

Machine Theory of Mind for Autonomous Cyber-Defence

Luke Swaby^{1,*}, Matthew Stewart¹, Daniel Harrold¹, Chris Willis¹, and Gregory Palmer¹

¹BAE Systems Applied Intelligence Labs, Chelmsford Office & Technology Park, Chelmsford, Essex, CM2 8HN.

*luke.swaby2@baesystems.com

ABSTRACT

Intelligent autonomous agents hold much potential for the domain of cyber security. However, due to many state-of-the-art approaches relying on uninterpretable black-box models, there is growing demand for methods that offer stakeholders clear and actionable insights into their latent beliefs and motivations. To address this, we evaluate Theory of Mind (ToM) approaches for Autonomous Cyber Operations. Upon learning a robust prior, ToM models can predict an agent's goals, behaviours, and contextual beliefs given only a handful of past behaviour observations. In this paper, we introduce a novel Graph Neural Network (GNN)-based ToM architecture tailored for cyber-defence, Graph-In, Graph-Out (GIGO)-ToM, which can accurately predict both the targets and attack trajectories of adversarial cyber agents over arbitrary computer network topologies. To evaluate the latter, we propose a novel extension of the Wasserstein distance for measuring the similarity of graph-based probability distributions. Whereas the standard Wasserstein distance lacks a fixed reference scale, we introduce a graph-theoretic normalization factor that enables a standardized comparison between networks of different sizes. We furnish this metric, which we term the *Network Transport Distance* (NTD), with a weighting function that emphasizes predictions according to custom node features, allowing network operators to explore arbitrary strategic considerations. Benchmarked against a Graph-In, Dense-Out (GIDO)-ToM architecture in an abstract cyber-defence environment, our empirical evaluations show that GIGO-ToM can accurately predict the goals and behaviours of various unseen cyber-attacking agents across a range of network topologies, as well as learn embeddings that can effectively characterize their policies.

Introduction

As autonomous systems are entrusted with increasingly critical tasks across various industries, the need for explainability has become paramount.¹⁻⁶ Recent advancements have established deep neural networks (DNNs) as the cornerstone of autonomous systems due to their proficiency in handling the complex, high-dimensional data required to scale to real-world environments.⁷⁻¹⁰ However, in safety-critical domains where decisions can have far-reaching consequences, the opacity of these 'black-box' models poses significant risks, raising concerns around trust, accountability, and safety.^{11,12} One such domain is that of cybersecurity, where the proliferation of complex computer networks with ever-growing attack surfaces has fueled a cyber arms race with both attackers and defenders increasingly leveraging DNNs for their respective strategies.¹³⁻¹⁹ Consequently, there is growing urgency for methods to interpret and analyze the latent decision-making processes of these models.

In recent years, a plethora of approaches have emerged aimed at demystifying DNNs. In the context of cybersecurity, these include Explainable Artificial Intelligence (XAI) techniques for malware and anomaly detection,²⁰ as well as solutions that rationalize the actions of deep reinforcement learning agents.^{21,22} However, these approaches typically focus on post-hoc, data-centric clarifications of *how* a model arrived at a particular decision rather than *why*, and thus lack the level of contextual detail that humans generally deem important when making high-stakes decisions.^{11,23,24,24-26} This limitation becomes especially pronounced in multi-agent decision-making scenarios where understanding the intentions, strategies, and circumstances of other entities can significantly impact outcomes. For example, a computer network defender capable only of identifying the presence and mechanics of a given attack is at a significant disadvantage compared to one who can additionally infer the adversary's objectives. Without understanding *why* certain actions occur, the defensive strategies of the former risk being reactive and insufficient in anticipating evolving threats.

In the psychological literature, the human ability to infer the latent mental states of others is known as a '*Theory of Mind*' (ToM).²⁷ An extensive corpus of cognitive studies has linked the absence of ToM with significant social deficiencies in humans.²⁸⁻³² This has led several researchers to investigate whether the intersubjective capabilities of AI systems can be improved with ToM-inspired model architectures. This topic was first explored in *Machine Theory of Mind* (MToM) by Rabinowitz et al.,³³ who implemented an observing agent—a parameterized model, dubbed '*ToMnet*'—that employs meta-learning to make predictions over the goals, strategies, and characters of observed agents based solely on past behaviour

observations.

The ToMnet architecture (Figure 1) features 3 submodules that respectively parse: i) a set of past behavior observations; ii) behavior observations from a current episode up to a timestep t ; and iii) an observation of the state of the environment at t . The outputs of each module are propagated forward to a unified prediction mechanism that predicts next-step action probabilities $\hat{\pi}$, the object that the observed agent is targeting \hat{c} , and the expected state occupancy, or, *successor representation*, \hat{SR} .³⁴ Trained and evaluated in simple gridworld environments, Rabinowitz et al.³³ found that ToMnet was able to characterize the behaviours and intentions of a broad taxonomy of agents, and showed strong generalization capabilities.

Intuitively, ToMnet offers a clear and effective paradigm for understanding the decisions made by agents in a cybersecurity context. However, whilst it has been applied to specialised cyber-defence problems before (see related work), our review of existing literature found little work on developing general, network-agnostic, and scalable ToMnet formulations for this domain. As such, our aim in this paper is to provide an abstract evaluation of whether ToMnet can be effectively utilized for cyber-defence.

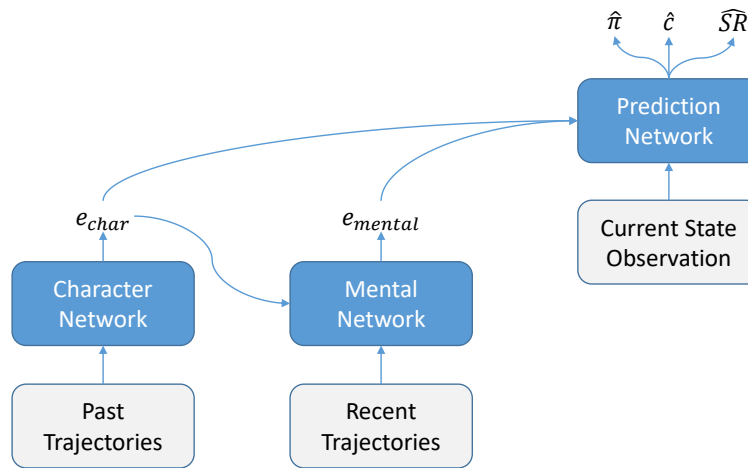


Figure 1. The original ToMnet architecture, adopted from Rabinowitz et al.³³ It consists of three components: a *character network*, a *mental network*, and a *prediction network*. The character network parses observations of past behaviour (e.g. previous episodes) to build a character embedding e_{char} . The mental network then uses this along with behavioural observations from a current episode to build a representation of the agent’s current beliefs and intentions: the mental state embedding e_{mental} . Finally, the predictor network utilizes both these embeddings along with a snapshot of the current state of the environment to forecast future behaviour: the next-step action probabilities $\hat{\pi}$, probabilities of whether certain objects will be consumed \hat{c} , and the predicted successor representations \hat{SR} .³⁴

Right away we can identify several limitations in ToMnet’s original formulation that present challenges when applying it to cyber-defence. First, in the original paper, the gridworld environment dimensions were fixed at 11×11 for all experiments, allowing ToMnet’s input and output layers to be fixed accordingly.³³ In cyber-defence scenarios, by contrast, the number of nodes in a network is not necessarily known *a priori*. The ToMnet architecture must therefore be adapted to handle variable sized inputs in order to qualify as a viable candidate for cyber applications. Furthermore, gridworld environments contain relatively few cells that lack individual features, and agents that can only perform a fixed, limited set of actions. Cyber-defence environments, in contrast, are typically dynamic, with significantly larger and more heterogeneous sets of nodes, states, actions, and features. This raises the question of whether ToMnet can scale to environments suffering the ‘curse-of-dimensionality’.

The domain-specific nature of our experiments also necessitates bespoke evaluation methodologies. For example, in the original paper, predicted successor representations were qualitatively evaluated against agents’ true behaviours.³³ The key strength of this approach is that it enables an intuitive understanding of precitive idiosyncrasies and the extent to which the model can approximate the nuances of human-like reasoning (in keeping with the exploratory nature of the study). However, our aims here are slightly different in that we intend to quantitatively evaluate ToMnet’s performance across various experimental settings in order to determine its viability for real-world cyber-defence applications. As such, there is a significant incentive for developing a quantitative metric for evaluating predicted successor representations in a cybersecurity context (i.e. ones that correspond to computer network topologies).

Considering the depth of insight provided by successor representations as compared to ToMnet’s other two outputs further highlights this incentive. In contrast to next-step action probabilities and final target predictions, which constitute single-step

predictions, successor representations capture the broader spatio-temporal aspects of an agent’s policy in a single predictive map. In the context of cyber-defence, this could represent, for instance, an attacker’s predicted trajectory through a network, unlocking several unique opportunities including subgoal discovery and the subsequent allocation of defensive resources to threatened intermediate services. The performance and effectiveness of such a capability would clearly demand rigorous and principled assessment.

We identify several requirements such a metric would have to fulfil to be of use to network administrators. It must, first of all, be interpretable. This rules out unbounded metrics that lack a fixed reference scale (e.g. the Mean-Squared Error), as human operators may not always possess the context or domain-expertise required to interpret their meaning, especially when comparing values for networks of different scales. Second, the metric must be able to handle sparse inputs. This is because cyber attackers can follow relatively narrow paths towards targets in the network, resulting in limited network exploration and, consequently, (near-)zero expected occupancies that can disproportionately affect the outputs of zero-sensitive metrics (e.g. the Kullback-Leibler Divergence, and, by extension, the Jensen-Shannon Divergence). Finally, given that computer networks, and therefore any corresponding successor representations, represent geometric spaces in which nodes adhere to specific topological configurations, having a metric that explicitly accounts for this is imperative. For instance, a predicted path that lies closer to a target path in terms of proximity in network space than another prediction is intuitively more useful, even if it is equally erroneous in terms of misplaced probability mass.

In summary, what is needed to adapt the MToM framework to the cybersecurity domain is:

1. A flexible ToMnet architecture that can maintain robust predictions over high-dimensional networks and feature spaces;
2. An interpretable, stable metric for quantitatively evaluating predicted successor representations that respects the spatial relationships between their elements.

Contributions Summary

In this paper, we present a novel graph neural network (GNN)-based ToMnet architecture along with a specialised metric for measuring similarity between graphical successor representations. Together these contributions address the above limitations to bring MToM a step closer to being applicable to the domain of cyber-defence.

The literature on using machine learning for autonomous cyber-defence has identified a graph representation as a natural fit for the domain.¹³ Therefore, we propose Graph-In, Graph-Out ToMnet (GIGO-ToM) for cyber-defence, for which both feature extractors and output layers are implemented using GNNs. GIGO-ToM provides a flexible formulation where the output dimensionality matches that of the presented graph-based observation, allowing it to parse inputs of variable dimensions. We train GIGO-ToM to reason over the likely behaviours and goals of cyber-attacking agents within a reputable, abstract cyber-defence environment,³⁵ and show how these can be accurately predicted across a range of network topologies.

For our main contribution, we present a novel extension of the Wasserstein Distance (WD)³⁶ for measuring similarity between graph-based probability distributions. The baseline Wasserstein Distance computes the minimum ‘cost’ of translating one distribution into another by optimally reallocating probability mass according to a specified cost matrix D . Therefore, by setting D to the matrix of pairwise shortest path lengths between nodes in the input graph, a value is generated representing the minimum amount of work required to equate input distributions by moving probability mass around the network. This value is upper-bounded by the network’s diameter. Therefore, dividing it by the network’s diameter yields a network-agnostic, unit-bounded metric representing a fraction of the worst-case scenario. We call this metric the *Network Transport Distance* (NTD), and use it to evaluate predicted successor representations throughout our experiments.

For additional customizability, we furnish the NTD with a weighting function that linearly combines an arbitrary set of node feature vectors with a corresponding set of user-specified parameters (including a set of weighting coefficients for each selected node feature and a floor parameter for the min-max scaling function that determines the overall influence of the weighting on the final score). This yields a composite weights vector that is subsequently used to rescale input distributions, thereby emphasizing areas of the network highlighted by the selected node features whilst preserving the metric’s key theoretical properties. We demonstrate how this method can be used to gain a strategically richer understanding of how prediction errors are distributed around the graph by analyzing discrepancies between the metric values produced by different weighting configurations.

Related Work

GNNs are a class of neural networks that can be applied to graphs— data structures that systematically model relationships between entities¹³—making them highly effective for applications where data points have explicit relationships. GNNs therefore represent a promising candidate for any MToM problem that can be represented graphically. For example, Wang et al.³⁷ used a GNN-based ToMnet to build socially intelligent agents. Shu et al.³⁸ evaluated if agents designed to reason about other

agents, including a GNN-based ToMnet approach, can learn or hold the core psychology principles that drive human reasoning. GNN-based ToM approaches have also been applied to learn the relationships between players within multi-player games such as Press Diplomacy,³⁹ and emergent adversarial communication.⁴⁰

As mentioned above, ToM has also previously been applied to cyber-security. Cheng et al.⁴¹ introduced a ToM-based stochastic game-theoretic approach to reason about the beliefs and behaviors of attackers by using a Bayesian attack graph to model multi-step attack scenarios. Malloy and Gonzalez⁴² present a novel model of human decision-making inspired by the cognitive faculties of instance-based learning theory, ToM, and transfer of learning, finding from experimentation in a simple Stackelberg Security Game⁴³ that theory of mind can improve transfer of learning in cognitive models. However, we have found no papers reviewing the general utility of ToMnet in the domain of cyber-defence. As such, to our knowledge, our broad focus on the development of network-agnostic, scalable GNN-based ToMnet formulations is unique.

Although the application of our evaluation metric is specialized, there do exist promising candidates off-the-shelf. In statistics, the Wasserstein Distance (WD), also known as the Earth Mover’s Distance, is a measure used primarily for comparing probability distributions represented over a metric space.³⁶ It quantifies the minimal cost required to transform one distribution into another by considering the work needed to relocate the distribution mass in the most efficient manner possible.

Owing to its flexibility (especially with respect to being able to define a custom metric space), The WD has been extensively applied with various extensions to adapt it to different problem domains. For instance, Wang et al.⁴⁴ introduced a normalized variant tailored for tiny object detection, which models bounding boxes as Gaussian distributions to improve the robustness of detection performance against minor localization errors. The normalization is key, as it converts a boundless distance metric to a standardised similarity measure (i.e. between 0 and 1, like IoU).

For graphical applications, Kim et al.⁴⁵ pioneered a two-step ARG matching algorithm that improved the robustness and performance of graph matching through nested WD computations, while Noels et al.⁴⁶ crafted an WD-based graph distance metric tailored for analyzing financial statements, enhancing tools for fraud detection and company benchmarking. Other closely related existing metrics in this domain include the Geometric Graph Distance (GGD),⁴⁷ which quantifies the minimum cost required to transform one geometric graph into another by modifying node positions and edge connections within a Euclidean space, and the Graph Mover’s Distance (GMD),⁴⁸ which offers a more computationally tractable extension of the GGD, simplifying the transformation process using the EMD and making it more feasible for large-scale applications.

While both GGD and GMD consider graph topology, their focus is primarily on geometric and structural transformations. Our method, the Network Transport Distance (NTD), distinguishes itself by focusing specifically on how information or resources flow through a network, directly incorporating paths and distances into the WD computation and subsequent normalization in a way that, to the best of our knowledge, has not been explored previously. Furthermore, in contrast to some of the more domain-specific WD extensions mentioned, the NTD offers a versatile metric that can be applied to any problem involving the comparison of graph-based probability distributions, and enriched with \mathcal{W} wherever one wishes to explore the relationship between node features and prediction errors.

