

See discussions, stats, and author profiles for this publication at: <https://www.researchgate.net/publication/330590120>

Reinforcement Learning for Intelligent Penetration Testing

Conference Paper · October 2018

DOI: 10.1109/WorldS4.2018.8611595

CITATIONS

31

READS

1,927

2 authors, including:



Mohamed chahine Ghanem

London Metropolitan University

11 PUBLICATIONS 103 CITATIONS

SEE PROFILE

Some of the authors of this publication are also working on these related projects:



Anti-Forensics: A Script-Based Tool for Extracting Evidence Hidden by Cryptographic and Steganographic Techniques [View project](#)

Article

Reinforcement Learning for Efficient Network Penetration Testing

Mohamed C. GHANEM and Thomas M. CHEN

¹ School of Mathematics Computer Science and Engineering; City, University of London;
mohamed.ghanem@city.ac.uk

² School of Mathematics Computer Science and Engineering; City, University of London;
tom.chen.1@city.ac.uk

Version November 21, 2019 submitted to Information

Abstract: Penetration testing (also known as pentesting or PT) is common practice for actively assessing the defences of a computer network by planning and executing all possible attacks to discover and exploit existing vulnerabilities. Current penetration testing methods are increasingly becoming non-standard, composite and resource consuming despite the use of evolving tools. In this paper, we propose and evaluate an AI-based pentesting system which makes use of machine learning techniques, namely reinforcement learning (RL) to learn and reproduce average and complex pentesting activities. The proposed system is named Intelligent Automated Penetration Testing System (IAPTS) and will be a module that integrates with industrial PT systems and frameworks to enable them to capture information, learn from experience and reproduce the test in future nearly similar testing cases. IAPTS aims to save human resources while producing much enhanced results in term of time consumption, reliability and frequency of testing. IAPTS takes the approach of modelling PT environments and tasks as a partially observed Markov decision process (POMDP) problem which is solved by POMDP-solver. Although this paper scope is limited to network infrastructures PT planning and not the entire practice, the obtained results support the hypothesis that RL can enhance PT beyond the capabilities of any human PT expert in terms of time consumed, covered attacking vectors, accuracy and reliability of the outputs. In addition, this work tackled the complex problem of expertise capturing and re-use by allowing the IAPTS learning module to store and re-use PT policies in the same way that a human PT expert would learn but in a more efficient way.

Keywords: penetration testing, artificial intelligence, machine learning, reinforcement learning, network security auditing, offensive cyber-security, vulnerability assessment.

1. Introduction

Computer networks are more than ever exposed to cyber threats of increasing frequency, complexity and sophistication [1]. Penetration Testing (shortly known as pentesting or PT) is a well-established proactive method to evaluate the security of digital assets, varying from a single computer to websites and networks, by actively searching for and exploiting the existing vulnerabilities. The practice is an emulation of the operational mode that hackers follow in real-world cyber attacks. In the current constantly evolving digital environment, PT is becoming a crucial and often mandatory component of cyber security auditing particularly after the introduction of the European GDPR (General Data Protection Regulation) for organizations and businesses. In addition to legal requirements, PT is considered by the cyber security community as the most effective method to assess the strength of security defences against skilled adversaries as well as the adherence to security policies [2]. In practical terms, PT is a multi-stage process that often requires a high degree of competence and expertise due to the complexity of digital assets such as medium and large networks. Naturally,

research has investigated the possibility of automating tools for the different PT stages (reconnaissance, identification, and exploitation) to relieve the human expert from the burden of repetitive tasks. However, automation by itself does not achieve much benefits in terms of time, resources and outputs because PT is a dynamic and interactive process of exploring and decision making, which requires advanced and critical cognitive skills that are hard to duplicate through automation.

A natural question arises in regard to the capability of AI to provide a potential solution that goes beyond simple automation to achieve expert-like output. In other research fields, AI proved very helpful to not only offload work from humans but also possibly handle depths and details that humans can not tackle fast enough or accurately enough. Rapid progress in the AI and notably machine learning (ML) sub-field led us to believe that an AI-based PT system utilizing well-grounded models and algorithms for making sequential decisions in uncertain environments can bridge the gap between automation and expertise that PT community experience. In this perspective, the existing PT systems and framework started shifting from executing experts' tasks to become more autonomous, intelligent and optimized aiming that all existing threats are checked systematically and efficiently without or with little human expert intervention. Furthermore, these systems should optimise the use of resources by eliminating time-consuming and irrelevant directions and ensure that no threat is overlooked.

In addition to the regular use of PT, the testing results (output) should be processed and stored to serve for further use. In fact, the main difference between human PT expert and automated systems is that humans learn alongside performing the tests and enrich their expertise throughout, while systems omit the re-usability of the data which is sometimes crucial especially when the testing is repeated such as for regular compliance tests. In practical terms, the vast majority of the assessed network configurations will not change considerably over a short period and therefore the output of previous tests could remain entirely or partly applicable for an eventual re-testing required after one or more of these following points occur:

- Network hardware, software, segments or applications were added, changed or removed
- Significant systems' upgrades or modifications are applied to infrastructure or applications
- Infrastructure changes including moving locations
- Security solutions or patches were installed or modified
- Security or users' policies were modified

Automation was and remains the best solution to save time and resources in any domain and PT is not an exception to this rule. Therefore, the offensive cyber security community accorded during the last decade a particular attention to the automation of the used systems. Such improvement permitted to save significant time, efforts and resources in performing the task. Given the particularity of PT practice, the increasing size and complexity of the tested assets along with the significant number of vulnerabilities, exploit and attacks' vectors which should cover by the tester, the blind automated system becomes powerless and often perform worse than manual practice pushing the researcher to focus on improving such systems by adopting a variety of solutions. This paper explores in deep the design and development of an ML-based PT system that allows intelligent, optimized and efficient testing by perceiving its environment and decide autonomously how to act in order to perform PT tasks as better as human experts and save time and resources along with improving accuracy and testing coverage and frequency.

1.1. Research Context

Performing a periodic offensive security testing and auditing is an essential process to ensure the resilience and the compliance of the assessed asset notably the confidentiality, availability and integrity. PT is reputed to be the best approach to assess the security of digital assets by identifying and exploiting its vulnerabilities. Currently, dozens of commercial and freeware systems, platforms and frameworks are being used by PT experts with some offering some automation features which nevertheless remain either local (specific to very limited context or tasks) or not optimized (blind

automation) and therefore creating significant accuracy and performance issues notably in case testing medium and large networks.

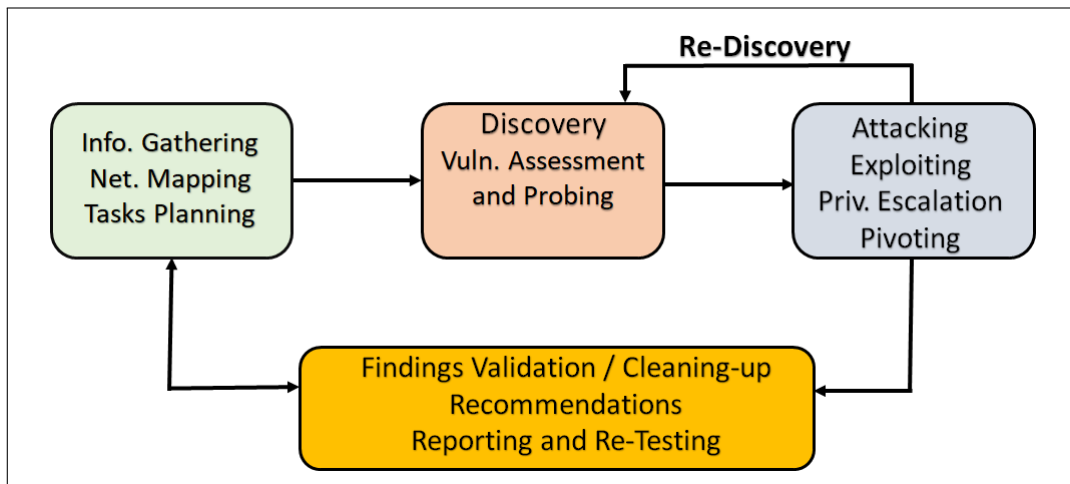


Figure 1. PT is a non-standard active method for assessing network defence by following a sequential and interactive multi-phase procedure starting by gathering information and ending by reporting the obtained results [1].

furthermore, others issues are usually related to the existing automated systems, notably the the congestion created in the assessed network triggering both security and performances issues along with the associated volume of data generated from the testing outputs that are often unexploited. Finally, PT environment is characterised as fast-changing and complex and the human experts are suffering from the complexity, repeatability and resemblance which in large and complex networks context such as large organisations using standard system and subsequently security protection. Performing PT in alike scenarios will create high degree of obfuscation and make it almost impossible to cover the whole asset properly [3-4].

