprint arXiv:1903.02710*, 2019.
- [Parker-Holder *et al.*, 2022] J. Parker-Holder, M. Jiang, M. Dennis, M. Samvelyan, J. Foerster, E. Grefenstette, and T. Rocktäschel. Evolving curricula with regret-based environment design. *arXiv preprint arXiv:2203.01302*, 2022.
- [Qi *et al.*, 2017] C. R. Qi, H. Su, K. Mo, and L. J. Guibas. Pointnet: Deep learning on point sets for 3d classification and segmentation. In *CVPR*, pages 652–660, 2017.
- [Rabinovich *et al.*, 2010] Z. Rabinovich, L. Dufton, K. Larsson, and N. R. Jennings. Cultivating desired behaviour: Policy teaching via environment-dynamics tweaks. In *AAMAS*, pages 1097–1104, 2010.
- [Rabinowitz *et al.*, 2018] N. Rabinowitz, F. Perbet, F. Song, C. Zhang, S.M. A. Eslami, and M. Botvinick. Machine theory of mind. In *ICML*, pages 4218–4227, 2018.
- [Rached *et al.*, 2004] Z. Rached, F. Alajaji, and L. L. Campbell. The kullback-leibler divergence rate between markov sources. *IEEE Transactions on Information Theory*, 50(5):917–921, 2004.
- [Raileanu *et al.*, 2018] R. Raileanu, E. Denton, A. Szlam, and R. Fergus. Modeling others using oneself in multi-agent reinforcement learning. In *ICML*, pages 4257–4266, 2018.
- [Shum *et al.*, 2019] M. Shum, M. Kleiman-Weiner, M. L. Littman, and J. B. Tenenbaum. Theory of minds: Understanding behavior in groups through inverse planning. In *AAAI*, pages 6163–6170, 2019.
- [Skianis *et al.*, 2020] K. Skianis, G. Nikolentzos, S. Limnios, and M. Vazirgiannis. Rep the set: Neural networks for learning set representations. In *AISTATS*, pages 1410–1420, 2020.
- [Tacchetti *et al.*, 2019] A. Tacchetti, F. H. Song, P. A. M. Mediano, V. Zambaldi, J. Kramár, N. C. Rabinowitz, T. Graepel, M. Botvinick, and P. W. Battaglia. Relational forward models for multi-agent learning. In *ICLR*, 2019.
- [Terry *et al.*, 2020] Justin K Terry, Nathaniel Grammel, Ananth Hari, Luis Santos, and Benjamin Black. Revisiting parameter sharing in multi-agent deep reinforcement learning. *arXiv preprint arXiv:2005.13625*, 2020.
- [Tobin *et al.*, 2017] J. Tobin, R. Fong, A. Ray, J. Schneider, W. Zaremba, and P. Abbeel. Domain randomization for transferring deep neural networks from simulation to the real world. In *IROS*, pages 23–30, 2017.
- [Xu *et al.*, 2021] H. Xu, R. Wang, L. Raizman, and Z. Rabinovich. Transferable environment poisoning: Training-time attack on reinforcement learning. In *Proceedings of the 20th International Conference on Autonomous Agents and MultiAgent Systems*, pages 1398–1406, 2021.

- 699 [Yaman *et al.*, 2022] Anil Yaman, Nicolas Bredeche, Onur
700 Çaylak, Joel Z Leibo, and Sang Wan Lee. Meta-control
701 of social learning strategies. *PLoS computational biology*,
702 18(2):e1009882, 2022.
- 703 [Yang *et al.*, 2020] J. Yang, A. Li, M. Farajtabar, P. Sune-
704 hag, E. Hughes, and H. Zha. Learning to incentivize other
705 learning agents. In *NeurIPS*, pages 15208–15219, 2020.
- 706 [Zaheer *et al.*, 2017] M. Zaheer, S. Kottur, S. Ravanbakhsh,
707 B. Poczos, R. R. Salakhutdinov, and A. J. Smola. Deep
708 sets. In *NeurIPS*, 2017.
- 709 [Zhang and Lesser, 2010] C. Zhang and V. Lesser. Multi-
710 agent learning with policy prediction. In *AAAI*, 2010.