In summary, whilst primarily geared at developing the field of cyber-defence, our methodology includes novel contributions to both the suite of existing ToMnet architectures as well as the network analysis toolbox that can be utilized beyond this domain.

Preliminaries

The Hot-Desking User Problem

In this report we evaluate the ability of ToMnet architectures to characterize agents situated within cyber-defence games. Our focus is on scenarios featuring two agents, a Red cyber-attacking agent and a Blue cyber-defence agent. For our specific use-case, we consider a scenario that we term the *hot-desking user problem*. The problem consists of a computer network \mathcal{G} with a static topology (i.e. the set of nodes $v \in \mathcal{V}$ and the edges $e \in \mathcal{E}$ connecting them remains the same for a given topology). In each episode we have changes with respect to:

- (i) the vulnerabilities of each node v (e.g. due to different services running);
- (ii) the location of the high-value nodes, which are determined based on the location of a set of users $u \in \mathcal{U}$.

We assume that the number of users $|\mathcal{U}|$ remains consistent across episodes. However, due to hot-desking, the locations of the users change after each episode. Therefore, in each episode we have a different set of tuples (u_i, v_j) , where i represents each user and j represents each high-value node. In our scenarios, a Red cyber-attacking agent is attempting to evade a Blue cyber-defence agent in order to reach the high-value nodes, over which Red has a preference. This preference could be due to, say, a certain user working on a classified topic of interest to the Red agent. Based on the log-files from previous attacks—the past trajectories $\tau^{(obs)}$ —our ToMnet architectures are tasked with predicting:

- (i) the exact node that Red is targeting: \hat{v} ;
- (ii) the attack trajectory followed by Red on the path towards \hat{v} (or, the *successor representation*): \hat{SR} .³⁴

An illustration of the hot-desking user problem is provided in Figure 2.

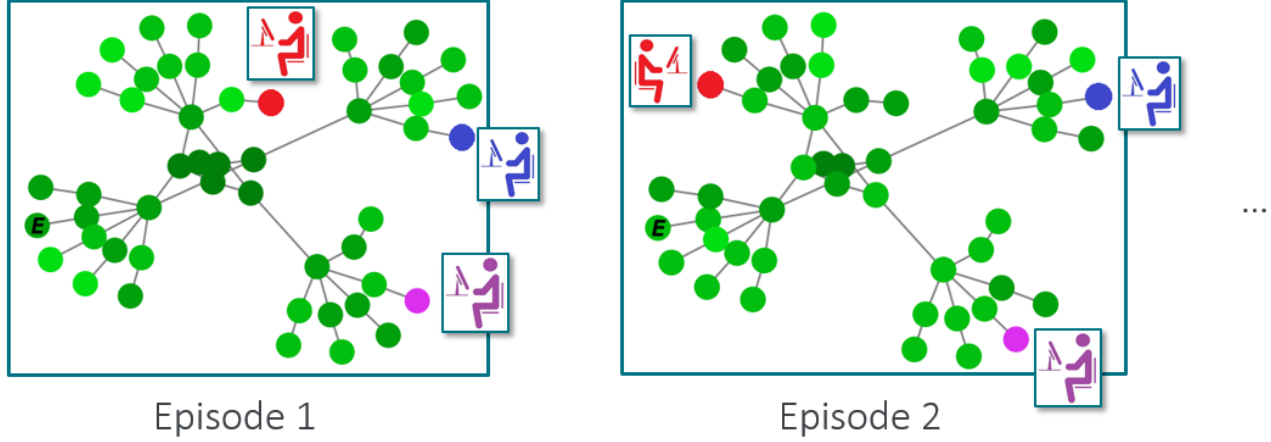


Figure 2. An illustration of the hot-desking user problem for cyber-defence. The node labelled ‘E’ denotes the entry node for the Red cyber-attacking agent. Blue, red, and pink nodes represent current user locations. The color weighting of the remaining nodes represents their vulnerability to attacks, with darker shades representing higher vulnerability scores. In each episode the users may relocate to a different desk. Therefore, ToM solutions are required that can generalize across settings with respect to different user locations in each episode.

Machine Theory of Mind

The hot-desking user problem can be thought of as a *partially observable Markov game* (POMG). Unlike in (fully observable) Markov Games, also known as stochastic games,⁴⁹ the full state of the environment is hidden from all players within a POMG. This closely reflects the conditions of real-world cyber-defence scenarios, wherein cyber-attacking agents will typically be unable to observe the complete layout of the network and all the services running on each node. Similarly, cyber-defence agents will rarely be able to observe the full state of the network given the computational overhead that would be required to relay all of this information at every time-step.

Formally, our task for ToMnet is to make predictions over a family of POMGs:

$$\mathcal{M} = \bigcup_j \mathcal{M}_j. \quad (1)$$

Each \mathcal{M}_j represents a POMG and is defined as a tuple $(\mathcal{N}, \mathcal{S}, \mathcal{A}, \mathcal{P}, \mathcal{R}, \gamma, \Omega, \mathcal{O})$, where:

- \mathcal{N} is a set of agents.
- \mathcal{S} is a finite state space.
- \mathcal{A} is a joint action space $(\mathcal{A}_1 \times \dots \times \mathcal{A}_n)$ for each state $s \in \mathcal{S}$, with \mathcal{A}_i being the number of actions available to player i within a state s .
- \mathcal{P} is a state transition function: $\mathcal{S}_t \times \mathcal{A}_1 \times \dots \times \mathcal{A}_n \times \mathcal{S}_{t+1} \rightarrow [0, 1]$, returning the probability of transitioning from a state s_t to s_{t+1} given an action profile $a_1 \times \dots \times a_n$. Here, each action a_i belongs to the set of actions available to agent i : $a_i \in \mathcal{A}_i$.
- Ω is a set of observations.
- \mathcal{O} is an observation probability function defined as $\mathcal{O} : \mathcal{S} \times \mathcal{A}_1 \times \dots \times \mathcal{A}_n \times \Omega \rightarrow [0, 1]$, where, for agent i , a distribution over observations o that may occur in state s is returned, given an action profile $a_1 \times \dots \times a_n$.

- \mathcal{R} is a reward function: $\mathcal{R} : \mathcal{S}_t \times \mathcal{A}_1 \times \dots \times \mathcal{A}_n \times \mathcal{S}_{t+1} \rightarrow \mathbb{R}$, that returns a reward r .
- γ is a discount rate defining the agent's preference over immediate and future rewards within the decision-making process. Larger values for γ define a preference for long-term rewards.

For the environments considered in this paper we also allow *terminal states* at which the game ends.

The ToMnet architectures in this report are tasked with making predictions over a range of network topologies with different node settings. We also assume we have a family of agents:

$$\mathcal{N} = \bigcup_i \mathcal{N}_i, \quad (2)$$

with each \mathcal{N}_i representing a specific agent. Each agent is implemented by a policy π_i , defining the agent's behaviour. We shall use the same problem formulation as Rabinowitz et al.,³³ and associate the reward functions, discount factors, and conditional observation functions with the agents \mathcal{N} rather than the POMGs \mathcal{M} . Therefore, for each POMG \mathcal{M}_j we have a tuple $(\mathcal{S}, \mathcal{A}, \mathcal{P})$, and for each agent \mathcal{N}_i we have $(\Omega, \mathcal{O}, \mathcal{R}, \gamma, \pi)$.

For our ToMnet approaches we assume access to (potentially partial and/or noisy) observations of the agents, via a state-observation function $\omega^{(obs)}(\cdot) : \mathcal{S} \rightarrow \Omega^{(obs)}$ and an action-observation function $\alpha^{(obs)}(\cdot) : \mathcal{A} \rightarrow \mathcal{A}^{(obs)}$. For each agent \mathcal{N}_i in a POMG \mathcal{M}_j , we assume our observer can see a set of *trajectories*, defined as sets of state-action pairs: $\tau_{i,j}^{(obs)} = \{(s_t^{(obs)}, a_t^{(obs)})\}_{t=0}^T$, where $a_t^{(obs)} = \alpha^{(obs)}(a_t)$ and $s_t^{(obs)} = \omega^{(obs)}(s_t)$. Critically, setting $\omega^{(obs)}(\cdot)$ and $\alpha^{(obs)}(\cdot)$ to the identity function grants the observer full observability of the POMG state along with all overt actions taken by the agents while keeping their parameters, reward functions, policies, and identifiers concealed. This formulation closely approximates the depth of knowledge that would be available to a human observer in a real-world scenario.

Our approaches are based on the original ToMnet³³ (Figure 1). To recap, ToMnet consists of three modules, a *character net*, *mental net*, and a *prediction net*. The character net f_θ yields a character embedding $e_{c,i}$ upon parsing N_{past} past trajectories $\tau^{(obs)}$ for an agent \mathcal{N}_i , $\{\tau_{i,j}^{(obs)}\}_{j=1}^{N_{past}}$. As in the original paper, we shall process the trajectories independently using a recurrent neural network, implemented with a Long Short-Term Memory (LSTM) network⁵⁰ f_θ , and sum the outputs:

$$e_{c,i} = \sum_{j=1}^{N_{past}} f_\theta(\tau_{i,j}^{(obs)}). \quad (3)$$

The mental (state) net g_ϑ meanwhile processes the current trajectory up to time-step $t-1$ along with $e_{c,i}$:

$$e_{m,i} = g_\vartheta([\tau_{i,j}^{(obs)}]_{0:t-1}, e_{c,i}). \quad (4)$$

The prediction net subsequently uses both embeddings and the current state s to predict the agent's behaviour, including: items that the agent will consume \hat{c} (e.g. a high-value target node that a Red cyber-attacking agent wants to gain access to), and; the successor representations $\hat{S}\mathcal{R}$ ³⁴ (i.e. the Red cyber-attacking agent's attack trajectory towards \hat{c}).

Graph-Based YAWNING-TITAN

We simulate the hot-desking user problem in an adapted version of the YAWNING-TITAN cyber-defence environment. This is an open-source framework originally built to train cyber-defence agents to defend arbitrary network topologies.³⁵ Each machine in the network has parameters that affect the extent to which they can be impacted by Blue and Red agent behaviour. These include vulnerability scores that determinine how easy it is for a node to be compromised. Our extension includes graph observations for both graph-based ToMnet architectures and graph-based Blue cyber-defence agents. Critically, these observations will contain different combinations of node features depending on the observability of the agent being trained. For example, a ToMnet architecture with full observability will receive observations containing relevant labels in the node features (such as high-value nodes), whereas these will be absent for cyber agents with partial observability.

Networks

Custom network configurations are required to systematically evaluate GIGO-ToM across different network sizes and topologies. Therefore, we have added a new network generator, *TreeNetwork*, which resembles a tree network with layers for the core, edge, aggregation, access, and subnet nodes and can be easily visualized (Figure 3a – 3e).

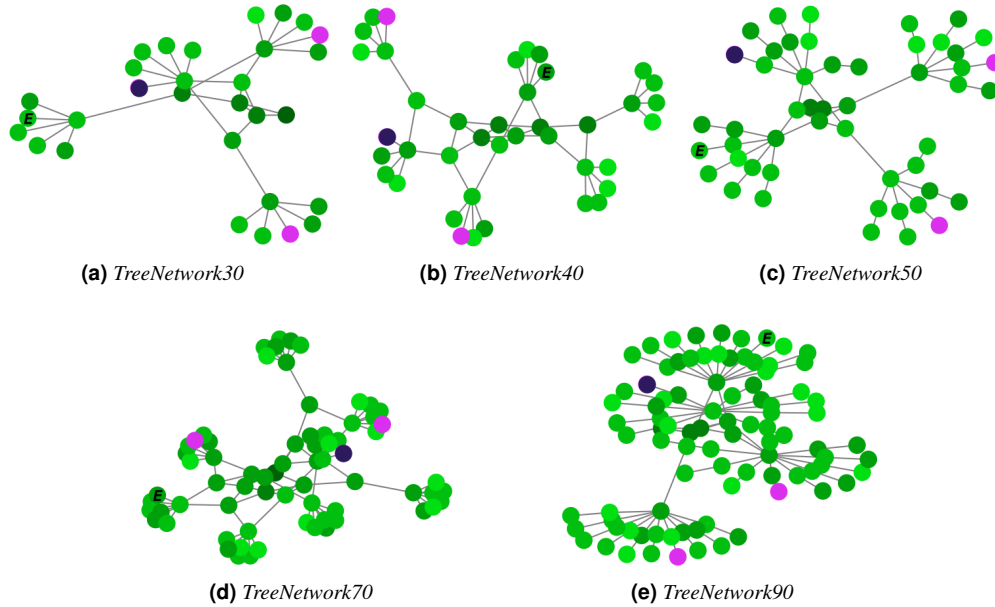


Figure 3. Example visualizations of the five custom YAWNING-TITAN *TreeNetwork* topologies used in our experiments. For each topology, the node labelled ‘E’ denotes the entry node for the Red agent. Pink nodes represent high-value nodes, with the dark blue node being the one ultimately targeted by a hypothetical Red agent.

Agents

To enable an extensive evaluation, we have additionally implemented a range of rule-based Blue cyber-defence and Red cyber-attacking agents to our graph-based YAWNING-TITAN framework. While our agent repertoire is broad (a comprehensive overview is provided in [Appendix B](#)), our evaluations suggest that only two were useful in generating interesting behaviour for our experiments:

- **BlueMSN-D:** A deterministic agent that removes the infection from a compromised node (via the `MakeSafeNode` action) if it is within three hops of any high-value node, otherwise scans the network for hidden compromised nodes (via the `Scan` action). The strategy here is to protect only the core of the network. Blue is content with sacrificing the edge nodes to do so.
- **RedHVTPreferenceSP:** Selects a particular high-value node to target at the start of the episode and takes the shortest path towards it. The selected high-value node is based on a high-value node preference vector sampled from a Dirichlet distribution $\pi \sim Dir(\alpha)$ as well as the shortest distance between high-value nodes and entry nodes. Therefore, following,³³ we can define an agent *species* as corresponding to a particular value of α , with its members being the various parameterizations of π sampled from $Dir(\alpha)$. This attacker has prior knowledge of the network structure, for example, due to insider information, meaning it can calculate path lengths and take the most direct path to its chosen high-value node.

The motivation for using rules-based agents over incentive-based learning approaches lies in their interpretability. For example, we can define a preference over high-value targets and actions, as well as reducing the amount of time required for gathering training and evaluation data for our experiments. However, we note that, in principle, ToM approaches are agnostic with respect to how the Blue and Red agent policies are obtained, and can therefore also be applied to learnt policies.

Methods

As previously mentioned, one of the challenges when applying ToMnet to cyber-defence scenarios is that the number of nodes in the network is not necessarily known in advance. This has implications for both the input and output layers of ToMnet. For our current problem domain, neural network layers are required that can handle variable sized inputs.¹³ There are now a number of solutions to this problem. For example, using transformers⁵¹ or applying one dimensional convolutions with global max pooling. However, we consider that a natural format for representing a cyber-defence scenario on a computer

network is to use a graph-based representation. Therefore, we have implemented a graph-based observation wrapper for the YAWNING-TITAN cyber-defence environment.³⁵

In order to process the graph-based state-observations, we introduce two graph-based ToMnet architectures for making predictions for cyber-defence scenarios: i.) a *Graph-in, Graph-Out* ToMnet (GIGO-ToM), where both input and output layers are implemented using graph neural networks, and; ii.) a *Graph-in, Dense-Out* ToMnet (GIDO-ToM), where inputs are processed by graph neural network layers and outputs are generated by subsequent dense neural network layers, more closely approximating the original ToMnet formulation to benchmark our evaluations.

Following evidence from the literature that graph attention layers⁵² are less susceptible to the ‘over-squashing’ of node features than other popular candidates,⁵³ these are used for node feature extraction end-to-end in both models.