During the last decade, the use of machine learning has intensified in the cyber security domain and especially in defensive applications such as intrusion detection and prevention systems (IDPS), Malware analysis and Anti-viruses solutions. Recently, MIT researchers developed a big data security framework named AI2 which combined security analyst expertise with machine learning to build an IDPS with active learning [2]. In the offensive cyber security domain, there have been few attempts to equip existing PT systems and frameworks with learning capabilities and thus for the obvious reason of complexity associated with PT practice notably into the modeling and design of an ML-led offensive security systems. In fact, it is natural to imagine that one or more machine learning techniques can be applied to different PT phases enabling systems to perform tests by learning and reproducing tests and thus improving efficiency and accuracy over the time [3] to reach systems capable of imitating human PT experts in performing intelligent and automated penetration [20].

In practical terms, incorporating ML in any PT system will at least reduce recurrent human errors due to tiredness, omission, and pressure. It will also boost system performance when performing different tests. ML-based automation will also relieve network congestion and downtime by reducing the number of tests by performing only relevant tests and doing that outside the regular business or office hours and thus avoid any type of assets' availability issues. Three core issues are expected to arise in a ML-based PT system. First, acquiring and generalising experience-use knowledge gained during the learning process for an optimal future use in similar situations. The second issue is adapting to the very particular context of learning that fulfil sequential decisions making with the rewarding process and approach (both automated and human-expert rewarding contexts). Finally, the exploration-exploitation trade-off aims to guarantee the best possible results within a reasonable use of resources. Furthermore, the training of such system will requires that the learning module be open

and able to interact directly with the expert to deal with complex situations by offering indications and suggestions which can be accepted or rejected by the PT expert.

1.2. Paper Outline

In this paper, we are mainly concerned with the network perspective of the PT practice and we will focus solely on the application of ML and specifically RL technique to the PT practice to make it intelligent and efficient. The proposed solution can be extended to other types of PT such as Web and application testing by introducing some changes in the core program. This paper will start with a brief background on PT practice and highlight the fact that ML is so crucial to today's PT frameworks and systems. The second section reviews relevant literature and surveys related works especially ones tackling the uses of AI and ML in the PT practice and the limits and drawbacks of current PT. The third section will briefly introduce the RL approach and justify this choice for the PT context along with presenting the first version of the proposed model and its different components. Section 4 describes the proposed system called Intelligent Automated Penetration Testing System (IAPTS), the adopted learning approach and the modeling of PT as RL problem. Section 5 will describe IAPTS in more detail as well as the performed tests and the obtained results within a specific context and test-bed network. Finally we analyze and discuss the obtained results and make the relevant conclusions along with highlighting future research works.

2. Literature Review

This work is rooted in a long line of applied research works on automating and optimising offensive cyber-security auditing processes and systems especially vulnerability assessment (VA) and PT [2, 4, 10]. Among the most significant contributions in this regard, we present here a summary of the previously completed research with a special focus on the adopted approaches and the contributions.

Initially, researchers were interested in the planning phase. Some works were implemented within the industrial PT systems and frameworks while others remained stuck at research ideas level [7-9]. As PT automation and enhancement domain is situated between both cyber-security and AI research fields, several axes of research were dressed started with the consideration of attack graphs and progressed throughout different research fields and methodologies of Automated Planning consequently sub-area of AI. Early research focused on modelling penetration as attack graphs and decision trees reflecting the view of PT practice as sequential decision making [4]. Practically, most of the works were more relevant to vulnerabilities assessment than to PT and among the most significant contributions in this regard, we present in this literature review section a summary of the previously completed research with a special focus on the adopted approaches and the contributions. For the purpose of clarity, we start by dressing the full picture of the research in this field and we proceed later into dividing the research axes by type, methodology, and approach [6].

2.1. Previous works on PT automation

Automation is an obvious approach to adopt for PT tasks when the objective is to produce highly-efficient PT systems. Nonetheless, automating all the whole process of testing including the versatile tasks and sub-tasks for each phase is challenging and often fails to reach the objective if done in inappropriate way notably the use of automated tools and systems which blindly perform all the possible and available tests without any optimisation or pre-processing [6-7]. The automated systems require the permanent control of a human PT expert and often fail to produce acceptable results in medium and large assets context because of the significant number of operations required to cover the entire network [8-11]. In addition to the required time which surpass realistic duration of tests, more others issues are created by automation such as the generated traffic (network congestion) and the high number of false positives alerts triggered on the asset defence solution such as IDPSs and FWs. giving what has been said, PT blind automation approach use was limited to a small network

and some medium size network with the use of customized scripts which are inconvenient requiring substantial effort as well [3-4].

Early research focused on improving PT system by optimising the planning phase which was modeled as attack graphs or decision trees problem which reflect the nature of PT practice as sequential decision making. Most of the works were nonetheless relevant to vulnerabilities assessment (VA) rather than PT because of the static nature of the proposed approach and its limitation to planning phase [6-8]. Amongst the most significant contributions, we find the modeling of VA as attack graphs in form of atomic components (actions), pre-condition and post-condition to narrow the targeted vulnerability [11] but this approach was more an application of classical planning methods in order to find the best attack graph. Further similar works were carried out on automating planning of PT tasks but alike blind automation did not address the problem of enhancing performances and only covered the planning phase of PT practice [3,12,16].

Nevertheless, a remarkable work on optimisation was introduced in [4] by modelling PT as Planning Domain Definition Language (PDDL) which for the first time accounted for attacking and post-attacking phases of PT in addition to the flexibility offered by the solution which enabled integration with some PT systems [4]. The proposed solution generates different type of attack plans (single and multi-paths) for real world PT scenarios which is then directly implemented within the attacking and exploiting system and executed in the due course along with interacting with information gathering tools for transforming the information acquired during that phase into input to a planning problem to be solved separately and then used by the attacking system for the purpose of optimisation. the only drawback of this approach was the scalability which was fatal as it was only limited to small and medium size networks [6].

AI was also considered to improve PT practice in some research [5-9] but most of the proposed modelling approach failed to deal with the persisting uncertainty in PT practice and especially the lack of accurate and complete knowledge about the assessed systems. An exception was the use of ML algorithms within a professional PT and VA system called Core-Impact in which researcher modelled PT planning phase as a partially observable Markov decision process (POMDP) which was then solved using external POMDP solver to determine the best testing plan in form for attack vectors. However, the proposed model itself is questionable as it obviously fails to model the full PT practice and thus can not cover the remaining testing phases and tasks especially the vulnerability assessment, testing and pivoting phases reputed to be highly interactive, sequential and non-standard compared with the planning and information gathering phases.

2.2. Drawbacks and limits of the current PT practice

In this subsection, we will present an overview of the domain of PT and the automation of the practice along with highlighting the limitations of the current (existing) automated frameworks, systems, and tools in dealing with the real-world situation. Penetration testing often involves routine and repetitive tasks which make it particularly slow on large networks. These tasks are unfortunately crucial for the practice and cannot just be dismissed although much of this routine can be automated. Although the proposed solutions were in theory very relevant and seemed to solve the problem, the PT practice demonstrated that the brought improvements were not enough to solve the core issue in the practice which time and resources. Some solutions were on the other hand, fundamentally unfit and inadequate for PT context. It is obvious that human capabilities and performance are limited when it comes to large and complex tasks compared with a machine especially with nowadays computing power.

The average penetration tester can spend days or weeks in testing a medium-size LAN (we are concerned here by comprehensive testing when the entire network is covered). In addition to the time and effort allocated, a considerable amount of systems downtime will be accounted as result of the performed tasks. The first two points will be added to the poor performances in term of results quality and accuracy including error and omission which could be crucial resulting from the fact that

human makes mistakes, change opinion and get bored. The penetration testing automation (automated systems and tools) were therefore presented as the magic solution to the named issues. A fully or even semi-automated solution was thus developed aiming to reduce human labor engaged in the testing, save time, increase testing coverage and testing frequency and allowing Broader tests by attempting more possibilities. The proposed solution was very diverse in term of adopted approach when some relied on automated planning (phase 1) by generating automatically attack plan (named attack graph) and then executing the attack in an automated manner. Others solution were more creative and attempted to mimic the whole process and make it automated so the system can carry out complex (chained) penetration testing tasks following different attack vector and use more exploits.