For both GIDO-ToM and GIGO-ToM, the character network processes sequences of graph-based state-observations $\tau_{ij}^{(obs)}$ called *trajectories*. From the trajectories, each graph is fed through two GATv2 layers combined with dropout layers ($p = 0.5$). The outputs from each layer are then subjected to both *global max* and *global average* pooling, with the outputs subsequently being concatenated. The resulting features for each graph are then sequentially fed into an LSTM that generates a character embedding $e_{c,i}$ for each individual trajectory i . The outputs from the LSTM are then summed as described in Equation 3, resulting in the character embeddings.

GATv2 layers are also used to extract features for the mental embedding e_m and for extracting features that are fed into the prediction network. While our GIDO-ToM and GIGO-ToM architectures share the same feature extraction for the character and mental network components (Figure 4), they differ with respect to the implementation of the prediction network and output layers. Below we describe these differences in more detail.

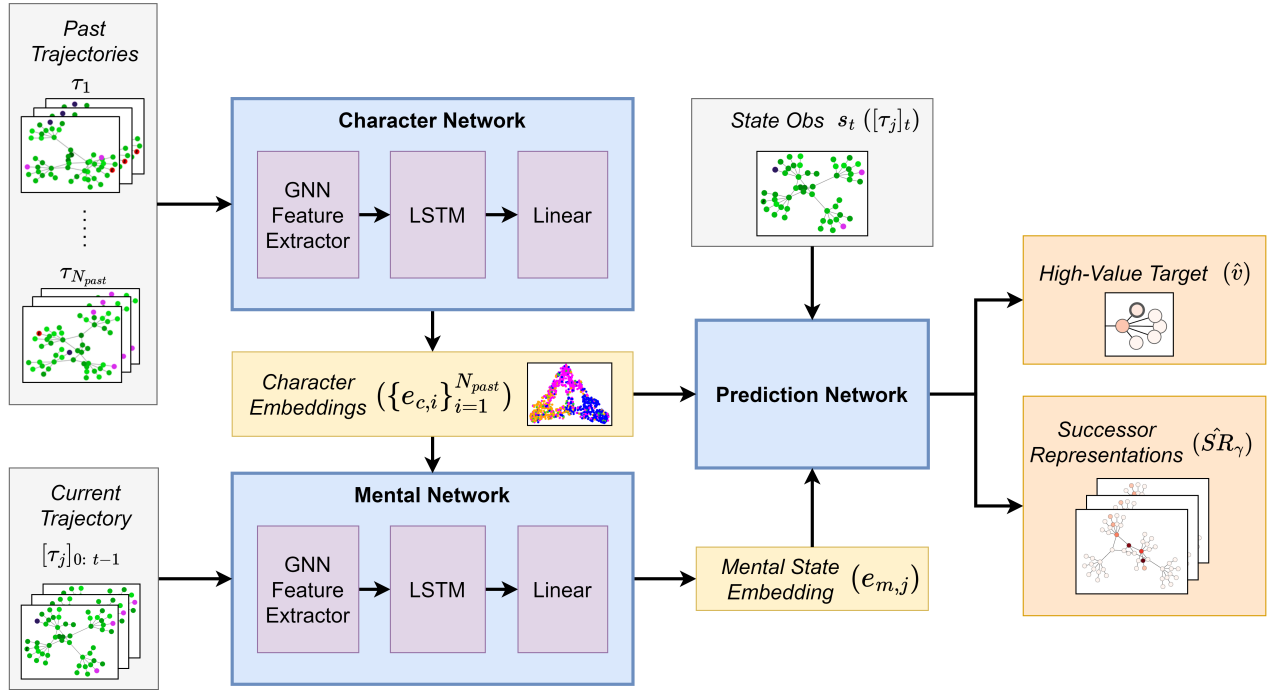


Figure 4. Shared architectural components of GIGO-ToM and GIDO-ToM.

Graph-In, Dense-Out ToMnet

Similar to the original ToMnet, for GIDO-ToM, the character and mental embeddings, e_c and e_m , along with the current state observation $s_t^{(obs)}$, are fed into a prediction network. The prediction network first uses GATv2 layers to extract features from the $s_t^{(obs)}$, before concatenating the resulting embedding with e_c and e_m . The concatenated embeddings are subsequently fed through a number of dense layers before being passed to the individual output layers: i.) a prediction over the consumable \hat{c} , representing a prediction over which high-value target the Red cyber-attacking agent wants to reach, and; ii.) the successor representation $\hat{S}R$.

As for the original ToMnet,³³ past and current trajectories will be sampled for each agent \mathcal{N}_i to obtain predictions. Each prediction provides a contribution to the loss, where the average of the respective losses across each of the agents in the minibatch is used to give an equal weighting to each loss component.

For the high-value target losses we use the negative log-likelihood loss for each high-value node n :

$$\mathcal{L}_{hvt,i} = \sum_n -\log p_{c_n}(c_n | s_t^{(obs)}, e_{c,i}, e_{m,i}). \quad (5)$$

Finally, we apply a soft label cross entropy loss \mathcal{L}_{sr} to optimise the network with respect to predicting the successor representations:

$$\mathcal{L}_{sr,i} = \sum_{\gamma} \sum_s -SR_{\gamma}(s) \log \hat{SR}_{\gamma}(s), \quad (6)$$

where s is a state-observation $s_t^{(obs)}$, and for each discount factor γ , we obtain the ground truth $SR_{\gamma}(s_t^{(obs)})$ through an empirical normalised discounted rollout, from the time-step t where the state-observation was gathered, onwards. The combined loss is:

$$\mathcal{L}_{total} = \mathcal{L}_{hvt} + \mathcal{L}_{sr}. \quad (7)$$

The disadvantage of the above approach is that we fix the max number of nodes per output layer through examining the maximum number of nodes present within the dataset. GIDO-ToM is depicted in [Figure 5](#).

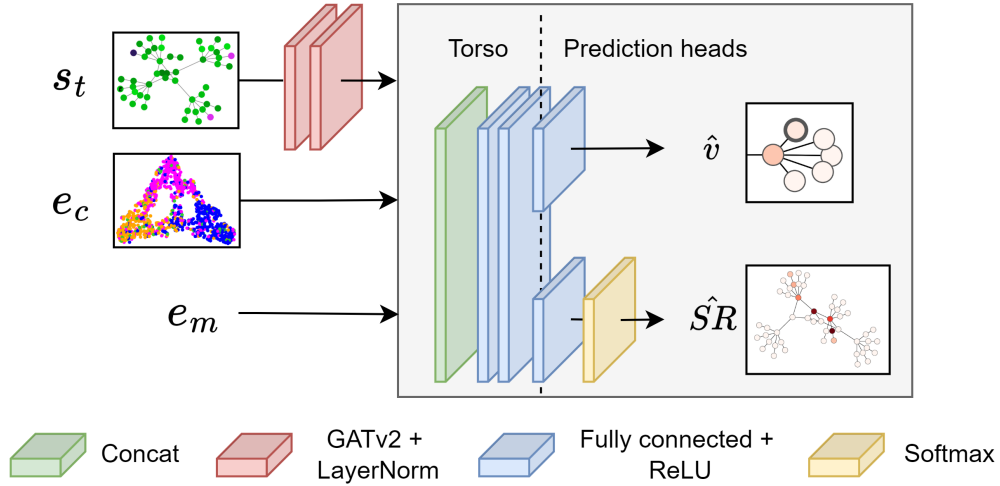


Figure 5. The Graph-In Dense-Out (GIDO-ToM) prediction network architecture. Output layers are implemented using dense neural network layers, necessitating a fixed number of output nodes.

Graph-In, Graph-Out ToMnet

In contrast to GIDO-ToM, for GIGO-ToM we concatenate the character and mental network embeddings directly with each node's features within the current state-observation. These are then directly processed by GNN layers before each of the prediction tasks. As mentioned above, for GIGO-ToM, the high-value target and successor representation predictions are generated using GATv2 layers. The final GNN layer of each respective prediction branch represents the outputs, i.e. a weighted binary cross entropy loss is applied to each node for optimizing high value node predictions:

$$\mathcal{L}_{hvt,i} = \sum_n -w_n (y_n \log x_n + (1 - y_n) \log(1 - x_n)). \quad (8)$$

where x represents the ground-truth with respect to whether or not a node was consumed at the end of the episode, and y is the prediction. The positive weighting term w_n is added to account for the positive / negative node imbalance *within* an individual sample, i.e., only one out of n nodes will (potentially) be consumed at the end of a given episode.

For the successor representation loss \mathcal{L}_{sr} we stick with the soft label cross entropy loss from Rabinowitz et al.³³ For batches where the number of nodes on each graph is inconsistent we zero-pad while computing the loss. We note this formulation only requires padding when computing the loss, meaning at inference time there are no superfluous nodes (in contrast to when specifying the max number of nodes for GIDO-ToM).

The approach outlined above provides a flexible formulation regarding the number of nodes in the network, with the output topology reflecting that of the current state-observation fed into the prediction net. GIGO-ToM is depicted in [Figure 6](#).

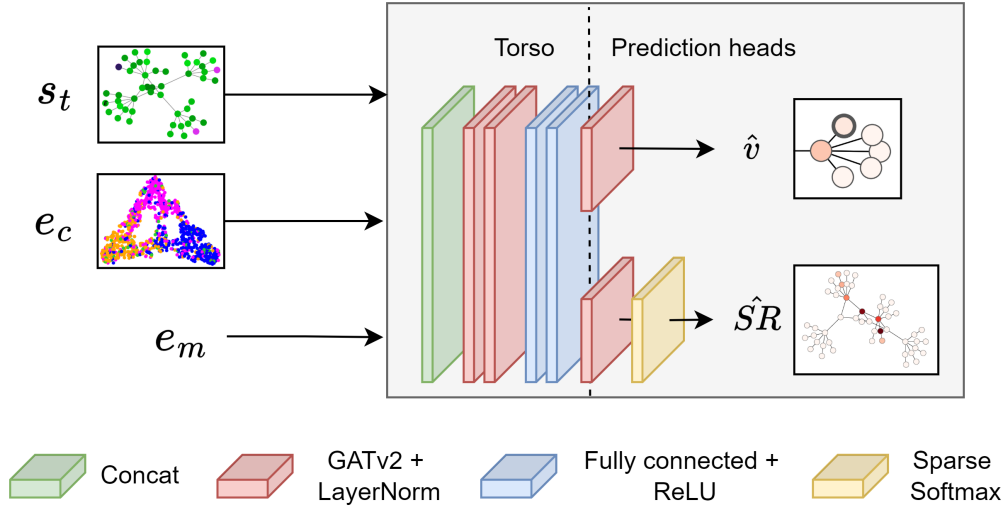


Figure 6. The Graph-In Graph-Out (GIGO-ToM) prediction network architecture. Both feature extractors and output layers are implemented using graph neural network layers, providing a flexible formulation that can parse inputs and generate outputs of variable dimensions.

The Network Transport Distance

Of the ToMnet outputs, one that is particularly salient for the depth of insights on offer to network administrators is the *successor representation*. In the field of reinforcement learning, these are formally defined as matrices representing the expected discounted future state occupancy following a given policy from a given state.³⁴ In the case of our hot-desking user problem, these represent predictions regarding the likely attack trajectory that Red will take against the current Blue cyber-defence agent. Predicted successor representations are often evaluated using generic metrics that assume input data points are independent and identically distributed (i.i.d.). However, in real-world cyber-defence scenarios, the nodes of a network will often have individual features granting them varying degrees of vulnerability, utility, and overall strategic importance,⁵⁴ and will adhere to a certain topological configuration. As such, a metric for evaluating predicted successor representations that captures these contextual nuances is desirable.

We achieve this with the Network Transport Distance (NTD); a unit-bounded extension of the Wasserstein Distance (WD) that includes a graph-theoretic normalization factor that represents the result as a fraction of the worst-case predictive scenario for the network in question.

The metric can be calculated as follows:

$$\text{NTD}(P, Q, D) = \frac{1}{\max(D)} \inf_{\mu \in M(P, Q)} \left[\int_{X \times X} d(i, j) d\mu(i, j) \right], \quad (9)$$

where:

- P and Q represent probability distributions in which individual elements correspond to nodes in a graph \mathcal{G} .
- D represents the matrix of pairwise shortest path lengths between nodes in \mathcal{G} . The diameter of \mathcal{G} is therefore $\max(D)$.
- $\inf_{\mu \in M(P, Q)}$ denotes the infimum (i.e. the greatest lower bound) over all possible transport plans μ . $M(P, Q)$ represents the set of all joint distributions whose marginals are P and Q .
- $\int_{X \times X} d(i, j) d\mu(i, j)$ calculates the total cost of transporting mass from distribution P to distribution Q . The function $d(i, j)$ is the ground metric on the graphical space X , which here is the shortest path length between nodes i and j in G . Therefore, $d(i, j) = D[i, j]$. $d\mu(i, j)$ denotes the amount of probability mass transported from i to j .

Like the Wasserstein Distance (WD), the NTD is symmetric and non-negative:

$$\text{NTD}(P, Q, D) = 0 \Leftrightarrow P = Q. \quad (10)$$

Unlike the WD, however, the NTD has an upper bound of 1:

$$\text{NTD}(P, Q, D) = 1 \Leftrightarrow \text{WD}(P, Q, D) = \max(D). \quad (11)$$

This makes the metric network-agnostic and more interpretable to a broad audience, including non-technical stakeholders. An illustrative example of why incorporating topological information into the NTD computation can be seen in [Figure 7](#).

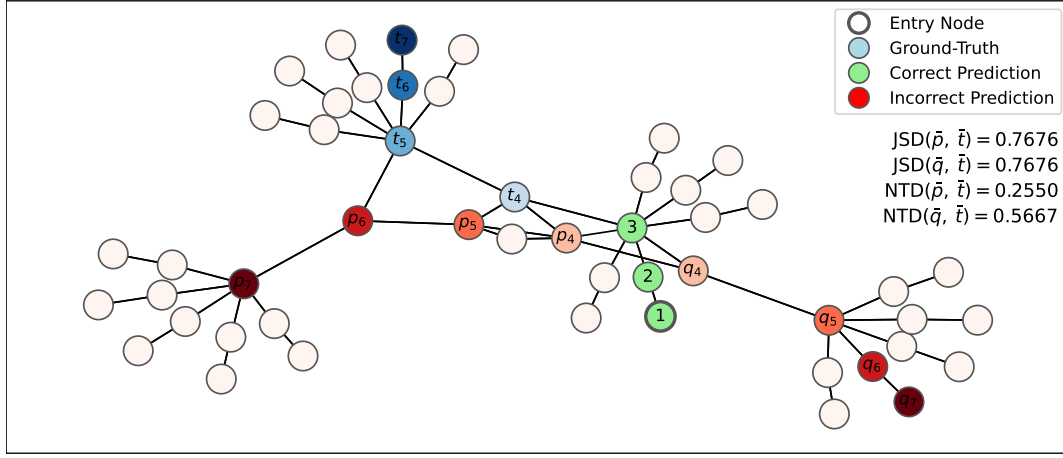


Figure 7. Hypothetical example demonstrating the importance of a topology-aware metric for evaluating graph-based successor representations. Shown is a 50-node network with three distinct attack paths overlaid: a ground-truth \bar{t} , and two predictions \bar{p} , \bar{q} . For demonstrative purposes, each path is the same length and each node along the paths is assigned the same probability score, with only their locations varying i.e. $\forall i, t_i = p_i = q_i$. Therefore, depicted is a scenario where \bar{p} , \bar{q} , and \bar{t} represent different graphical translations of the same probability distribution. Green and white nodes represent locations where probability masses for all three distributions are identical, with green denoting shared non-zero values and white shared zero values. Erroneous predictions are displayed in red, and the ground-truth in blue. Darker shades indicate higher probability scores. NTD scores for each predicted path are displayed on the right, along with Jensen-Shannon distance (JSD) scores for a non-parametric comparison. All three paths coincide for for the first 3 nodes before diverging into different regions of the network. However, the erroneous segment of \bar{p} is considerably ‘closer’ to \bar{t} than that of \bar{q} in terms of node proximity in the graph space. In contrast to the JSD, which ignores this information, the NTD explicitly leverages the graph’s structure to yield $\text{NTD}(\bar{p}, \bar{t}) < \text{NTD}(\bar{q}, \bar{t})$, as desired.