Cyber security research community start questioning the limits of the existing PT systems, frameworks and tools which are expected to become more automated and perform most of PT tasks with little or without human intervention and especially during the first 2 phases of PT; information gathering and vulnerabilities' discovering. Organisations with constant need to internal security auditing are, on the other hand, interested in more efficient PT systems which are fully automated and optimised to perform basics and repetitive PT tasks without human intervention and therefore alighting PT experts from that burden and dedicate them to more advanced tasks such testing advanced, complex and non-common attacks [4-6]. Nonetheless, researchers were struggling with the automation as PT practice is a complicate process which human barely master and therefore designing a machine that replace PT experts in conducting tests is a challenging work giving the multi-phases nature of PT practice with high-dependency between the different phases and tasks. alongside to the complexity of the PT practice, the information handled is another major issue as PT reconnaissance and information gathering phase usually produce incomplete profile of the assessed system and fail to yield a complete knowledge and leave a certain amount of residual uncertainty, this issue is often dealt with by expert by repeating some tasks, changing approach or simply making assumptions and continuing the tests.

On the top of the classic complexity associated to the PT, modern attacks upgraded their capabilities by adopting evasive technique and complex attacking path that allow them to evade network and systems defences. skilled attackers would usually seek to achieve their goals through the exploitation of a series of vulnerabilities as part of a chain of sub-attack which enable them to can take advantage of hidden (non-obvious) and composite vulnerabilities (composed of a chain of harmless flaws when together become an exploitable vulnerability) in networks. Each part of the infrastructure or systems may be approved to be secured when considered alone, but it/their combination and interaction can often provide a pathway for an opportunistic attacker. The ability to detect and analyze the interaction of multiple potential vulnerabilities on a system or network of systems leads to a better understanding of the overall vulnerability of the assessed system [6]. Finally, PT output data is a crucial issue because it is currently not used properly during retesting or future tasks and simply discarded after the PT report generation. In cyber security context, only few security configurations and systems' architecture change over short and medium term and therefore most of the previous tests output remain applicable when a re-testing is required and this particular problem constitute one of the key motivation of our research.

2.3. *Motivation and Contribution*

As a matter of fact, complexity is the worst enemy of control and thus security, computers networks do not constitute an exception to this unanimous rule. During the last decade, protecting and defending networks and critical digital assets from cyber threats required the security professional to consider less classic approach (avoiding the trap of bolting on more and more security layers and policies) and they turned their attention toward the offensive security. As with the real advances in technology and thus cyber-criminality, cyber-security researchers were confronted with the need of an intelligent PT system and framework to support human expert into dealing with high-demand on PT and the associated complexity and risks by allowing systems to take over human and conduct some or all of the PT tasks notably reconnaissances, information gathering, vulnerabilities assessment and

exploiting and therefore leave experts focusing on more complex issues such as post-exploiting and testing complex attacks.

Giving what the aforementioned facts about PT practice, no other technique or approach rather than ML seem to be fit to answer to our problem. In fact, several AI techniques were initially considered and following a comprehensive suitability research only Machine Learning through Reinforcement Learning was selected as the most prosperous option to allow an automated PT system to behave like real tester in term of operation mode and gain gradually the skills along with practice and thus gathering information, assessing and exploiting in an intelligent and optimised manner allowing the the discovery of all relevant and unlikely to be detected vulnerabilities and attacks to be tested along with pivoting between different asset to mimic the work of the real hackers. This research comes to bridge the gap in the current PT practice and will aim to resolve the following issues:

- Reducing the cost of systematic testing and regular re-testing due to human labor cost,
- Reduce the impact on the assessed Network notably the security exposure, performances and downtime during the testing ,
- Alight human experts from the boring tasks repeatability during test and assign them to more challenging tasks,
- Dealing more effectively with cyber threats' high emergence and fast changing rate (Short Lived Patterns) by allowing flexibility and adaptability,
- Perform more broad tests by covering a wide variety of attack vectors and also consider complex and evasive attacking paths which are hard to identify and investigate for human testers,

To sum up, cyber hackers seeks to achieve specific goals through the exploitation of a series of vulnerabilities as part of a chain of sub-attacks. Skilled attackers can take advantage of hidden (non-obvious) and complex vulnerabilities (composed of a chain of harmless flaws) in a network infrastructure or segment. Each part of the infrastructure or systems may be approved to be secured when considered alone, but the combination or the interaction can result in opening a pathway for an attacker. for this reason, the ability to assess and analyse the interaction of multiple potential vulnerabilities on a system or network is becoming crucial in PT practice

3. Reinforcement Learning Approach

Cyber security system often categorised under two types; expert-driven or automated system utilising unsupervised machine learning [5]. Expert-driven systems such AVs, FWs, IDPSs and SIEMs rely on security experts' input and usually lead to high rates of errors until Reinforcement Learning (RL) techniques were used to give existence to more goal-directed learning systems that provide autonomous or semi-autonomous decision making which accurately reflect real-world context of cyber-security and especially offensive security domain such as vulnerabilities assessment and PT context [12]. The main reason behind our choice of RL are:

- Effective autonomous learning and improving by allowing constant interaction with the environment.
- Rewarding based learning and existing flexible rewarding schemes which might be delayed to enable RL agent to maximize a long-term goal.
- Richness of the RL environment which help in capturing all major characteristics of PT including the uncertainty and complexity.

As shown in Fig. 2, RL allows an agent to learn from its own behaviour within the RL environment by exploring it and learning how to act basing on rewards received from performing actions undertaken. This decision policy can be learned once and for all, or be improved or adapted if better results is encountered in the future. If the problem is appropriately modelled, some RL algorithms can converge to the global optimum which is the ideal behaviour that maximises the the overall reward.

RL learning scheme exclude the need for a significant intervention from human who is expert in the domain of application. In addition, RL implication will mean that less time is allocated for the

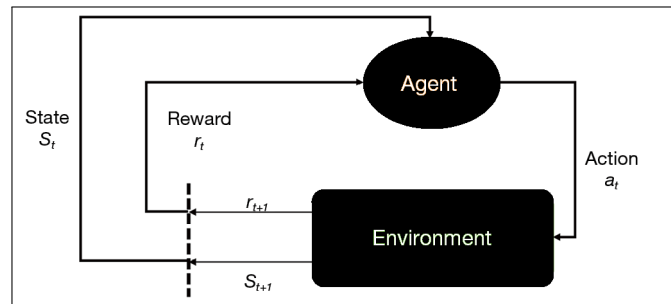


Figure 2. RL agent observes the state of the environment $x(t)$ at time t , selects an action $a(t)$ based on its action selection policy, and transitions to state $x(t+1)$ at time $t+1$ and receives a reward $r(t+1)$. Over time, the agent learns to pursue actions that lead to the greatest cumulative reward [5].

learning and customisation as it is the case with ML and expert systems (ES) respectively. In addition to what has been said about the suitability of RL for enhancing the automation of PT solutions, RL branch is a very active domain of research and several new algorithms have been introduced recently along with some very efficient toolboxes and implementations with the ability of solving complex RL problem under constrain resources and producing great results [19].

3.0.1. Towards a POMDP modelling of PT

In PT, an attack is a set of tasks which are launched and executed, manually by a human tester or automatically by a PT platform, following a certain order in order to fulfill a goal or reach an objective. Depending on the context the goal can be predefined or unknown and also can vary throughout the attack. The goal (or also known as the objective) of the attack is known within the PT community as the target which can be either logical or physical entity. Often, the target is a computer (physical or virtual machine with an OS and running applications) or a computer network or some information hosted on a computer such as files, DBs or web-servers. The attack target can also switch during an attack if a more valuable or easily exploitable target is identified to serve as a pivot later on. Furthermore, it is also common that an attack has no specific target with the example of script kiddie hacker running a set of exploit against all reachable machines regardless relevance in order to find one or more vulnerable to that specific attack [20].

The starting point for this research is an automated PT system which lack for efficiency and optimization which in term of number of covered tests and the consumed resources and time as any PT test should not last forever and consume an excessive amount into performing or exploring irrelevant tests along with ensuring that no threat is ignored or underestimated. Therefore, the aim is developing RL-led autonomous PT system which utilise RL and other techniques at different levels of the practice to improve performance, efficiency, testing coverage and reliability [20]

3.1. POMDP Solving Algorithm

RL algorithms are methods for solving real-world problems modeled in form of MDPs or POMDPs which usually involve complex and sequences of decisions in which each decision affects what opportunities are available later and running for sequences long-term goals. In this work, we are not concerned with the development of improvement of a new RL solving algorithm or methods, but only with finding the appropriate algorithm relevant to our problem and which produce acceptable results.