For added contextual awareness, we weight the Network Transport Distance by selecting a set of node features from a network based on their relevance to the strategic importance of node attacks, then linearly combining them with a corresponding set of user-specified coefficients using a customizable weighting function that can be applied to arbitrary metrics. We denote the end-to-end metric NTD_θ .

Mathematically, for a network \mathcal{G} with n nodes, we have:

- A set of m user-selected node feature vectors: $\mathcal{X} = [\bar{x}_0, \dots, \bar{x}_m] \in \mathbb{R}^{n \times m}$.
- A corresponding set of m user-specified feature weights: $\mathcal{C} = [c_0, \dots, c_m] \in [-1, 1]^m$.

Our weighting function $\mathcal{W} : \mathbb{R}^{n \times m} \times [-1, 1]^m \times [0, 1] \rightarrow [0, 1]^n$ linearly combines these and normalizes the result:

$$\mathcal{W}(\mathcal{X}, \mathcal{C}, f) = \left\| \sum_{i=0}^m c_i \|\bar{x}_i\|_f \right\|_f = \bar{w}, \quad (12)$$

where:

- $\bar{x}_i \in \mathcal{X}$ is a node feature vector.
- $c_i \in \mathcal{C}$ is a user-specified constant that determines the weight of \bar{x}_i in the final metric ($-1 \leq c_i \leq 1$).
- $\|\cdot\|_f$ is the min-max scaling function with ‘floor’ parameter f , with f defining the minimum weight that can be assigned to an input value ($0 \leq f \leq 1$). This forces all outputs to lie in the interval $[f, 1]$, thereby guaranteeing a minimum level of influence on the output metric:

$$\|\bar{x}\|_f = \frac{(\bar{x} - \min(\bar{x}))(1 - f)}{\max(\bar{x}) - \min(\bar{x})} + f. \quad (13)$$

The computed weighting \bar{w} is used to rescale the input probability distributions before passing them into the Network Transport Distance. Therefore, for two \mathcal{G} -based probability distributions $P, Q \in \mathbb{R}^n$, we have:

$$\text{NTD}_\theta(P, Q) = \text{NTD}(P_\theta, Q_\theta), \quad (14)$$

where $\bar{x}_\theta = \frac{\bar{w} \cdot \bar{x}}{\sum \bar{w} \cdot \bar{x}}$. It follows that \mathcal{W} is theoretically decoupled from the metric it is applied to, and can therefore be applied to any metric that is amenable to input weighting. An illustration of influence of our weighting function \mathcal{W} on output NTD_θ scores under different node feature weighting parameterizations can be seen in Figure 8.

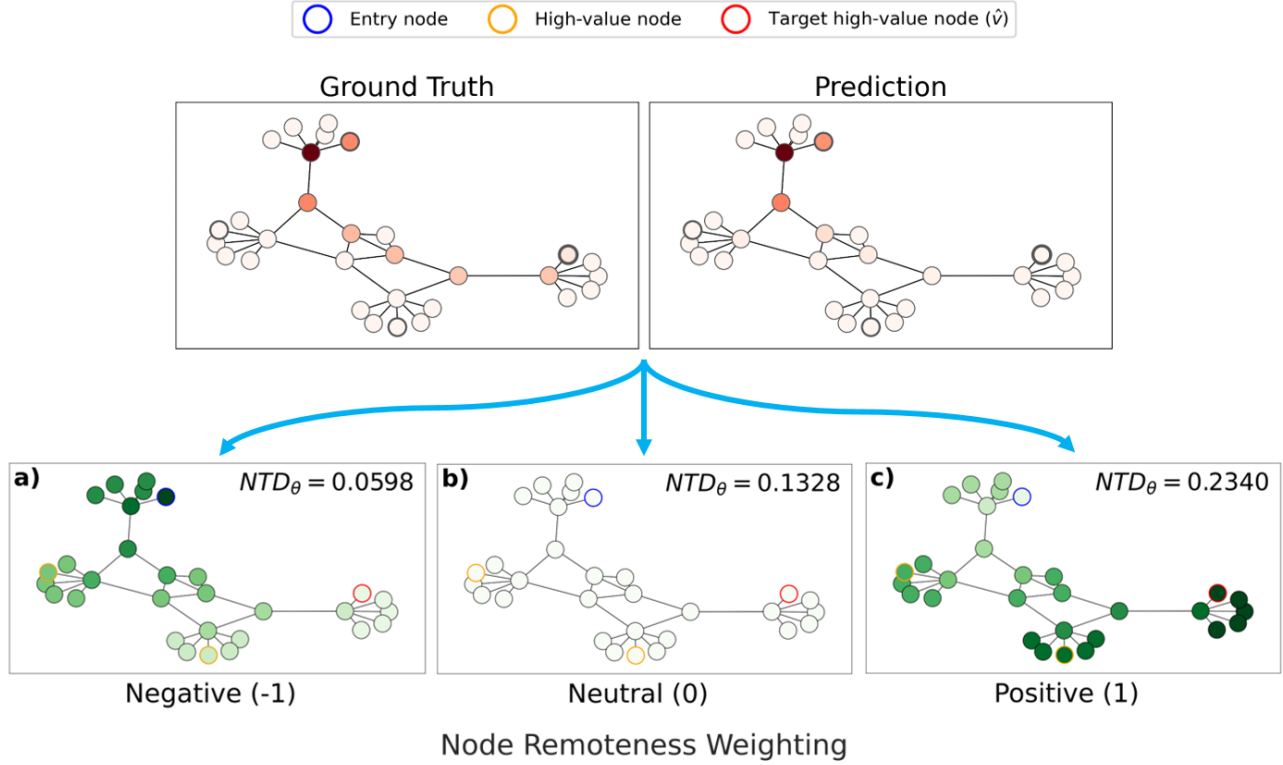


Figure 8. An illustrative example showing the influence of our weighting function \mathcal{W} on output NTD_θ scores under different node feature weighting parameterizations. The top figures display a ground truth and predicted attack path ($\hat{S}R_{true}$ and $\hat{S}R_{pred}$) for *TreeNetwork30*. Darker shades of red represent higher likelihoods of a node coming under attack. The three plots beneath (a-c) show the composite node weights $w_i \in \bar{w}$ and corresponding $\text{NTD}_\theta(\hat{S}R_{true}, \hat{S}R_{pred})$ scores under different weighting parameterizations θ . Higher weights are observable via darker shades of green. For this example, we define a single *node remoteness* feature as the shortest path length from each node to the entry node. Varying the coefficient of this feature in \mathcal{W} therefore emphasizes different parts of the network based on proximity to the entry node. Subfigure a) shows a negative weighting, which emphasizes predictions closer to the entry node; b) shows the unweighted NTD , and; c) shows a positive weighting that emphasizes predictions made further away from the entry node. The significant NTD_θ increase observed when inverting the coefficient from negative to positive reveals that $\hat{S}R_{pred}$ diverges further from $\hat{S}R_{true}$ with distance from the entry node, as can be visually confirmed in the upper figures. Such tests over NTD_θ scores averaged across larger batches of predictions yield more general insights into predictive idiosyncrasies.

Experiments

In this section, we empirically evaluate GIGO-ToM’s ability to provide actionable insights into the characters, goals, and behaviours of various cyber-attacking agents in the YAWNING-TITAN environment, using the hot-desking user problem as our task formulation.

First, we explore GIGO-ToM’s ability to characterize the policies of unseen cyber-attacking agents. We next investigate how well GIGO-ToM can predict the services targeted by these agents as well as their attack routes towards them.

Benchmarked against our dense-layered ToMnet implementation (**GIDO-ToM**), we find that, through observing past episodes of Blue and Red agents interacting with each other, GIGO-ToM can learn character embeddings that can effectively

differentiate between Red agents’ preferred outcomes, providing a means to characterize their policies. We also observe that GIGO-ToM can accurately predict Red’s preference over high-value nodes as well as their attack trajectories towards them, occasionally predicting multiple attack paths in cases where the target node is uncertain.

We now describe our experimental set-up before presenting these results in detail.

Games: For the following experiments, we consider games consisting of the `BlueMSN-D` cyber-defence agent against 1,000 possible parameterizations of the `RedHVTPreferenceSP` cyber-attacking agent. The policy of each Red agent \mathcal{N}_i is defined by a fixed vector π_i specifying its preference over high-value targets. These are sampled from a single agent species: a Dirichlet distribution with concentration parameter $\alpha = 0.01$, ensuring a diverse range of sparse policies. For generality, we build a mixed topology setting named *TreeNetworkMixed* for which each game is rolled out in one of the following network topologies: $\{\textit{TreeNetwork30}, \textit{TreeNetwork40}, \textit{TreeNetwork50}, \textit{TreeNetwork70}, \textit{TreeNetwork90}\}$. Each network has three high-value users/nodes. The location of the high-value nodes (and therefore the location of the high-value users) are randomly selected at the start of each episode. One entry node is provided to the Red cyber-attacking agent. Given that our focus is on our models’ ability to predict high-value node locations, we simplify the task with respect to the entry node by keeping it consistent across episodes. To ensure a dataset of distinctive trajectories for our experiments, high-value nodes are always situated on leaf nodes.

Data: To build a ToMnet dataset \mathcal{D}_{ToM} , a set of Blue cyber-defence agents \mathcal{N}^{blue} and a set of Red cyber-attacking agents \mathcal{N}^{red} is defined, along with a set of network configurations \mathcal{M} for them to compete within. The Cartesian product of these three sets defines a set of game configurations: $G = \mathcal{N}^{blue} \times \mathcal{N}^{red} \times \mathcal{M}$. For each game $g \in G$, n_c current episodes are generated, for each of which a distinct set of n_p further past episodes is generated so that the sets of past trajectories \mathcal{P}_i for all samples $s_i \in \mathcal{D}_{ToM}$ are mutually exclusive: $\bigcap_{i=1}^{|\mathcal{D}_{ToM}|} \mathcal{P}_i = \emptyset$. This precludes data leakage between training and validation sets inflating performance metrics. Thus, $n_{total} = n_c + n_c n_p$ episodes are generated for each game configuration $g \in G$, resulting in a database of trajectories: $\bigcup_{i=1}^{|G|} \{\tau_{ij}\}_{j=1}^{n_{total}}$. We fix $n_c = 3$ and $n_p = 8$ for all experiments, generating 27 total trajectories for each game configuration $g \in G$. Therefore, using a single `BlueMSN-D` agent against 1,000 parameterizations of the `RedHVTPreferenceSP` agent in the *TreeNetworkMixed* topology setting, a set of 3,000 current episodes are generated along with a pool of 24,000 episodes from which past and current trajectories could be extracted for building data samples. Each data sample $s_i \in \mathcal{D}_{ToM}$ consists of N_{past} past trajectories for the character network, a current trajectory up to a certain time-step t for the mental network, and the current state observation at that step $s_t^{(obs)}$ for the prediction network, along with any ground-truth variables required for evaluation (e.g. \hat{v}_{true} and $\hat{S}R_{true}$). Given that trajectories often comprise hundreds of time-steps under our settings, we sample state-observations from past trajectories at regular intervals to reduce wall-times. Five state-observations were sampled from each past trajectory for the following experiments. For each current trajectory parsed by the prediction network, we uniformly select state observations $s_t^{(obs)}$ from the first two time-steps of the current episode. We split our resulting dataset into train and validation sets according to the ratio of 75% / 25%.

Evaluation: To evaluate GIGO-ToM, we systematically examine each of its modules’ outputs over a hold-out test set of 200 unseen Red agents. Character networks are qualitatively evaluated through inspecting the embeddings e_{char} they generate for each agent. We collect the embeddings obtained from our test samples and apply t-SNE⁵⁵ in order to obtain 2D visualizations of the clusters that emerge. Prediction networks are assessed with respect to their outputs: the *high-value node* \hat{v} and the *successor representation* $\hat{S}R$. Predictions regarding targeted high-value nodes \hat{v}_{pred} are evaluated against ground-truth labels \hat{v}_{true} using weighted F1 scores and confusion matrices. Predicted successor representations $\hat{S}R_{pred}$ are evaluated against true attack paths $\hat{S}R_{true}$ using our novel distance metric: the Network Transport Distance (NTD). Given that NTD scores are computed on a per-sample basis, violin and box plots are used to visualize their distribution and summary statistics over the test set.

For each prediction task, we explore the influence of both past trajectory exposure (N_{past}) and network size on predictive performance. To investigate how $\hat{S}R$ predictions vary along spatial and temporal dimensions, we add two further evaluation conditions. First, we evaluate predicted successor representations over three discount factors $\gamma \in \{0.5, 0.95, 0.999\}$ to emphasize rewards at various points in the future ($\gamma = 0.5$ represents short-term attack trajectories, while $\gamma = 0.999$ represents long-term trajectories). Second, to better understand any irregularities in GIGO-ToM’s predictions or Red agent behaviour, we equip the NTD with our novel weighting function \mathcal{W} (Equation 12) to evaluate how accuracy varies with a custom *node remoteness* feature, defined as the length of the shortest path from a given node to the entry or targeted high-value node. This allows us to investigate if and when predicted or ground-truth attack paths venture into remote regions of the network that lie away from the path between the entry and high-value node. Weighted NTD scores are denoted NTD_θ . We fix $f = 0.1$ wherever \mathcal{W} is used throughout this report. For brevity, since we are only using a single node feature in our weighting configuration, θ will refer to the node remoteness weighting coefficient for the remainder of the report.

For each prediction task, GIGO-ToM’s performance is benchmarked against our dense-layered ToMnet variant, GIDO-ToM, on the same test set. Results for each of these are summarized below.

How well can GIDO-ToM/GIGO-ToM characterize various cyber-attacking agents?

Figure 9 shows clustered character network embeddings e_{char} for both GIGO-ToM and our benchmark model, GIDO-ToM, over hold-out test sets of unseen Red agents when given access to varying numbers of past behaviour observations ($N_{past} \in \{1, 2, 3, 4\}$). We find that while GIDO-ToM generally produces compact and clearly separable clusters, they exhibit poor intra-cluster consistency, indicating confusion between the policies of observed agents. By contrast, GIGO-ToM exhibits strong generalization capabilities across all evaluation settings, consistently producing coherent clusters that are each predominantly characterised by one of the target users. For GIDO-ToM, performance peaks at $N_{past} = 3$ and degrades again for $N_{past} = 4$, heavily overlapping the clusters for two of the users. For GIGO-ToM, by contrast, the consistency, compactness, and separateness of clusters improves with higher values of N_{past} , demonstrating its ability to learn agents' policies and generalize effectively in a few-shot manner. We accordingly hypothesize that GIGO-ToM's performance on this task would further improve with access to more past data samples.