When it comes to solving a large and complex RL problem is the often complicated and therefore an adequate choices of the solving algorithms and approach should be made. Therefore, for solving the PT POMDP complex environment the IAPTS should rely on different solving algorithms rather than simply one, in fact, depending on the context IAPTS will adapt to utilise to most adequate solving approach. Furthermore, the choice of different algorithms is justified by the constraints IAPTS may face in term of the available resources (time, memory and computational) which make the use of one solving

algorithm challenging and thus adopting a flexible approach where the accuracy is often sacrificed to acceptance. Finally, it is important to remind that large environment can also cause challenges to solving algorithm especially when dealing with a large number of transitions and observations or opting for a static rewarding schemes [20-24].

most of RL solving algorithms fall under to major categories; the reward (value) oriented solving and policy search solving. The reward approach allows an RL agent evolving within the environment to select the sequences of actions that lead to maximising the overall received reward or minimise received sanctions in the long term run and not only in the immediate future, this approach aims to dress an optimised and comprehensive rewarding function which rely on the atomic reward values associated with the RL environment to determine and an optimal (best possible) rewarding scheme (function) for each transition and observation. In term of efficiency, this solving approach is often complex and time consuming with several cases of an infinite horizon if the problem representation is not enough consistent and optimised. The second approach, namely policy search seeks to construct a decision policy graph which is in practice done by learning the internal state/action mapping of the environment and uses direct search method for identifying policies that maximizes long term reward, optimal policy is reached when all the states and all the actions are tried and allocated a sufficient amount of time to find the best possible associated policies. in this research we opted for the use of both reward-optimisation and policy-search approaches. Nonetheless, for the purpose of implementing policy-search algorithms we found that it is useful to include both On-policy and Off-policy implementation to allow a better evaluation in term of policies quality. The IAPTS POMDP solving module will use a powerful off-the-shelf POMDP-solver allowing the use of different solving approaches and state-of-the-art algorithm [19] to allow the exclusion of all external factor when it comes to evaluating different solving algorithms performances. Initially the following algorithm were shortlisted:

3.2. *PERSEUS algorithm*

PERSEUS is a randomized point-based Value Iteration for POMDPs proposed by [5] performs approximate value backup stages to ensure that, in each stage, the value of each point in the belief set is improved. The strength of this algorithm is its capacity of searching through the space of stochastic finite-state by performing policy-iteration alongside to the single backup which improve the value of the belief points. Perseus backs up also a random basis by selecting a subset of points within the belief set which are enough to improve the value of each belief point in the global set. In practice, PERSEUS is reputed to be very efficient because of the approximate solving nature and is the best candidate for solving large size POMDP problems as it operates on a large belief set sampled by simulating decisions sequences from the belief space leading to significant acceleration in the solving process.

3.3. *GIP algorithm*

GIP (generalized incremental pruning) is a variant of POMDP exact solving algorithm family relying on incremental pruning. GIP algorithm replaces the LPs that were used in several exact POMDP solution methods to check for dominating vectors. GIP is mainly based on a Benders decomposition and uses only a small fraction of the constraints in the original LP. GIP was proven in [19] that it outperforms commonly used vector pruning algorithms for POMDPs and it reduces significantly the overall running time and memory usage especially in large POMDP environment context. The latest version of GIP is, to the best of our knowledge, the fastest optimal pruning-based POMDP [21].

3.4. *PEGASUS algorithm*

PEGASUS is policy-search algorithm dedicated to solving large MDPs and POMDPs and was initially proposed by [13] and adopts a different approach to the problem of searching a space of policies given a predefined model as any MDP or POMDP is first, transformed into an equivalent POMDP in which all state transitions (given the current state and action) are deterministic and thus reducing

the general problem of policy search to one in which only POMDPs with deterministic transitions are considered. Later, an estimation value of all policies is calculated making the Policy-search simply performed by searching for a policy with high estimated value. This algorithm has already demonstrated huge potential as it produces a polynomial rather than exponential dependence on the horizon time making it an ideal candidate to the penetration testing POMDP solving.

3.5. Other candidates

In addition to the candidates, other RL algorithms will be considered such as Backwards Induction and Finite Grid, this last is instance of point-based value iteration (PBVI) and will be mainly utilized in determining the shortest attack-path when more than one policy is found. Some of the proposed algorithms are already part of the POMDP-solver software and an optimized implementation is provided by the contributor and constantly improved over the versions. Nonetheless, some algorithm was implemented and integrated for the sake of bench-marking. Initially, and as the research focus was to dress a high-quality POMDP model representation for the PT practice bridging the gap between the theoretical research and real-world situation facing the industry professional, the use of such “ready solution” was highly recommended and was hopeful in advancing the research and also for the impact of the results obtained.

3.6. POMDP solving choices

PT is a complex practice in which the targets can be known or unknown, global or local, simple or composite and each phase is a sequence of non-standard tasks in which the order is a crucial factor. Therefore, the IAPTS should reflect to the best the real-world domain of PT and RL approach here is meant to address the kind of learning and decision-making problems that allow the PT system to capture, reproduce and store expertise in the whole PT tasks and sub-tasks relying on well established RL solving algorithms elected to be the fit to PT context and produces acceptable results [23-24]. The PT practice is thus represented as POMDP environment and serve as an input to the off-the-shield solver in which a decision-making agent will be exploring its environment to aiming to maximize the cumulative reward it receives or finding the optimal policies graph (PGs) through the RL agent which perceives the environment and solve the problem by estimating the value function to to dress the best decision policies or rewarding function [20].

4. IAPTS Design and Functioning

The proposed Intelligent Automated Penetration Testing System (IAPTS) functioning diagram is detailed in Figure 3. Python scripts were developed to perform the pre-processing from the raw data and then the produced results is used into optimising the representation of the PT domain in form of POMDP problem. The IAPTS knowledge base (memory) will be initially handled manually and a human PT expert will decide on the storage of the obtained results (policies extracted after applying the generalization) along with the management of tasks related to expertise extracting and storing. In other words, the extracted expertise will be performed manually until the IAPTS reach a pseudo-maturity state in which it will be in charge of capturing, assessing and storing the expertise will be implemented and embedded within the IAPTS expertise memory. The projected IAPTS will be an independent module that can be embedded with the industrial PT framework. the current version of IAPTS is associated with Metasploit Framework (MSF) as external module communicating via customised python scripts with MSFRPC API. The purpose of such configuration is to avoid modifying the core component of the MSF and allowing us, for research purpose, to measure the IAPTS performances away from the PT framework.

4.1. IAPTS operative modes

IAPTS will evolve through different levels of automation and intelligence to reach the pseudo-maturity level in which it should be able to perform an entire PT on networks. Overall,

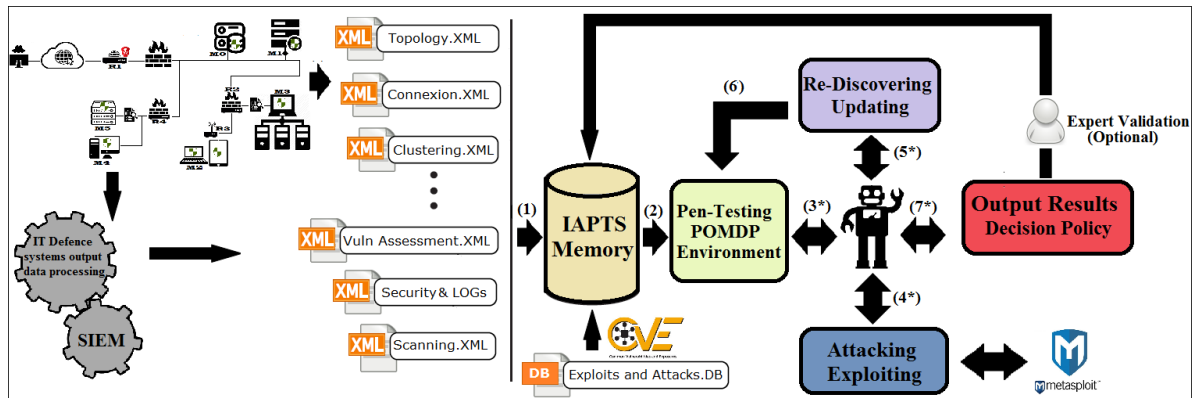


Figure 3. IAPTS functional diagram.

IAPTS can operate in four different levels which are dictated by the development of the system knowledge base in term of captured and generalised expertise as follow:

1. Fully autonomous; IAPTS entirely in control of testing after achieving maturity so it can perform PT tasks in the same way that human expert will do with some minor issues that will be reported for expert review.
2. Partially autonomous; the most common mode of IAPTS and reflect first weeks or months of professional use when IAPTS will be performing tests under constant and continuous supervision of a high-calibre PT expert.
3. Decision-making assistant; IAPTS will shadow human expert and assist him/her by providing pinpoint decision on scenarios identical to those saved into the expertise base and thus alight tester from repetitive tasks.
4. Expertise building; IAPTS running in the background while human tester perform tests and capture the decisions made in form of expertise and proceed to the generalisation and of the extracted experience and build the expertise base for future use.

4.2. From PT to a Reinforcement Learning Problem

We present here an improved version of the modelling of PT practice as a POMDP problem which constitute the core module of IAPTS. for simplicity purpose, we use an illustrative example to introduce the different steps towards the representation of a PT domain in form of POMDP problems. In the context of PT, we believe that there is no need to represent the entire network topology and security configurations in the RL environment but only representing specific data judged relevant from the PT point of view and thus alighting the RL environment [20]. The RL representation will capture the following information about the assessed network: machines and networking equipment architecture, connectivity and reachability, network defence and security configurations. The aforementioned information will be used to dress a PT-style view of the network without encumbering the RL environment and impact the performances. In addition, we used pre-processing output relevant data to be included within the environment or to serve in enhancing RL learning algorithms to acquire such as proxy server logs, web-server logs, database logs, routing device logs, apps and other security logs.

4.3. representing Network PT as RL Environment

We describe here the process of elaborating an RL environment starting from a given PT example. The overall extraction and elaboration process is explained, in mirror with the PT diagram, in Figure 6 in which we build upon the IAPTS logic into converting PT domain into POMDP representation along with respecting real world PT and adopting the same approach into the elaboration of the POMDP environment sections. noting that all the sections are dynamic and allow high frequency changes as

the PT progress and information are updated or upgraded. The following are the different components of the RL POMDP environment:

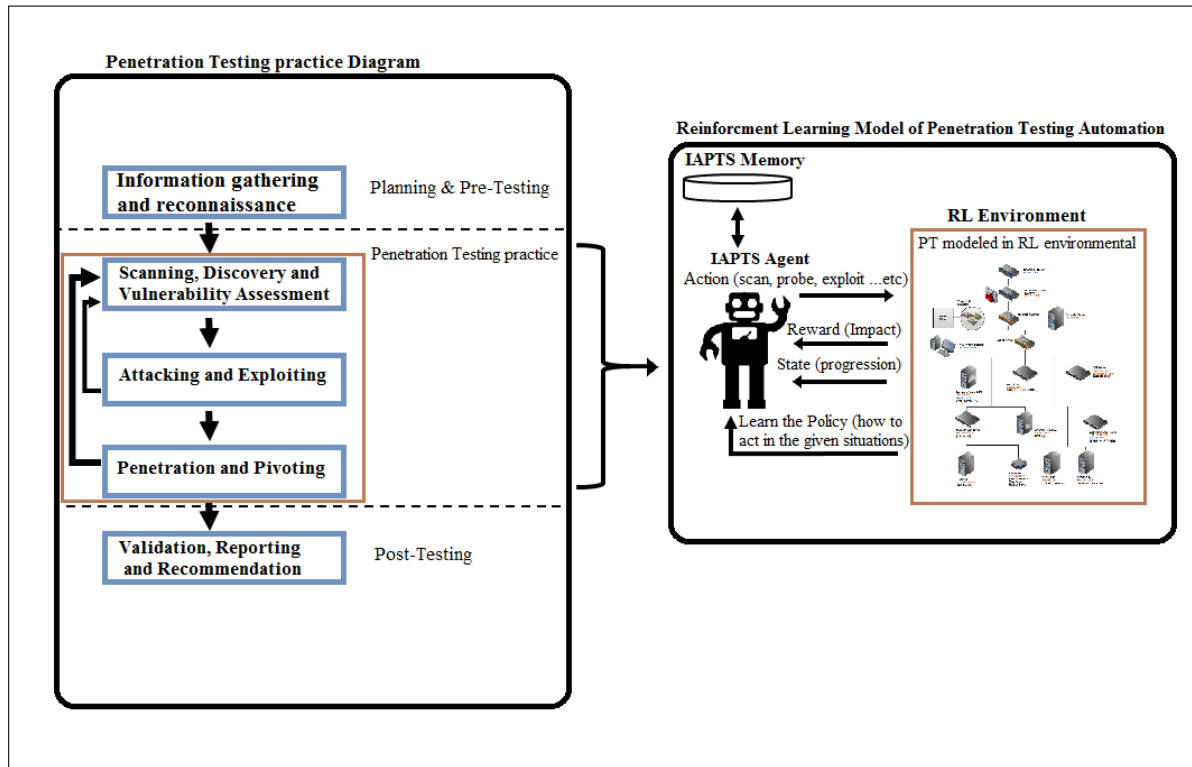


Figure 4. IAPTS modelling of PT as RL problem diagram.

State space: contain all relevant information, from PT expert view, about the assessed network. It will include information about any software or hardware machine including virtual and networking equipment that run an OS. the information are OS parameters, port, services and applications, patches in addition to relevant security and connectivity information. These information are represented in POMDP language using a special notation that aims to minimise the size of the file but remain concise, clear and precise. In practice most of the Action space is dressed at early level as modern PT rely on initial information gathered during the first phases. nonetheless, some information will remain missing or not accurate enough and thus represented in a probabilistic way after being enhanced by information coming from the pre-processed output to avoid redundant or useless representations [20]. Any machine or device within the network will be assigned a number "i" and will be represented as M_i or R_i and the remaining associated information are represented in, but not limited to, the following way M_i -OS1-Port80-ServiceABC or R_i -OS2-Port443-SerciceXXX. These information will be continuously updated as the discovering and scanning tasks progress to confirm previous probabilistic information or to add a new one. Furthermore, modern network Routers are more than simple transmission equipment, in fact they can run Operating Systems and embed one or more security isolation and protection mechanisms notably FWs, AVs, IDPSs, VLANs and others. Following this logic, network and firewalls can be considered as machines (running OS and thus having vulnerabilities) or just security isolation boarder for clustering purposes detailed later.

In addition to the machine and devices information, state space will include information about the networking and security configuration of the assessed network such as connectivity, security isolation (sub-net, virtual LAN) and defence restrictions. the purpose of such representation is to enhance and optimise the input for the POMDP solving algorithms so a better RL environment is represented. The following example summarise the information captured about two machines M_i and M_j as M_i - M_j -TCP-SSH-0". Only relevant security and networking configurations information are

considered and machines that belong to the same segment and have the same protection should be represented together then we represent other segments' machines.

Action space: POMDP model action are an exact reflection of the PT actions performed by testers and thus en-globe all PT tasks and sub-tasks following a certain notation. as with any RL problem the number of action is known, static and limited and PT does not fall out of this logic and we include in this space as variety of Pt related actions such as Probe, Detect, Connect, Scan, Fingerprint, VulnAssess, Exploit, PrivEscal, Pivot in addition to some generic action that will be used for control purpose by RL agent.) that the expert can perform is huge and cannot be totally represented within the RL action space such as Terminate, Repeat and Give-Up and others as detailed previously in [20]. Furthermore, as in PT domain successful or failed action might require further or repeating actions we defined some additional actions in order to differentiate between the original action and the others action. in practice, the purpose of such re[presentation is to deal with the special and complex scenarios notably:

- a failed action to fully (root) control a machine that lead to further action attempting user-session or escalate privileges or switching to other attack paths;
- dealing with action relying on uncertain information and fail because of the assumption made and require further actions when additional information become available and might be successful;
- actions prevented or stopped by security defence (Ws and IDPSs) which may be re-attempted following different circumstances.

4.4. POMDP Transitions and Observations probabilities

in the first phase of this work, transitions and observations were uniformly sampled. nonetheless, after multiple attempts we found-out that in the particular context of PT, it is far more efficient and reasonable to use real-world data built from IAPTS past tests and enhanced by the human-expert initially meant to passively supervising the IAPTS. the data used to artificially simulate testing environment is captured and stored by IAPTS during the regular testing but is carefully inspected by the authors who will rely on their expertise to only include the adequate data and discard the rest of the data. in addition to the regular output of the past experiences, failed or incomplete testing scenarios will be of a crucial use during the retesting process. in fact, as IAPTS aims to gradually replace human expert in PT, the system should act as human in dealing with failure into performing some PT tasks or successfully carrying-out tests. similar to human IAPTS will uses an evaluation procedure to recognise that what have been done could be useful in another context or with minor amendments for the similar context. in IAPTS, we rely initially on human expert interaction to provide a feedback on the failed and incomplete testing to select and store the highly prominent ones for future use even if they ultimately failed. in term of data, IAPTS will be mainly dealing with the Policies stored into the PG file which constitute the outcome of the POMDP problem solver [18-20].