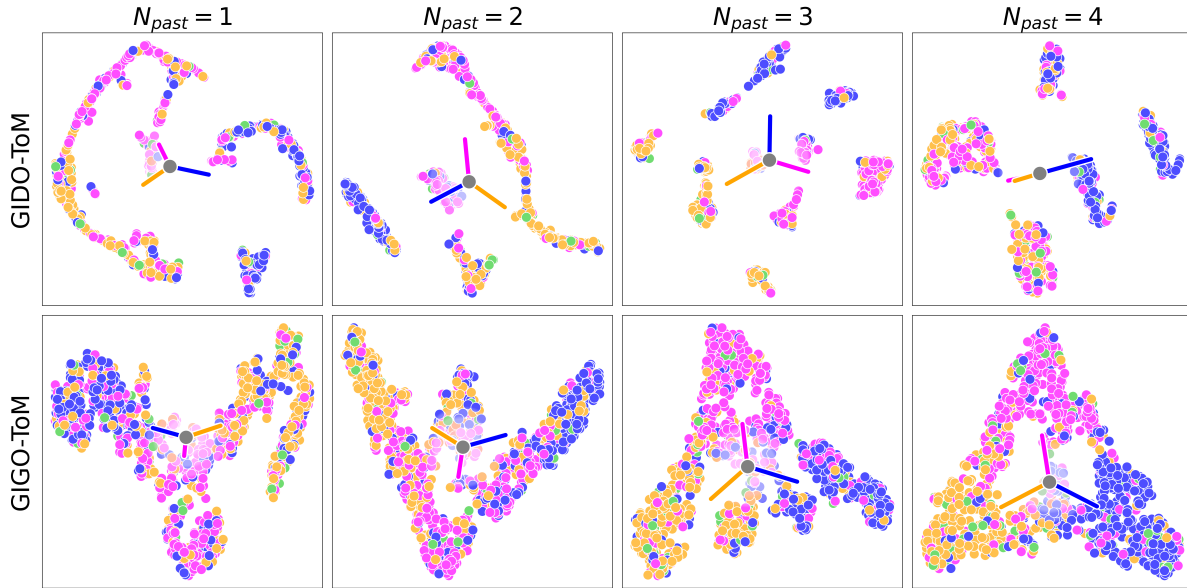


Figure 9. Test set character embeddings e_{char} for 200 unseen Red agents in the *TreeNetworkMixed* network setting. The character network was exposed to a different number of past trajectories (N_{past}) for each training run. Points are colour-coded according to the agents' policies over high-value target preference. Green points represent instances where Blue prevents Red from reaching the high-value target. Colored lines are drawn from the overall data mean to the centroid of each corresponding cluster. Longer, more divergent lines therefore indicate better class separation.

How well can GIDO-ToM/GIGO-ToM predict the services targeted by cyber-attacking agents?

For target high-value node prediction we observe that GIGO-ToM vastly outperforms the benchmark, consistently achieving weighted F1 scores exceeding an order of magnitude higher across all evaluation settings (Figure 10a). At first glance, GIGO-ToM's peak F1 score of 0.6893 ($N_{past} = 4$) suggests moderate performance. However, given the complexity of the classification problem, with 60 possible high-value nodes across the *TreeNetworkMixed* topology, this score indicates a strong predictive capability. Such is clear from Figure 11, which depicts row-normalized confusion matrices over the set of all viable high-value nodes across the *TreeNetworkMixed* topologies. The clearly distinguishable diagonals for each evaluation setting demonstrate minimal confusion between predicted target nodes.

We initially hypothesized that the class proliferation introduced by the larger networks within *TreeNetworkMixed* would degrade GIGO-ToM's performance on this task. However, by evaluating GIGO-ToM's performance on the individual network topologies within *TreeNetworkMixed*, we found no clear correlation between the number of nodes present in a network and predictive performance (Figure 10b). What did affect performance, however, was the number of branches (and therefore possible high-value target locations) present in the network. This is clear from the performance drop observed for the *TreeNetwork40* and *TreeNetwork70* topologies, which had 6 and 8 branches respectively, compared with just 4 for the others. This suggests that while GIGO-ToM's high-value target prediction capabilities are robust to high-dimensional feature spaces within the range tested in this paper, it is vulnerable to higher rates of misclassification as the number of network branches increases.

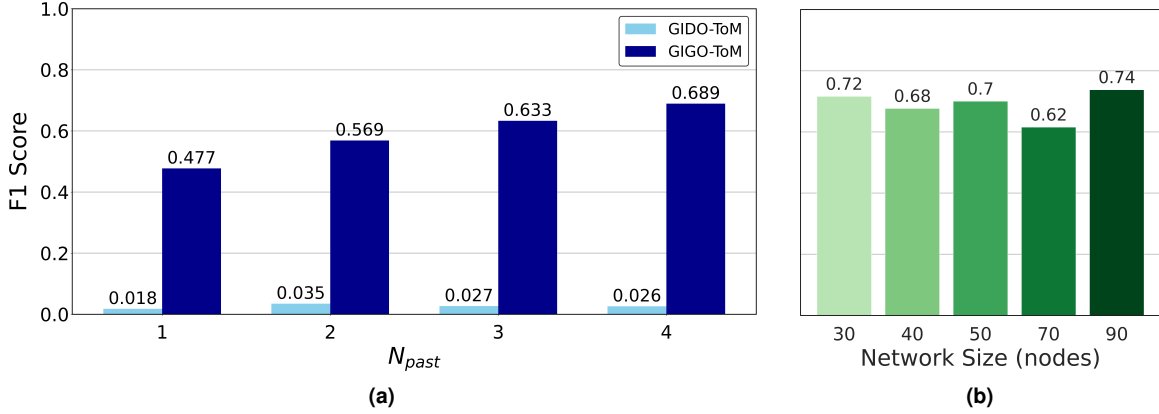


Figure 10. Figure 10a: high-value node prediction F1 scores for both GIGO-ToM and GIDO-ToM over the entire test set. Results are recorded for runs across different values of N_{past} (number of past trajectories presented to the character network). Figure 10b: Stratified F1 scores for the GIGO-ToM, $N_{past} = 4$ run, with the *TreeNetworkMixed* test set partitioned by network topology.

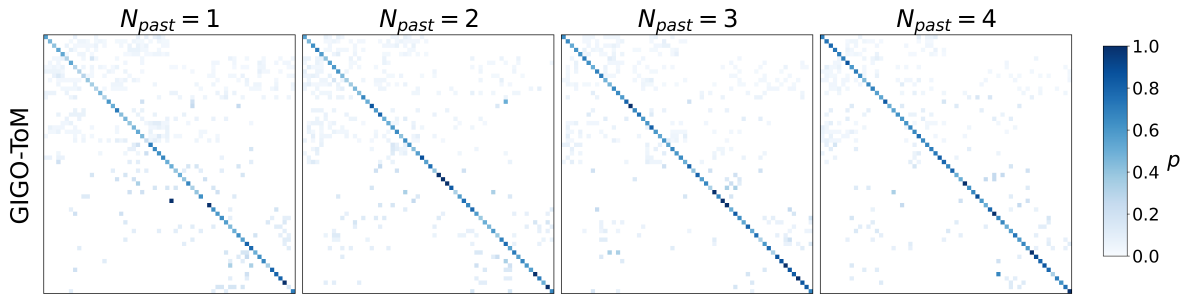


Figure 11. GIDO-ToM and GIGO-ToM confusion matrices for high-value node prediction across different values of N_{past} (number of past trajectories presented to the character network) for *TreeNetworkMixed*. For each confusion matrix, columns represent predicted target nodes and rows true target nodes. Rows are normalized so that the activation of each cell represents the proportion of predictions made for the corresponding true class (p).

How well can GIDO-ToM/GIGO-ToM predict the attack trajectories of cyber-attacking agents?

In this section, we evaluate GIGO-ToM’s ability to predict Red agent attack routes given a handful of past behaviour observations. Figure 12a displays mean baseline (unweighted) NTD scores for each N_{past} run for both GIGO-ToM and our benchmark model, GIDO-ToM. We find that GIGO-ToM consistently predicts more accurate attack trajectories than GIDO-ToM, with performance improving with increased past behaviour exposure, as for our previous experiments. Furthermore, we observe a negative correlation between GIGO-ToM’s predictive performance and network size up to the 50-node mark, after which median NTD scores plateau and variability improves slightly (Figure 12b). This finding suggests that GIGO-ToM’s successor representation prediction capabilities are robust to the larger networks in our test range despite having to predict progressively longer attack routes.

To investigate how NTD scores vary across temporal and spatial dimensions, we plot the distribution and summary statistics of NTD_{θ} scores across different values of γ (predictive time horizon) and different values of the coefficient of our *node remoteness* feature (spatial proximity to the entry or targeted high-value node in graph space) in our metric weighting function \mathcal{W} in Figure 13a.

We firstly find that performance is significantly affected by one’s choice of γ . Unsurprisingly, as before, short-term predictions ($\gamma = 0.5$) are best. This is due to the relative complexity of the YAWNING-TITAN environment, which presents a significantly more challenging learning problem when the model is forced to account for more distant future states. The same is true spatially, underscored by the significant NTD_{θ} deterioration observed when inverting the node remoteness weighting from -1 to 1 (thereby switching the emphasis from nodes closer to the Red agent’s start and end points to those further away).

This could indicate one of two possibilities: either GIGO-ToM predicts paths into remote regions of the network that the Red agent does not in fact visit, or the Red agent visits remote regions of the network that GIGO-ToM fails to predict. One

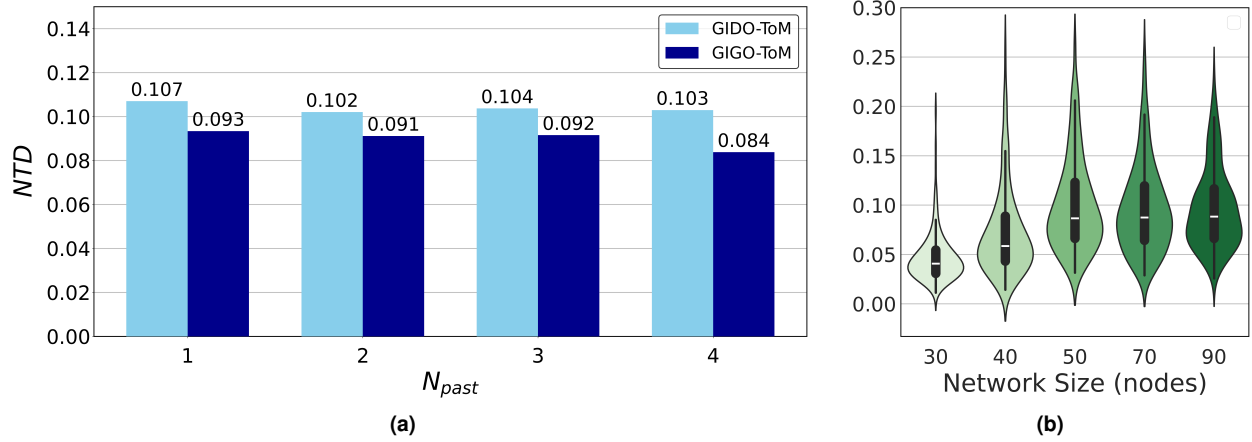


Figure 12. Figure 10a: Mean test set attack path prediction (\hat{SR}_{pred}) NTD scores for both GIGO-ToM and GIDO-ToM across different values of N_{past} (number of past trajectories presented to the character network). Figure 10b: Violin plot with overlaid box plot displaying test set NTD distribution and summary statistics for the GIGO-ToM, $N_{past} = 4$ run, stratified by network topology.

would expect to see instances of the latter if the Red agent were often pushed off-course by the Blue agent’s defensive activity. However, given the strong deterministic policies of the RedHVTPreferenceSP agents used here and our inclusion of only episodes that result in Red agent victory, the former is more likely.

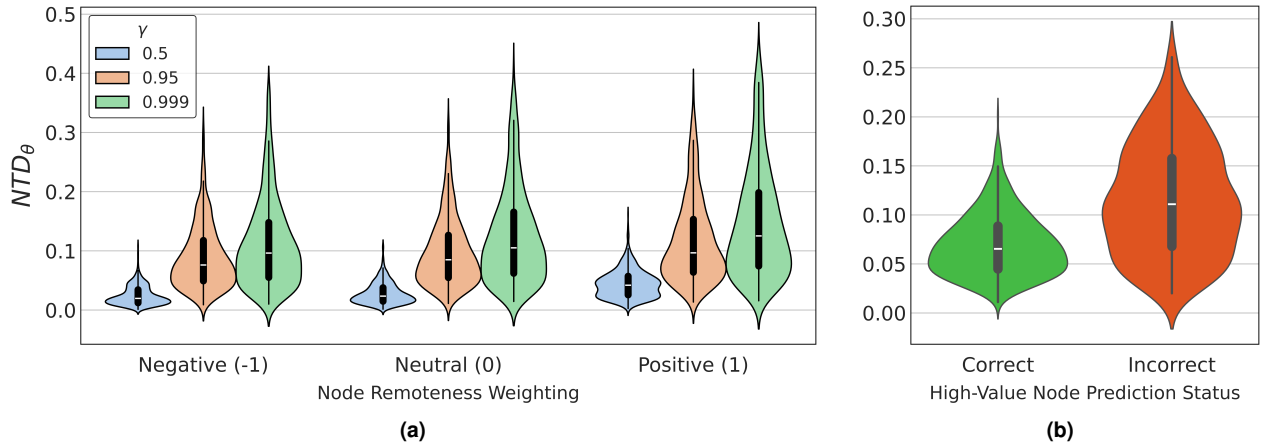


Figure 13. Figure 13a: Violin plot with overlaid box plot displaying our spatio-temporal evaluation of predicted test set attack paths: Weighted NTD_{θ} scores (GIGO-ToM, *TreeNetworkMixed*, $N_{past} = 4$), stratified by node remoteness weighting (spatial) and discount factor γ (temporal). ‘Neutral’ represents the baseline/unweighted NTD. Figure 13b: Test set NTD scores stratified by high-value node prediction status.

It should be noted, however, that, this does not necessarily entail that GIGO-ToM has simply predicted a single path to the wrong high-value node. Indeed, we do observe that GIGO-ToM’s predicted attack paths deviate significantly further from the ground-truth for samples where the incorrect target high-value node is predicted, confirming that the two outputs are largely aligned (Figure 13b). However, the interesting predictive idiosyncrasies that are possible here warrant further investigation.

As a principled next step, we inspect ground-truth vs. predicted successor representations for the test sample with the widest NTD range spanned across our node remoteness weighting settings (Figure 14). This will reveal, in simple terms, the sample for which GIGO-ToM’s predictive accuracy varies most with respect to proximity to the entry and preferred high-value node. To challenge GIGO-ToM, we consider only samples for the hardest evaluation setting within the *TreeNetworkMixed*, $N_{past} = 4$ run: *TreeNetwork90*, $\gamma = 0.999$.