In this research, the probabilistic output of PT action (scanning, fingerprinting, exploiting) was a crucial factor we considered doing allocation the adequate probabilities for Transitions and Observations in order to mirror the real-world PT practice. therefore, we opted for a cross validated method using two well-established and standard sources respectively NIST National Vulnerability Database (CVSS) [18] and Common Vulnerability and Exploits (CVE) which constitute a reliable online catalog for all known proven vulnerabilities associated with different type of Operating systems, software and Applications. The use of such sources is motivated by the rich content, easy accessibility, regular update and the available scoring function and mechanism such as CVSSv3 and the calculation of the Probabilities associated to each transition or observation is detailed in [20].

4.5. Rewarding schema

On the other hand, IAPTS Rewarding will be twofold depending on the system maturity. In early stage, IAPTS will rely solely on the rewards allocated by the PT expert supervising it along with some default rewarding values. rewarding the performed actions will be predefined by human expert who will have to decide on the adequate reward for each action performed depending on his/her overall

sight he got on the practice, experience and testing achievements. Afterward, IAPTS will alight the human expert from the rewarding task and only request human decision on the global PG (attack policies). IAPTS reward function will be utilised and thus the reward for the performed actions will be calculated following a well-established criteria such as: reaching a terminal state; achieving a final (global) target or local goal (controlling an intermediates machines); or failing to reach any goal. The criteria for the choice of rewards will mainly be: the estimated value of the achievement, the time consumed; and the associated risk of detection as detailed in [20].

4.6. IAPTS memory, expertise management and pre-processing

This research is all about applying RL learning into medium and large LANs which subsequently mean that the projected system IAPTS will need to deal with big amount of intimation described as complex and redundant amongst the cyber-security community. modeling and representing the PT as POMDP environment is particular complex and will result in producing a huge POMDP environment and thus make it impossible to solve giving the restriction in time and computing power (memory). Therefore, an optimising and smart use of resources is required and the problem modeling is where all start. The system memory as shown in Fig. 5 is used for dynamically storing the data handled by the system such as the environment's attributes (States, Action, Observation, transition, Reward) and agent's memory (data regarding the Policy and Acquired knowledge and experiences that an agent gains by acting within the environment). In fact, the first part of this research will focus exclusively on searching the policy as an agent acts within the environment in a particular state and receives a reward from the environment. Initially, for a purpose of research facilitation, the reward value will be pre-defined by a human expert so no reward function will be used. Moreover, generalize experience output for further use – use knowledge gained in similar situations to Equip penetration testing system with “expert knowledge” will be completely done by this module in the future. in practical term, IAPTS will solve the RL problem, extract PGs and instruct MSFRPC API which will execute the testing plan and keep updating the IAPTS of the outcome on real-time base especially at vulnerability assessment and exploiting phases. this will enable IAPTS to adjust and adapt the tests as well as the the post-exploitation tasks such as pivoting or privilege escalation

In addition to the RL framework on which the system will operate, a parallel knowledge-based expert system will be implemented and constantly (with every practice) enhanced and fed. This pseudo-system will serve as RL initial belief. This system will capture details of the performed (manually) action by the human tester and also extract knowledge from the output of the information gathering phase and Security system data (Firewalls, AVs, IDPSs, SIEMs) and structure the relevant details. One can say that such system will be useless alongside with the RL system which is a legitimate interrogation. The answer will be that giving the known limits of RL in multi-dimensional state context along with the important size of the RL components, including initial belief detail will only slow down the system performance. Furthermore, the crucial information extracted from the security data will otherwise be omitted.

The only remaining issue is the human intuition (the ability to acquire knowledge without inference and/or the use of reason) which a system will not be able to substitute. Intuition provides penetration tester expert with beliefs that cannot be justified in every case and human can sometimes solve some brilliant problem without the use of any reasoning. Artificial decision making is the ultimate aims of the use of AI but still cannot model the intuition. As results, this issue will be sorted out by allowing the controlling human to interact with the system. In other words, a mechanism to obtain feedback from the expert tester (security analysts) should be utilized to overcome this issue. The feedback (along with the surrounding context) will be stored in the system memory for future use. the system memory will incorporate policies assessment and generalization features and experience-replay form previous test where the human expertise is extracted and defined as a policy automatically by the system (direct learning) along with the management of the input and output data such as the initial belief and reporting.

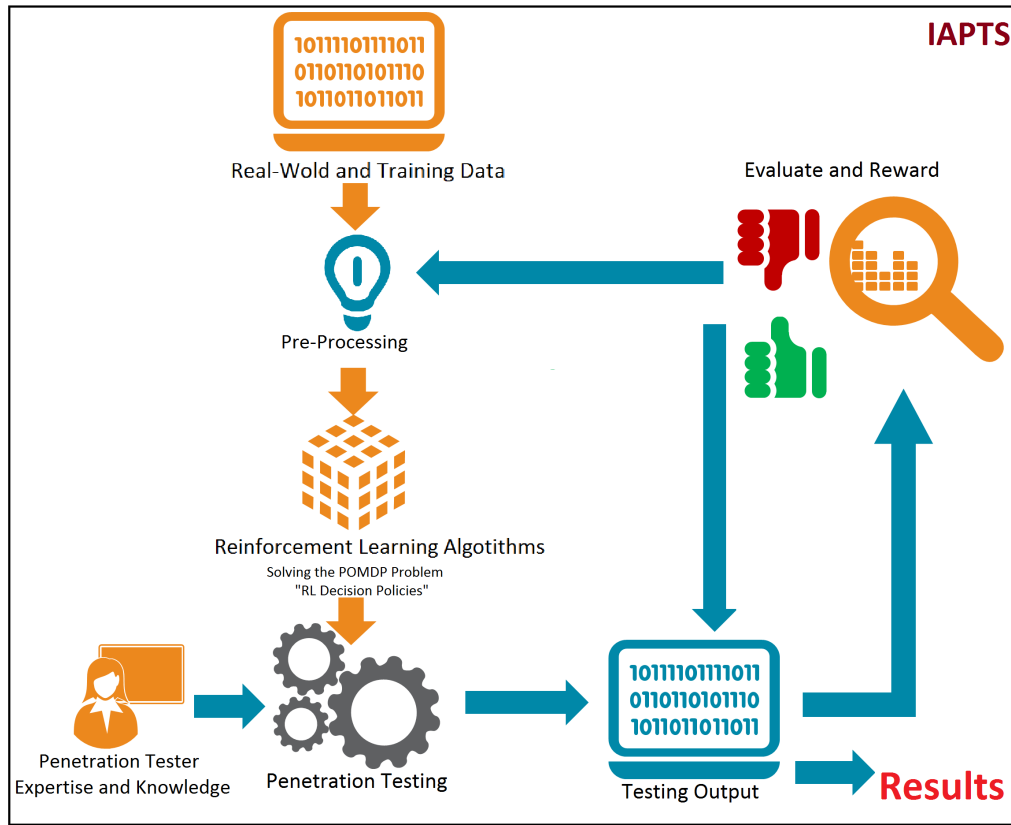


Figure 5. IAPTS learning, expertise extraction and validation procedure

Prioritised experiences' replay is an effective approach to improve the learning and thus efficiency in RL algorithms. In this work, we adopted this approach, but introduced some modifications for technical reasons, in order to enable RL algorithms to prioritise the use of certain sequences of transitions over others in order to enhance the learning of the IAPTS RL agent. In addition to selecting the most plausible and relevant policies (state-action decision sequences), we injected some other artificially construct transition sequences using information gathered from previous tests which were validated by a human PT expert. These sequences, when replayed, allow value function information to trickle down to smaller selection of POMDP transitions and observations, thereby improving solving algorithm efficiency in term of consumed time and memory. all the proposed customization were implemented within a modified version of the standard POMDP solving GIP LPsolve algorithm we called "with Initial belief".