Figure 14 reveals interesting behaviour. Specifically, GIGO-ToM seems to hedge it’s bets with respect to the Red agent’s

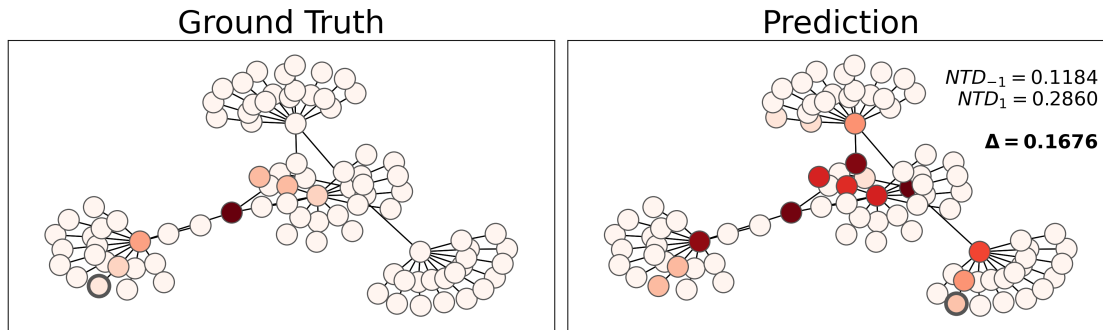


Figure 14. Evidence of GIGO-ToM’s ‘hedging’ behaviour: $\hat{S}R_{true}$ vs. $\hat{S}R_{pred}$ for the test sample with the largest discrepancy between NTD_{-1} and NTD_1 (i.e. for which the node remoteness weighting produces the greatest difference in NTD_θ score). \hat{v}_{true} and \hat{v}_{pred} are emboldened in each plot, respectively. To challenge GIGO-ToM, only samples from the most challenging evaluation sub-setting (*TreeNetwork90*, $\gamma = 0.999$) were considered here. Δ represents the absolute difference between the scores produced by the positive and negative node remoteness weighting coefficient. Displayed in this plot is the largest encountered over the test subset. We observe that GIGO-ToM predicts paths to multiple high-value nodes.

attack path, mapping routes from the entry node to every high-value node rather than committing to any single one. Indeterminate attack path predictions have little utility to network operators. Therefore, to investigate the prevalence of these non-committal predictions in our test set, we ran K-means over the predicted successor representations for these samples with $k = 4$ (3 high-value nodes plus a ‘miscellaneous’ bin). This exposed groupings based on approximate attack path, revealing that these indeterminate predictions constitute roughly 20% of the those made over the test set for *TreeNetwork90*, $\gamma = 0.999$ (Figure 15).

In summary, GIGO-ToM predicts attack paths to a level of accuracy that is certainly useful from a cyber-defence perspective. Even under our most challenging evaluation settings, distinctive attack trajectories are predicted for approximately 80% of test samples, averaging a mean NTD of 0.08. This performance is robust to high-dimensional network spaces, and improves with exposure to more past behaviour observations.

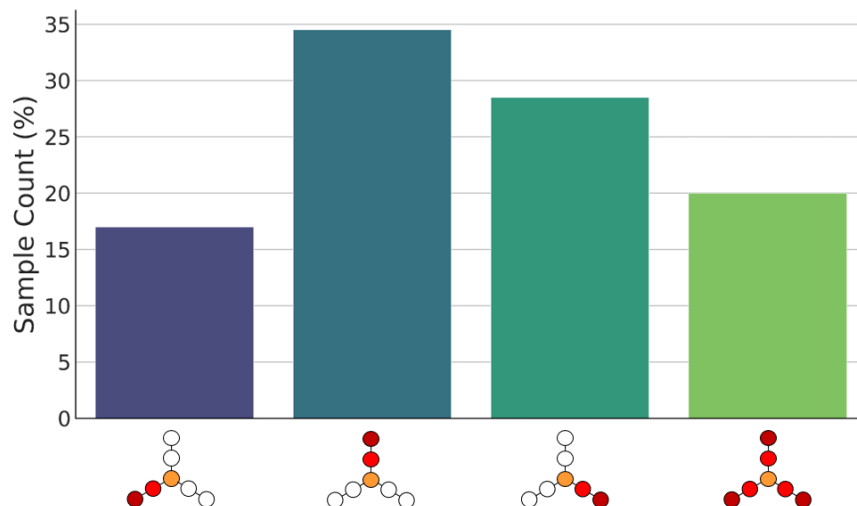


Figure 15. Histogram displaying the distribution of test samples for the *TreeNetwork90*, $\gamma = 0.999$ condition across different k-means groupings. K-means over the predicted successor representations for these samples with $k = 4$ (3 high-value nodes plus a ‘miscellaneous’ bin) exposed groupings based on approximate attack path. The first three bins correspond to determinate attack paths to a single network branch, and therefore represent ‘confident’ predictions. The final bin contains samples for which no single path was discernible. We hypothesize that hedged predictions result from observing a set of path trajectories where multiple users were consumed by the Red agent.

Discussion

In this work, we evaluate the efficacy of Theory of Mind-inspired metacognitive model architectures in providing actionable insights into the goals and behavioural patterns of adversarial cyber agents. We find that, through observing past episodes of Blue and Red cyber agents interacting with each other within the YAWNING-TITAN environment, our graph-based Machine Theory of Mind approach, GIGO-ToM, can learn character embeddings that can be clustered based on the Red agent's preferred outcome, thereby providing a means to characterize observed offensive cyber agents. We also find that our architecture can accurately predict Red's preference over high-value nodes, as well as the intermediate nodes that are likely to come under attack as the episode unfolds.

We note that there are numerous avenues for future research in this area, including a number of experiments that one could run with our framework. For example, we are keen to design representative 'Sally Anne' test^{28,56} scenarios for cyber-defence to evaluate the extent to which GIGO-ToM can attribute false beliefs to observed agents. This would serve the dual purpose of both stress-testing the mental network, which was deactivated for our experiments, and conclusively determining whether GIGO-ToM makes its predictions based on a true theory of mind rather than learned shortcuts.⁵⁷

ToM models are also often used as opponent models,⁵⁸ raising questions as to whether GIGO-ToM could be extended to provide an explicit belief model for partially observable environments with a graph-based observation space.

Finally, we consider that our novel similarity measure, the Network Transport Distance (NTD), holds potential for the direct optimization as well as evaluation of predicted successor representations. The computation of the Wasserstein Distance (on which the NTD is based) requires the resolution of a linear program, and is therefore both computationally intensive and non-differentiable, precluding its use as a loss function. However, differentiable approximations exist. The Sinkhorn Distance, for example, introduces an entropic regularization term to the optimal transport problem that smooths the objective function, permitting well-defined gradients that can be computed efficiently.⁵⁹ This raises the possibility of repurposing the Network Transport Distance as a loss function against which predicted graph-based successor representation can be optimized directly. To broach this topic, we developed a Sinkhorn-based NTD loss function and ran the preliminary test of comparing evaluation/test set NTD scores produced by this loss and the soft label cross-entropy loss when optimizing GIGO-ToM *only* with respect to successor representation predictions (i.e. dropping \mathcal{L}_{hvt} from Equation 7). Our initial findings suggest that the NTD may significantly improve predicted successor representations when used as a loss function. Preliminary results obtained using this loss function are provided in Appendix C.

In summary, our work shows that GNN-based ToMnet architectures show promise for the field of cyber defence. Under the abstract experimental conditions of this paper, the breadth and accuracy of GIGO-ToM's predictions are of a sufficient caliber to be leveraged by cyber-defence specialists to anticipate the goals and behaviours of observed cyber attackers in real time, enabling focused and timely defensive measures. We deem these findings compelling enough to warrant further investigation, and it is our hope that they are sufficient to inspire others to carry this work forward to more challenging and realistic cyber defence scenarios.

We hope to have also demonstrated that our graph-based similarity metric, the Network Transport Distance, provides a more flexible, intuitive, and interpretable means of comparing graph-based probability distributions than anything found during our review of existing literature. Beyond the obvious merits of a bounded and network-agnostic evaluation measure, our weighting method enables creative strategies for fine-tuning one's understanding of network vulnerabilities and attacker behaviour. We hold that these constitute valuable additions to the network analysis toolbox, with potential applications beyond the realm of cyber defence.

We hope that this work will raise awareness regarding the potential benefits of graph-based ToM for cyber-defence and that our work will inspire readers to attempt to answer some of the research questions that we have raised.

References

1. Pendleton, S. D. *et al.* Perception, planning, control, and coordination for autonomous vehicles (2017).
2. Scharre, P. *Army of none: Autonomous weapons and the future of war* (WW Norton & Company, 2018).
3. Yang, G.-Z. *et al.* Medical robotics—regulatory, ethical, and legal considerations for increasing levels of autonomy (2017).
4. Benitti, F. B. V. Exploring the educational potential of robotics in schools: A systematic review (2012).
5. Preece, A., Harborne, D., Braines, D., Tomsett, R. & Chakraborty, S. Stakeholders in explainable ai (2018).
6. Arulkumaran, K., Deisenroth, M. P., Brundage, M. & Bharath, A. A. Deep reinforcement learning: A brief survey, DOI: [10.1109/msp.2017.2743240](https://doi.org/10.1109/msp.2017.2743240) (2017).
7. Tang, Y. *et al.* Perception and navigation in autonomous systems in the era of learning: A survey (2022).
8. Deng, Y. *et al.* Deep learning-based autonomous driving systems: A survey of attacks and defenses (2021).

9. Adawadkar, A. M. K. & Kulkarni, N. Cyber-security and reinforcement learning—a brief survey (2022).
10. Li, C. & Qiu, M. *Reinforcement learning for cyber-physical systems: with cybersecurity case studies* (Chapman and Hall/CRC, 2019).
11. Rudin, C. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead (2019).
12. Du, M., Liu, N. & Hu, X. Techniques for interpretable machine learning (2019). [1808.00033](#).
13. Palmer, G., Parry, C., Harrold, D. J. & Willis, C. Deep reinforcement learning for autonomous cyber operations: A survey (2023).
14. Miles, I. *et al.* Reinforcement learning for autonomous resilient cyber defence (2024).
15. Deng, S. *et al.* Edge intelligence: The confluence of edge computing and artificial intelligence (2020).
16. NCSC. The near-term impact of ai on the cyber threat (2024). Accessed: 25/01/2024.
17. Nisioti, A., Mylonas, A., Yoo, P. D. & Katos, V. From intrusion detection to attacker attribution: A comprehensive survey of unsupervised methods (2018).
18. Kala, E. S. M. Critical role of cyber security in global economy (2023).
19. Police, A. F., Commission, A. C. I. *et al.* Acsc annual cyber threat report: July 2019 to june 2020 (2020).
20. Moustafa, N., Koroniotis, N., Keshk, M., Zomaya, A. Y. & Tari, Z. Explainable intrusion detection for cyber defences in the internet of things: Opportunities and solutions (2023).
21. Greydanus, S., Koul, A., Dodge, J. & Fern, A. Visualizing and understanding atari agents. In *International conference on machine learning, 1792–1801* (PMLR, 2018).
22. Burke, A. Robust artificial intelligence for active cyber defence (2020).
23. Antoniadi, A. M. *et al.* Current challenges and future opportunities for xai in machine learning-based clinical decision support systems: a systematic review (2021).
24. Miller, T. Explanation in artificial intelligence: Insights from the social sciences (2019).
25. Lipton, Z. C. The mythos of model interpretability (2018).
26. Samek, W., Montavon, G., Vedaldi, A., Hansen, L. K. & Müller, K.-R. *Explainable AI: interpreting, explaining and visualizing deep learning*, vol. 11700 (Springer Nature, 2019).
27. Premack, D. & Woodruff, G. Does the chimpanzee have a theory of mind? (1978).
28. Baron-Cohen, S., Leslie, A. M. & Frith, U. Does the autistic child have a “theory of mind”? (1985).
29. Garfield, J. L., Peterson, C. C. & Perry, T. Social cognition, language acquisition and the development of the theory of mind (2001).
30. Hamilton, A. F. d. C. Research review: Goals, intentions and mental states: Challenges for theories of autism (2009).
31. Ahmed, F. S. & Stephen Miller, L. Executive function mechanisms of theory of mind (2011).
32. Nichols, S. & Stich, S. P. *Mindreading: An integrated account of pretence, self-awareness, and understanding other minds* (Oxford University Press, 2003).
33. Rabinowitz, N. *et al.* Machine theory of mind. In *International conference on machine learning*, 4218–4227 (PMLR, 2018).
34. Dayan, P. Improving generalization for temporal difference learning: The successor representation (1993).
35. Andrew, A., Spillard, S., Collyer, J. & Dhir, N. Developing optimal causal cyber-defence agents via cyber security simulation. In *Workshop on Machine Learning for Cybersecurity (MLACyber)* (2022).
36. Rubner, Y., Tomasi, C. & Guibas, L. J. A metric for distributions with applications to image databases. In *Sixth international conference on computer vision (IEEE Cat. No. 98CH36271)*, 59–66 (IEEE, 1998).
37. Wang, Y., Zhong, F., Xu, J. & Wang, Y. Tom2c: Target-oriented multi-agent communication and cooperation with theory of mind (2021).
38. Shu, T. *et al.* Agent: A benchmark for core psychological reasoning. In *International Conference on Machine Learning*, 9614–9625 (PMLR, 2021).

39. Jeon, H., Oh, S., You, W., Jung, H. & Kim, K.-J. Inferring relationship using theory of mind in press diplomacy. In *ICML 2022 Workshop AI for Agent-Based Modelling* (2022).
40. Piazza, N. & Behzadan, V. A theory of mind approach as test-time mitigation against emergent adversarial communication (2023).
41. Cheng, Q., Wu, C., Hu, B., Kong, D. & Zhou, B. Think that attackers think: Using first-order theory of mind in intrusion response system. In *2019 IEEE Global Communications Conference (GLOBECOM)*, 1–6, DOI: [10.1109/GLOBECOM38437.2019.9013291](https://doi.org/10.1109/GLOBECOM38437.2019.9013291) (2019).
42. Malloy, T. & Gonzalez, C. Learning to defend by attacking (and vice-versa): Transfer of learning in cybersecurity games. In *2023 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, 458–464, DOI: [10.1109/EuroSPW59978.2023.00056](https://doi.org/10.1109/EuroSPW59978.2023.00056) (2023).
43. Von Stackelberg, H. *Market structure and equilibrium* (Springer Science & Business Media, 2010).
44. Wang, J., Xu, C., Yang, W. & Yu, L. A normalized gaussian wasserstein distance for tiny object detection (2021).
45. Kim, D. H., Yun, I. D. & Lee, S. U. A new attributed relational graph matching algorithm using the nested structure of earth mover’s distance. In *Proceedings of the 17th International Conference on Pattern Recognition, 2004. ICPR 2004.*, vol. 1, 48–51 (IEEE, 2004).
46. Noels, S., Vandermarliere, B., Bastiaensen, K. & De Bie, T. An earth mover’s distance based graph distance metric for financial statements. In *2022 IEEE Symposium on Computational Intelligence for Financial Engineering and Economics (CIFER)*, 1–8 (IEEE, 2022).
47. Cheong, O., Gudmundsson, J., Kim, H.-S., Schymura, D. & Stehn, F. Measuring the similarity of geometric graphs. In *Experimental Algorithms: 8th International Symposium, SEA 2009, Dortmund, Germany, June 4-6, 2009. Proceedings 8*, 101–112 (Springer, 2009).
48. Majhi, S. Graph mover’s distance: An efficiently computable distance measure for geometric graphs (2023).
49. Shapley, L. S. Stochastic games (1953).
50. Hochreiter, S. & Schmidhuber, J. Long short-term memory (1997).
51. Vaswani, A. *et al.* Attention is all you need (2017).
52. Veličković, P. *et al.* Graph attention networks. In *International Conference on Learning Representations* (2018).
53. Alon, U. & Yahav, E. On the bottleneck of graph neural networks and its practical implications (2020).
54. Doshi-Velez, F. & Kim, B. Towards a rigorous science of interpretable machine learning (2017).
55. Van der Maaten, L. & Hinton, G. Visualizing data using t-SNE. (2008).
56. Wimmer, H. & Perner, J. Beliefs about beliefs: Representation and constraining function of wrong beliefs in young children’s understanding of deception (1983).
57. Geirhos, R. *et al.* Shortcut learning in deep neural networks (2020).
58. Piazza, N., Behzadan, V. & Sarkadi, S. Limitations of theory of mind defenses against deception in multi-agent systems (2023).
59. Sinkhorn, R. & Knopp, P. Concerning nonnegative matrices and doubly stochastic matrices (1967).
60. Viehmann, T. Implementation of batched sinkhorn iterations for entropy-regularized wasserstein loss (2019).

Acknowledgements

Research funded by Frazer-Nash Consultancy Ltd. on behalf of the Defence Science and Technology Laboratory (Dstl) which is an executive agency of the UK Ministry of Defence providing world class expertise and delivering cutting-edge science and technology for the benefit of the nation and allies. The research supports the Autonomous Resilient Cyber Defence (ARCD) project within the Dstl Cyber Defence Enhancement programme.

Author contributions statement

All authors provided critical feedback and helped shape the research. L.S and G.P contributed to the development and evaluation of the GIDO-ToM and GIGO-ToM approaches, along with the experiment design. L.S formulated the Network Transport Distance for evaluating successor representations within the context of cyber-defence scenarios. D.H. was responsible for

implementing the graph-based observation wrapper for YAWNING-Titan. D.H. and M.S were responsible for implementing rules-based cyber-attack and cyber-defence agents for data gathering and evaluation. C.W. provided an extensive technical review. G.P. coordinated and supervised the project.

Data availability

The framework and data that support the findings of this study are proprietary to BAE Systems and restrictions apply to the availability of these items which are not publicly available. Request for access to these items can be sent to Luke Swaby (luke.swaby2@baesystems.com) and Gregory Palmer (gregory.palmer@baesystems.com). Requests will be subject to BAE Systems standard data protection policies and processes.

Appendices

A Additional Custom YAWNING-TITAN Network Topologies

In addition to the custom network topologies described in the [Preliminaries](#) section, we developed two more complex ones for more rigorous future testing:

- *ForestNetwork*: A network that resembles a forest layout, consisting of 72 nodes, shown in [Figure 16a](#).
- *OpticalCoreNetwork*: A network consisting of 54 nodes, where in the default setting the high value nodes are situated on one of our four servers at the centre of the network depicted in [Figure 16b](#).

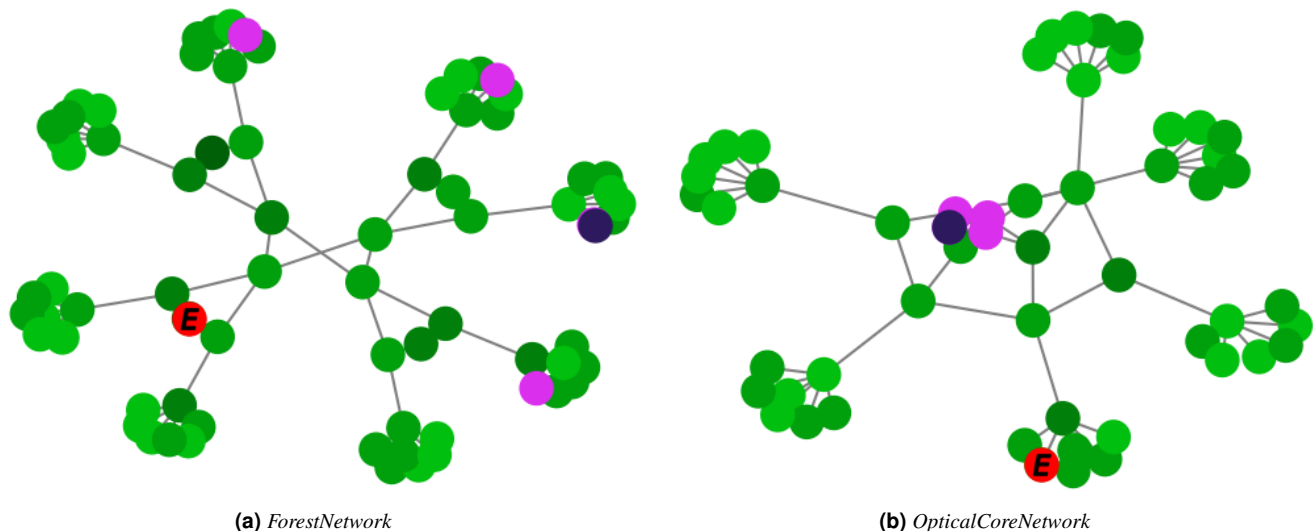


Figure 16. Additional YAWNING-TITAN Network Topologies

B Rule-Based Blue Cyber-Defence and Red Attacking Agents

In addition to the agents used for the experiments in this report, we initially developed a suite of other Blue cyber-defence and Red cyber-attacking agents for our evaluations. We first provide an overview of the actions available to each agent, before outlining the agents that we implemented. The goal of the Blue agent is to stop the Red agent's attack and survive for 500 time-steps. The actions that the Blue cyber-defence agents can take to this end are:

- *DoNothing*: performs no action.
- *Scan*: reveals any hidden compromised nodes.
- *MakeSafeNode*: removes the infection from a node.
- *ReduceNodeVulnerability*: lowers the vulnerability score of that node.

- `Restore`: uses `MakeSafeNode` and resets the node vulnerability at the start of the episode.
- `Isolate`: removes all connections from a connected node.
- `Reconnect`: returns all connections from an isolated node.

The goal of the Red cyber-attacking agent, on the other hand, is to evade the Blue agent in order to capture a high-value node within 500 time-steps, starting from a specific entry node. Red has its own set of offensive actions, including:

- `DoNothing`: performs no action.
- `BasicAttack`: attempts to infect a node.
- `RandomMove`: moves the red agent to another connected node.
- `ZeroDayAttack`: guaranteed infection on a node. Red agents have a `ZeroDayAttack` budget, gaining an additional `ZeroDayAttack` at a defined rate per time-step.
- `Spread`: attempts to infect all nodes connected to any compromised node.
- `Intrude`: attempts to infect all nodes.

For our Blue agents, we have:

- `BlueSleep`: Takes the `DoNothing` action at every step, used to benchmark the red agent's performance as if there were no blue agent present on the network
- `BlueRandom`: Selects a random action from the action-space, and then a random node if that action requires a target. This is used to benchmark against a completely random blue policy with no defence strategy.
- `BlueRandomSmart`: Selects a random action from the entire action-space, but chooses an appropriate node for that action. For example, it will only use `MakeSafeNode` and `Restore` on nodes that the agent can observe are compromised. This is used to observe how a blue agent would perform if it had no strategy in regards to its action selection, but made logical choices for its targets.
- `BlueIsolate`: Takes the `Isolate` action on the entry and high-value nodes at the start of the episode, cutting all possible routes the red agent has to the high value nodes. It will then take `Scan` for the rest of the episode until a new compromised node is revealed, where it will then take `MakeSafeNode` on the compromised nodes. This exploits the rules of the game for a guaranteed blue victory, but can be used to benchmark the score an agent would get through following this scorched-earth policy.

We have six agents that utilise the `MakeSafeNode` or `Restore` actions more intelligently to mitigate the attack. Each of these operate by selecting a node based on its distance from any high-value node, assigning greater priority to compromised nodes that are close to a high-value node:

- `BlueMSN-D`: A deterministic agent that takes the `MakeSafeNode` action if a compromised node is within three hops of any high value node, otherwise takes `Scan`. The strategy here is to protect only the core of the network. Blue is content with sacrificing the edge nodes to do so.
- `BlueMSN-S`: A stochastic agent. Takes `MakeSafeNode` and `Scan` based on a probability. The `MakeSafeNode` probability increases if a compromised node is close to a high-value node. This agent's strategy is to protect the entire network, but will look to remove infections closer to the high value nodes with greater urgency.
- `BlueRestore`: The same as `BlueMSN-S`, but taking `Restore` rather than `MakeSafeNode` to observe the difference in performance between the `MakeSafeNode` and `Restore` actions.
- `BlueMSN-RNV`: The same as `BlueMSN-S`, but will uniformly select between `Scan` or `ReduceNodeVulnerability` of the most vulnerable node if not taking `MakeSafeNode`. This is to assess if it is better for the agent to prioritise monitoring the network, or to make the network more resilient to attacks.
- `BlueMSN-Restore`: The same as `BlueMSN-S`, but will randomly select `MakeSafeNode` or `Restore` if not taking `Scan`. This can be used to train learning agents under what conditions it would be better for the blue agent to take `MakeSafeNode` or `Restore`.

- **BlueMSN-RNV-Restore:** This agent uniformly selects between `MakeSafeNode` and `Restore` when removing infections, and between `Scan` and `ReduceNodeVulnerability` otherwise to detect and protect against red activity respectively. As with the other stochastic agents, the probability that it selects between these two sets of actions corresponds to the distance the nearest compromised node is to a high value node. This allows the agent to sample the full range of defensive capabilities, excluding the `Isolate` action which was shown to be too powerful and restrictive as an action. Combines `MakeSafeNode – ReduceNodeVulnerability` and `MakeSafeNode – Restore` to take all four actions based on the action probabilities.

It is worth noting that only the `BlueRandom` and `BlueRandomSmart` take the `Reconnect` action. This is because the action is only useful to reconnect isolated nodes, and no other agents take the `Isolate` action due to its power relative to the other defensive actions. Using `Reconnect` can improve the agent’s score as it is punished for each isolated node in the reward function, but is redundant without `Isolate`.

For our cyber-attacking agents, we seven Red agents with action probabilities sampled from a Dirichlet distribution $\pi \sim Dir(\alpha)$:

- **RedRandomSimple:** Selects an action weighted by above action probabilities and selects random valid targets. This agent simulates an unskilled attacker with no prior knowledge or plan of attack. It occasionally wastes time trying to use `ZeroDayAttack`, even when unavailable.
- **RedRandomSmart:** First selects an action and target exactly like `RedRandomSimple` but if `ZeroDayAttack` is chosen and, if none are available, uses a `BasicAttack` instead. This is a more proficient version of `RedRandomSimple`.
- **RedTargetConnected:** Selects an action exactly like `RedRandomSmart` but always selects the target with most connections. This agent uses a disruptive attack strategy as its targets are highly connected, trying to reveal as much of the network as possible.
- **RedTargetUnconnected:** Selects an action exactly like `RedRandomSmart` but always selects the target with the least connections. Using the opposite strategy to `RedTargetConnected`, this agent makes the assumption that the most valuable data would be stored on the most secluded nodes.
- **RedTargetVulnerable:** Selects an action exactly like `RedRandomSmart` but always selects the target with highest vulnerability. The strategy of this agent is to compromise as many nodes as quickly as possible because the more vulnerable the target, the more likely the attack is to succeed.
- **RedTargetResilient:** Selects an action exactly like `RedRandomSmart` but always selects the target with lowest vulnerability. Using the opposite strategy to `RedTargetVulnerable`, this agent perhaps makes the assumption that the most valuable data would be stored on the most resilient nodes.
- **RedHVTSimple:** Behaves exactly like `RedRandomSimple` until a high value node becomes connected to a compromised node, at which point actions and target choices become deterministic – `ZeroDayAttack` if available, else `BasicAttack` on a high value node. This agent still has no prior knowledge of the network but knows the importance of compromising the high-value node.

There are two red agents with high value node preferences, instead of action probabilities, sampled from a Dirichlet distribution:

- **RedHVTPreference:** Selects a particular high value node to aim for at the start of the episode. The selected high-value node is based on the high value node preferences as well as the shortest distance between high value nodes and entry nodes. This attacker has prior knowledge of the network structure, for example insider information, meaning it can calculate path lengths and is then likely to take the most direct path to its chosen high-value node.
- **RedHVTPreferenceSP:** Selects a particular high value node to aim for, exactly like `RedHVTPreference`, but always takes the shortest path. This makes it a more ruthless version of `RedHVTPreference`.

Rule-Based Agent Evaluation

Upon implementing the agents described in the above section we conducted an extensive evaluation, allowing us to answer the question: ‘What are the optimal mixtures over our respective sets of agents that Red and Blue can sample from for different game settings?’. Here, we conduct an empirical evaluation using three network topologies: *TreeNetwork30*; *OpticalCoreNetwork* and *ForestNetwork*. Each network is implemented with three entry nodes. As always, high-value nodes are situated on leaf nodes.

We conduct 100 evaluation episodes for each Blue, Red joint-policy profile, and subsequently evaluate the agents with respect to Blue rewards, Blue win rate, and episode duration. For the latter, a 500 step episode represents a Blue victory.

We provide detailed tables with the results from each of the selected topologies in tables 4 – 12. Tables 1 – 3 meanwhile provide the averages. It should not come as a surprise that the strongest Blue cyber-defence agent with respect to win rate and duration is BlueIsolate (See Tables 2 and 3). However, isolating nodes comes at a cost with respect to the Blue reward (See Table 1). When discarding BlueIsolate it becomes apparent that the Red cyber-attacking RedHVTPreferenceSP presents a particular challenge to the remaining Blue cyber-defence agents. Blue’s best response against RedHVTPreferenceSP is BlueRandomSmart with respect to average duration (140 steps, see Table 3) and win rate (0.21, see Table 2). However, due to selecting the isolate, BlueRandomSmart performs poorly with respect to the Blue reward function (-942.48, see Table 1). In contrast BlueMSN-D achieves the best result with respect to Blue rewards (-58.37) against RedHVTPreferenceSP, representing the pure strategy Nash equilibrium within the Blue rewards payoff table. Other strong performers against RedHVTPreferenceSP with respect to the reward criteria include BlueMSN-S (-66.23), BlueMSN-RNV (-69.13), BlueMSN-RNV-Restore (-78.6) and BlueMSN-Restore (-78.01).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	-3162.04	-3160.1	-3224.54	-3224.76	-3184.8	-3183.72	-3186.57	-3184.71	-3187.16
BlueMSNDeterministic	73.13	-58.37	23.93	49.47	66.35	71.49	72.76	65.88	72.53
BlueMSNRNV	148.31	-69.13	79.05	91.68	130.99	118.46	125.89	129.02	124.08
BlueMSNRNVRestore	84.27	-78.6	37.7	50.26	69.68	64.44	67.07	75.75	68.58
BlueMSNRestore	37.86	-78.01	15.32	26.88	35.65	35.41	37.07	38.35	35.72
BlueMSNStochastic	131.29	-66.23	65.24	80.75	114.8	109.86	111.34	117.49	113.62
BlueRandom	-896.23	-310.67	-2031.29	-2346.51	-1184.87	-1091.91	-1063.61	-1261.43	-1111.15
BlueRandomSmart	-1611.56	-942.48	-1822.24	-1927.68	-1820.9	-1789.49	-2004.85	-1693.16	-2006.49
BlueRestore	-53.6	-91.03	-31.73	-28.06	-42.45	-41.26	-41.98	-43.57	-41.12
BlueSleep	-217.7	-97.13	-203.09	-321.19	-260.9	-265.46	-262.07	-234.79	-276.14

Table 1. Rewards for deterministic red agents ($\alpha = 0.01$). Averaged across the three networks.

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
BlueMSNDeterministic	0.8	0.01	0.72	0.97	0.91	0.94	0.96	0.9	0.97
BlueMSNRNV	0.93	0.0	0.89	0.99	0.99	0.94	0.97	0.98	0.96
BlueMSNRNVRestore	0.92	0.0	0.86	0.99	0.96	0.95	0.95	0.97	0.96
BlueMSNRestore	0.9	0.0	0.87	1.0	0.94	0.95	0.98	0.98	0.96
BlueMSNStochastic	0.91	0.01	0.86	0.99	0.95	0.94	0.96	0.96	0.98
BlueRandom	0.03	0.0	0.07	0.07	0.02	0.02	0.01	0.03	0.02
BlueRandomSmart	0.34	0.21	0.84	0.87	0.54	0.48	0.59	0.58	0.63
BlueRestore	0.87	0.01	0.88	1.0	0.97	0.95	0.96	0.97	0.97
BlueSleep	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Table 2. Win rates for deterministic red agents ($\alpha = 0.01$). Averaged across the three networks.