Finally, it is important to introduce our modified GIP LPsolve algorithm which was meant to improve the performance of IAPTS and also allow the IAPTS to capture the appropriate expertise in form of decision policy) process it to make it general decision rule and store it within IAPTS memory for future use. the simplest way to illustrate the importance of the learning on the long-term PT practice by adopting a test scenario inspired from the real-world situation of re-testing the same network after some updates or upgrades. In the retesting process, one or more machine configuration will be changed but not all of the machines and therefore IAPTS will re-use already acquired PG when it comes to repeat PT on the partially modified network with the use as initial belief the output of the previous tests.

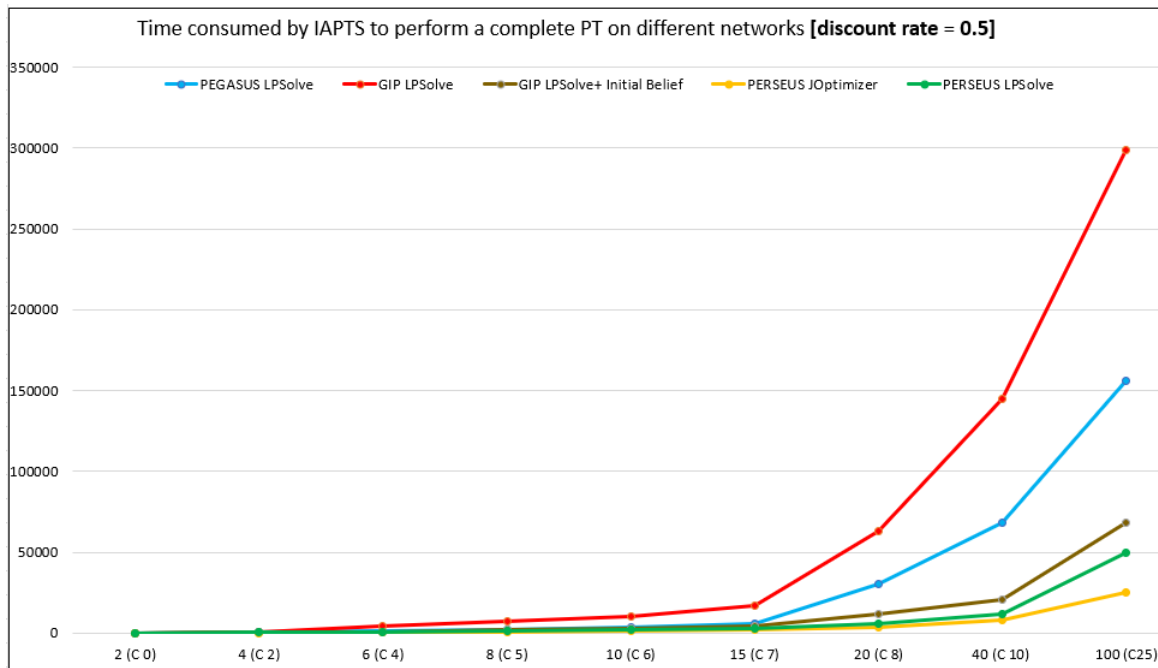


Figure 7. Time in seconds required by IAPTS to complete PT tasks on different LANs' sizes with a Discount rate of 0.5

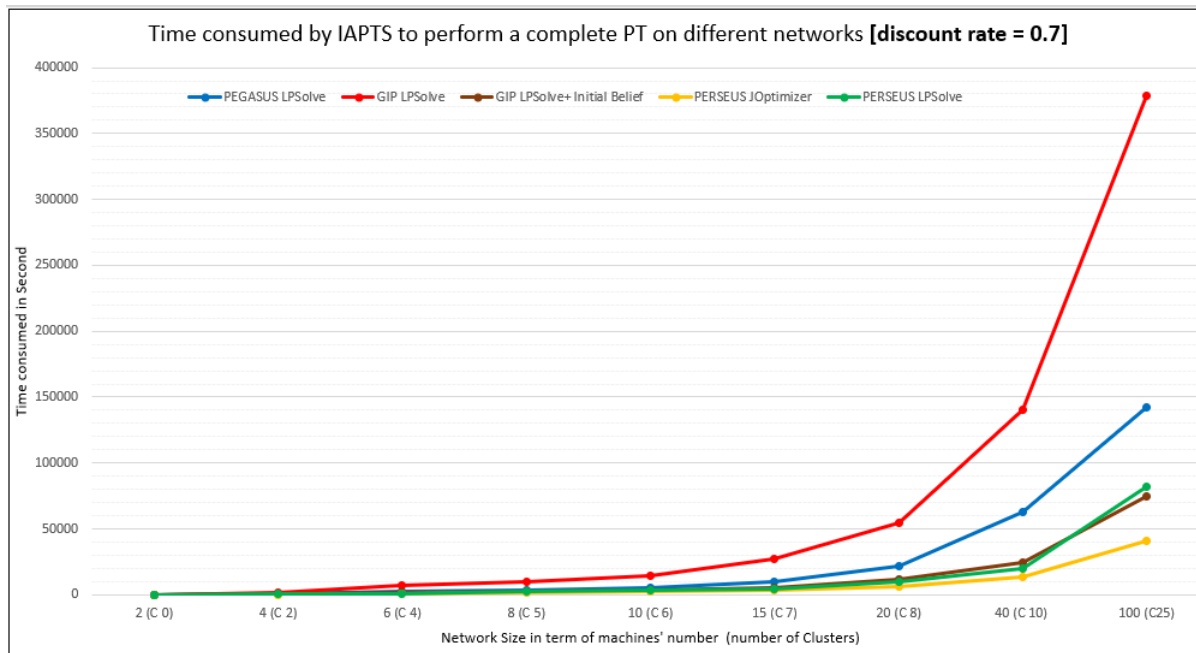


Figure 8. Time in seconds required by IAPTS to complete PT tasks on different LANs' sizes with a Discount rate of 0.7

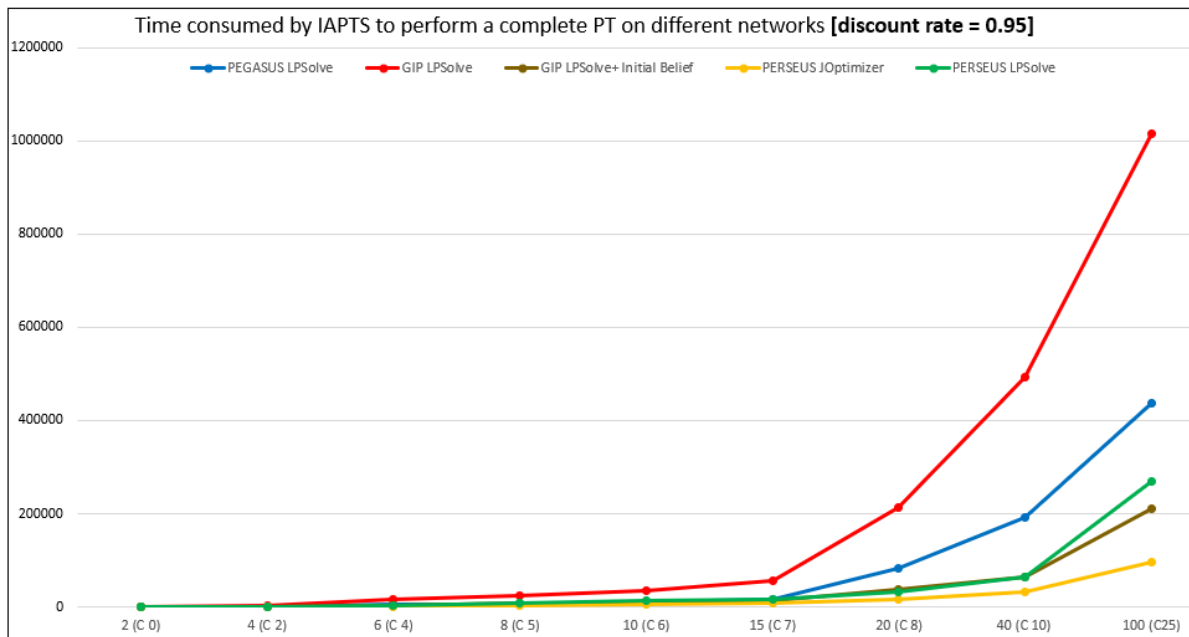


Figure 9. Time in seconds required by IAPTS to complete PT tasks on different LANs' sizes with a Discount rate of 0.95

Following the obtained results, we decided to introduce some changes within the solving algorithm and notably GIP aiming a better performance from IAPTS on short term basis. we opted for prioritized Transitions and Observations through the manipulation of the associate probabilities along with introducing some customisation into the initial beliefs sampling. the obtained results were surprisingly excellent and the new variant of GIP which we named GIP-LPSolve with Initial Belief performed much more better than the classic GIP in both time consumed and PG accuracy as shown in Figures 7, 8 and 9. Furthermore, in order to assess IAPTS performance expertise extraction and storing capabilities and the impact of performance enhancement we proceeded to the re-testing the same network with or without introducing minor or major changes to different number of machine configuration. the obtained results in the context of a 10 machines LAN were near to perfect as the performance enhancement was huge especially when re-testing the very same network as shown in Figure 10.