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0
BlueMSNDeterministic	449.56	93.08	420.57	494.35	478.29	483.69	489.04	475.2	491.38
BlueMSNRNV	479.8	63.14	466.51	497.94	499.63	477.65	490.29	493.28	488.06
BlueMSNRNVRestore	478.71	63.01	454.77	495.78	490.37	484.46	485.07	492.36	490.36
BlueMSNRestore	472.6	74.8	459.0	500.0	486.34	487.57	498.01	497.14	490.38
BlueMSNStochastic	480.32	74.4	459.33	497.89	489.79	485.24	489.48	490.65	495.17
BlueRandom	83.74	30.29	143.59	168.45	104.1	99.85	99.34	102.85	101.81
BlueRandomSmart	222.36	140.0	431.41	448.49	305.39	280.16	327.97	327.22	347.86
BlueRestore	464.6	75.0	467.64	499.49	494.09	487.17	487.92	491.51	492.01
BlueSleep	59.68	15.27	77.87	110.43	75.07	76.0	77.79	69.98	78.41

Table 3. Durations for deterministic red agents ($\alpha = 0.01$). Averaged across the three networks.

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0
BlueMSNDeterministic	489.09	114.04	481.36	500.0	496.66	496.42	496.66	493.63	497.05
BlueMSNRNV	497.58	78.7	500.0	500.0	500.0	495.75	496.38	500.0	495.72
BlueMSNRNVRestore	497.06	86.04	500.0	500.0	495.15	492.84	500.0	500.0	500.0
BlueMSNRestore	500.0	98.67	490.22	500.0	498.15	500.0	500.0	500.0	500.0
BlueMSNStochastic	500.0	103.24	498.83	500.0	500.0	498.07	500.0	500.0	500.0
BlueRandom	115.55	46.74	206.22	238.0	142.91	137.6	124.21	160.49	134.45
BlueRandomSmart	309.68	199.91	472.02	474.57	381.63	311.78	385.71	393.23	391.36
BlueRestore	500.0	109.11	500.0	500.0	497.22	497.07	498.44	500.0	500.0
BlueSleep	85.84	19.19	116.43	159.25	109.08	111.73	114.76	99.46	112.51

Table 4. ForestNetwork Durations for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	-3326.51	-3324.55	-3388.72	-3389.75	-3351.34	-3349.15	-3350.66	-3348.62	-3349.05
BlueMSNDeterministic	61.69	-49.9	26.19	29.28	47.19	51.41	52.88	46.24	50.9
BlueMSNRNV	164.33	-60.56	96.27	97.44	137.19	133.21	135.36	139.25	134.96
BlueMSNRNVRestore	99.25	-70.38	57.76	56.46	78.21	74.92	81.34	85.75	80.55
BlueMSNRestore	45.61	-70.58	23.1	26.86	38.82	37.8	37.69	39.16	39.06
BlueMSNStochastic	139.9	-53.31	80.95	81.08	120.31	116.78	117.61	120.49	116.47
BlueRandom	-1424.54	-534.28	-3342.88	-4057.49	-1993.61	-1776.46	-1459.18	-2403.73	-1679.39
BlueRandomSmart	-2726.55	-1552.57	-2243.33	-2189.74	-2571.26	-2366.15	-2587.6	-2289.04	-2423.3
BlueRestore	-49.97	-87.2	-27.35	-28.01	-42.17	-39.82	-40.71	-42.19	-39.3
BlueSleep	-329.09	-98.58	-324.32	-544.88	-419.92	-434.26	-428.0	-370.02	-430.67

Table 5. ForestNetwork Rewards for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
BlueMSNDeterministic	0.94	0.02	0.9	1.0	0.98	0.98	0.98	0.97	0.99
BlueMSNRNV	0.99	0.0	1.0	1.0	1.0	0.99	0.99	1.0	0.99
BlueMSNRNVRestore	0.99	0.0	1.0	1.0	0.99	0.98	1.0	1.0	1.0
BlueMSNRestore	1.0	0.0	0.97	1.0	0.99	1.0	1.0	1.0	1.0
BlueMSNStochastic	1.0	0.01	0.99	1.0	1.0	0.98	1.0	1.0	1.0
BlueRandom	0.04	0.0	0.13	0.16	0.04	0.04	0.01	0.06	0.01
BlueRandomSmart	0.5	0.3	0.93	0.93	0.69	0.53	0.71	0.72	0.72
BlueRestore	1.0	0.02	1.0	1.0	0.99	0.98	0.99	1.0	1.0
BlueSleep	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Table 6. ForestNetwork Results for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0
BlueMSNDeterministic	396.95	77.3	354.84	483.97	453.53	465.29	477.26	448.41	480.57
BlueMSNRNV	455.59	42.39	413.55	493.81	498.91	454.31	478.2	484.2	473.48
BlueMSNRNVRestore	451.12	36.26	380.7	487.34	480.33	475.23	455.21	477.75	474.48
BlueMSNRestore	425.16	42.88	402.8	500.0	466.32	477.06	494.03	491.5	471.41
BlueMSNStochastic	448.11	44.76	391.6	493.67	470.93	461.36	472.49	473.35	489.18
BlueRandom	63.96	20.7	86.29	118.74	77.49	62.45	69.53	51.95	75.28
BlueRandomSmart	135.24	99.6	387.01	402.58	197.6	214.59	255.34	246.6	299.48
BlueRestore	404.04	45.81	412.19	498.46	486.0	470.5	468.49	476.09	479.13
BlueSleep	32.69	11.86	36.85	67.32	42.24	41.68	43.49	40.72	41.52

Table 7. TreeNetwork Durations for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	-2829.86	-2827.55	-2893.26	-2893.41	-2851.54	-2848.26	-2856.28	-2852.96	-2858.54
BlueMSNDeterministic	79.76	-64.6	18.24	67.78	84.96	93.09	94.65	81.95	94.09
BlueMSNRNV	130.03	-79.93	55.39	87.18	126.07	105.13	114.59	120.99	111.48
BlueMSNRNVRestore	68.53	-87.83	10.57	45.49	62.27	56.84	50.97	65.74	58.32
BlueMSNRestore	23.93	-87.26	1.59	27.6	32.0	34.74	37.32	37.18	30.7
BlueMSNStochastic	118.59	-79.5	40.69	80.55	109.88	101.89	106.47	113.33	112.43
BlueRandom	-620.59	-189.51	-851.85	-1152.75	-709.15	-537.54	-626.11	-377.44	-672.24
BlueRandomSmart	-680.27	-547.88	-1450.26	-1533.96	-894.75	-1036.93	-1270.16	-1049.37	-1429.55
BlueRestore	-61.03	-94.48	-39.79	-29.3	-45.41	-45.14	-45.86	-47.39	-44.91
BlueSleep	-121.0	-96.04	-102.48	-154.69	-135.26	-130.76	-134.46	-125.67	-132.28

Table 8. TreeNetwork Rewards for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
BlueMSNDeterministic	0.63	0.01	0.57	0.93	0.81	0.89	0.92	0.78	0.94
BlueMSNRNV	0.87	0.0	0.74	0.98	0.97	0.87	0.92	0.94	0.91
BlueMSNRNVRestore	0.81	0.0	0.67	0.97	0.9	0.9	0.86	0.93	0.9
BlueMSNRestore	0.74	0.0	0.68	1.0	0.87	0.88	0.95	0.94	0.88
BlueMSNStochastic	0.77	0.01	0.63	0.97	0.87	0.86	0.9	0.89	0.95
BlueRandom	0.03	0.0	0.02	0.03	0.01	0.02	0.02	0.0	0.03
BlueRandomSmart	0.19	0.16	0.75	0.77	0.32	0.36	0.45	0.42	0.54
BlueRestore	0.65	0.0	0.68	0.99	0.95	0.89	0.91	0.93	0.92
BlueSleep	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Table 9. TreeNetwork Results for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0	500.0
BlueMSNDeterministic	462.64	87.91	425.52	499.09	484.69	489.36	493.2	483.56	496.52
BlueMSNRNV	486.24	68.34	485.98	500.0	499.98	482.89	496.28	495.63	494.98
BlueMSNRNVRestore	487.96	66.74	483.6	500.0	495.63	485.3	500.0	499.33	496.61
BlueMSNRestore	492.64	82.86	483.97	500.0	494.56	485.66	500.0	499.93	499.74
BlueMSNStochastic	492.84	75.2	487.57	500.0	498.44	496.28	495.96	498.6	496.32
BlueRandom	71.7	23.44	138.25	148.6	91.9	99.51	104.29	96.12	95.71
BlueRandomSmart	222.15	120.5	435.21	468.32	336.94	314.1	342.87	341.83	352.74
BlueRestore	489.76	70.09	490.74	500.0	499.05	493.93	496.82	498.43	496.91
BlueSleep	60.51	14.77	80.33	104.73	73.89	74.58	75.13	69.75	81.21

Table 10. OpticalCoreNetwork Durations for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	-3329.77	-3328.19	-3391.65	-3391.14	-3351.51	-3353.75	-3352.77	-3352.54	-3353.9
BlueMSNDeterministic	77.95	-60.62	27.38	51.34	66.9	69.96	70.73	69.45	72.58
BlueMSNRNV	150.58	-66.89	85.48	90.43	129.72	117.03	127.73	126.83	125.81
BlueMSNRNVRestore	85.04	-77.58	44.76	48.83	68.55	61.57	68.88	75.77	66.88
BlueMSNRestore	44.05	-76.2	21.26	26.17	36.12	33.7	36.21	38.7	37.41
BlueMSNStochastic	135.37	-65.9	74.08	80.63	114.22	110.9	109.93	118.66	111.96
BlueRandom	-643.57	-208.21	-1899.14	-1829.28	-851.86	-961.73	-1105.55	-1003.12	-981.82
BlueRandomSmart	-1427.86	-727.01	-1773.13	-2059.36	-1996.68	-1965.39	-2156.79	-1741.08	-2166.61
BlueRestore	-49.8	-91.42	-28.04	-26.89	-39.78	-38.81	-39.36	-41.12	-39.13
BlueSleep	-203.02	-96.78	-182.47	-264.0	-227.51	-231.36	-223.75	-208.67	-265.48

Table 11. OpticalCoreNetwork Rewards for deterministic red agents ($\alpha = 0.01$).

	RedHVTPreference	RedHVTPreferenceSP	RedHVTSimple	RedRandomSimple	RedRandomSmart	RedTargetConnected	RedTargetResilient	RedTargetUnconnected	RedTargetVulnerable
BlueIsolate	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0	1.0
BlueMSNDeterministic	0.83	0.01	0.69	0.98	0.94	0.95	0.98	0.94	0.98
BlueMSNRNV	0.94	0.0	0.93	1.0	0.99	0.96	0.99	0.99	0.98
BlueMSNRNVRestore	0.96	0.0	0.91	1.0	0.98	0.96	1.0	0.99	0.99
BlueMSNRestore	0.96	0.0	0.95	1.0	0.97	0.96	1.0	0.99	0.99
BlueMSNStochastic	0.97	0.0	0.96	1.0	0.99	0.99	0.98	0.99	0.99
BlueRandom	0.01	0.0	0.07	0.03	0.0	0.01	0.0	0.03	0.02
BlueRandomSmart	0.32	0.17	0.84	0.92	0.6	0.55	0.6	0.6	0.63
BlueRestore	0.96	0.0	0.97	1.0	0.98	0.98	0.99	0.99	0.98
BlueSleep	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0

Table 12. OpticalCoreNetwork Results for deterministic red agents ($\alpha = 0.01$).

C The Network Transport Distance as a Loss Function

In our discussion of this work, we posit that the Network Transport Distance (NTD) holds potential for the direct optimization of predicted successor representations as well as evaluation. While a comprehensive exploration lies beyond the scope of the current study, we here provide some preliminary experimental insights into this avenue.

In order to recast the NTD as a loss function, we took inspiration from an open-source PyTorch implementation of the Sinkhorn loss⁶⁰ and replaced the NTD’s enclosed linear program with Sinkhorn-Knopp iterations. This modification not only preserves computational efficiency but also ensures differentiability, thereby enabling its incorporation into gradient-based learning frameworks.

Mathematically, the extension can be formalized as as follows:

$$\text{NTD}_{\mathcal{L}}(P, Q, D, \lambda) = \frac{1}{\max(D)} \inf_{\mu \in \mathcal{M}(P, Q)} \left[\int_{X \times X} d(i, j) d\mu(i, j) - \lambda H(\mu) \right], \quad (15)$$

where:

- $\lambda > 0$ is a user-specified regularization parameter.
- $H(\mu) = -\sum_{i,j} \mu(i, j) \log \mu(i, j)$ is the entropy of the transport plan μ .

The optimal transport plan μ^* is obtained via Sinkhorn-Knopp iterations, which iteratively update the transport plan μ to satisfy the marginal constraints P and Q :

$$\mu^{(k+1)} = \text{diag}(u^{(k+1)}) \cdot K \cdot \text{diag}(v^{(k+1)}), \quad (16)$$

where:

- $K_{ij} = \exp\left(-\frac{D[i,j]}{\lambda}\right)$ is a kernel matrix.
- The scaling factors u and v are updated as:

$$u^{(k+1)} = \frac{P}{K v^{(k)}}, \quad v^{(k+1)} = \frac{Q}{K^T u^{(k+1)}}. \quad (17)$$

We then ran the preliminary test of comparing evaluation/test set NTD scores produced by this loss and the soft label cross-entropy loss when optimizing GIGO-ToM *only* with respect to successor representation predictions (i.e. dropping \mathcal{L}_{hvt} from Equation 7 and instead having $\mathcal{L}_{total} = \mathcal{L}_{sr}$). Our results suggest that predicted successor representations are significantly improved when optimized with the $\text{NTD}_{\mathcal{L}}$ over the SXE (Figure 17).

Successor Representation Loss Function (\mathcal{L}_{sr}) Comparison

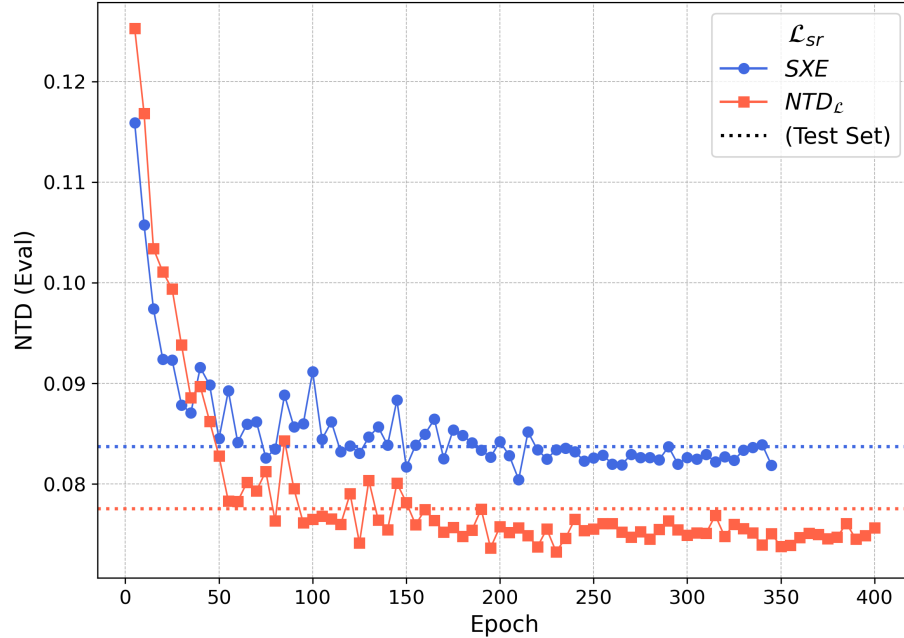


Figure 17. Evaluation and test set NTD scores for GIGO-ToM’s predicted successor representations over *TreeNetworkMixed* when training GIGO-ToM only to optimize SR predictions using our NTD loss function, benchmarked against the soft label cross-entropy loss (SXE). Mean NTD scores are computed across each $\gamma \in \{0.5, 0.95, 0.999\}$. Hold-out test set NTD scores are overlaid with dotted lines.