Finally, on the top of the overall performance enhancement and notably when using GIP LPSolve with initial belief algorithm, the quality of the produced decision policies was beyond human expertise especially in the case of 10 machines network when IAPTS highlighted two additional attack vectors which an average human PT expert would easily omit and illustrated in Figure 11.

5.3. Discussion and future works

the obtained results consolidate prior thoughts on the role of ML and specifically RL in the performance enhancement and resources-use optimization in PT. Commercial and open-source PT systems and frameworks were deigned initially to work either under human instructions or in a blindly automated manner, but both approaches fail to address the current environment in which PT practice is evolving notably the increasing size and complexity of the networks, the high number of vulnerabilities and the composite testing scenarios which mimic modern hackers operating approaches. RL revealed very efficient when used properly and IAPTS results are an additional evidence as in addition to the drastic performances enhancement comparing to an average human testers, several other positive points were noticed notably the pertinence of the produces result (acting policies) in term of relevance, coverage and accuracy. In practical term, using the adequate RL algorithm and adopting a new learning schemes enabled IAPTS to produce a very optimised attacking policies when targeting

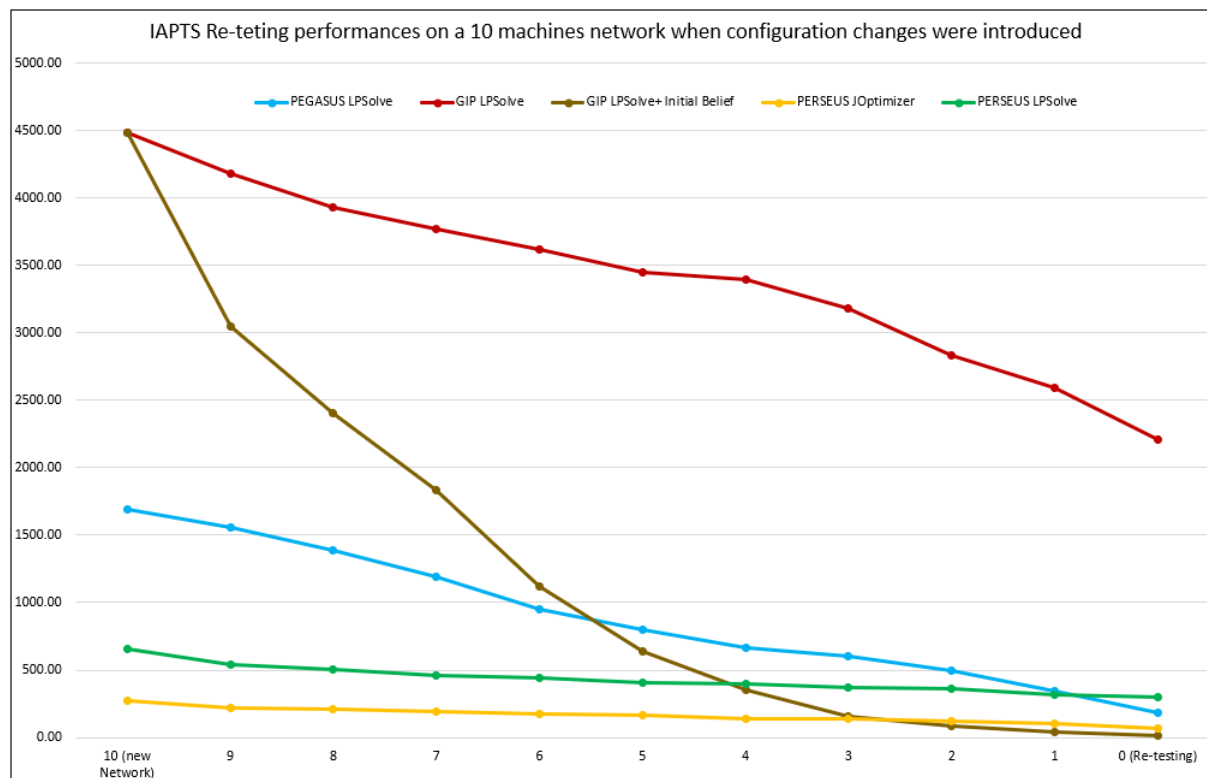


Figure 10. IAPTS re-testing performances' enhancement by algorithm on 10 Machines LAN

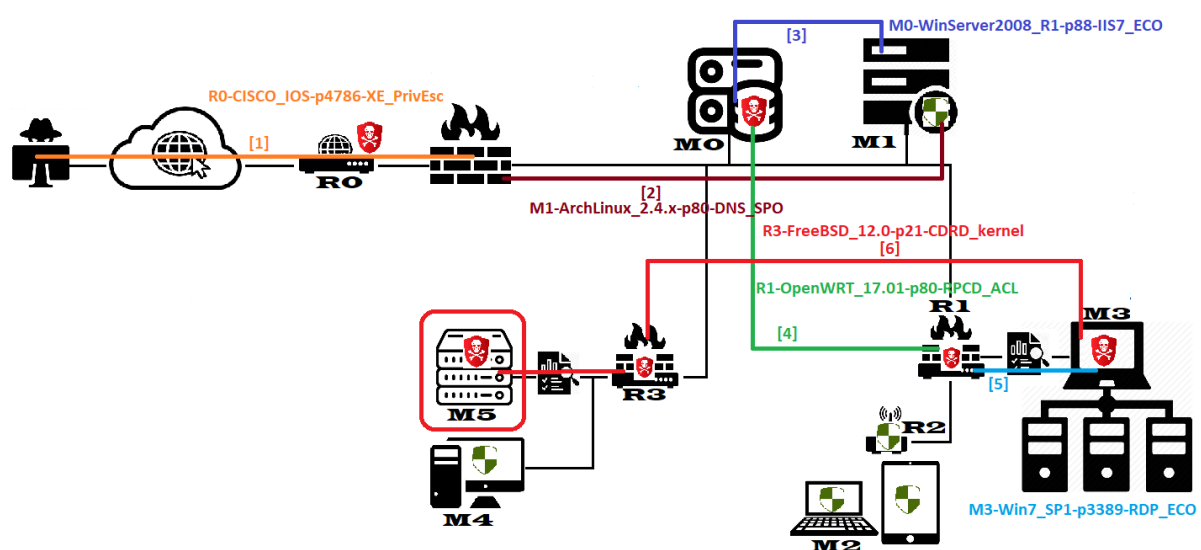


Figure 11. Example of IAPTS output PT policy translated into attacking vectors

the Machine M5 suspected to contain sensitive information and defined as the most secured machine within the test-bed network as illustrated in Fig. 11. Indeed, the produced policy is from an attacker point of view obvious but getting an automated system to opt for such attacking vectors despite being not minimal in term of cost of the exploits and consumed time is the novelty in IAPTS which is able to sacrifice simplicity for a higher objective. IAPTS exploring and large coverage capabilities was able to find a very complex and non-obvious attacking path in medium size networks where, relying on authors experience, no human tester will be tempted to adopt and possibly neglect in spite of being very relevant and constitute a possible attack path which a real hacker can chance it.

Furthermore, the proposed enhanced GIP-LPSolve which utilise a new mechanism in creating and managing POMDP initial belief was proved very efficient especially in small and medium size LANs. In fact GIP LPSolve is a variant of an exact solving RL algorithm which are often labeled as good in results quality but bad in performances, but the introduced changes in initial belief sampling and managing along with prioritising some decision sequence over others enabled the new variant to perform much more better and even outperform other RL approximate solving algorithms. On the other, the re-testing of the same network after the introduction of minor changes in few machine permitted to appreciate the full contribution of RL to PT practice by cutting drastically the consumed time and thus allowing a fast and reliable re-testing which is often the case in PT when periodic re-testing is compulsory despite the lack of any significant configuration changes within the networks systems.

Finally, we noticed that IAPTS performances on large size LANs decreases sharply and this is mainly due to the complexity which impact the size of the POMDP environments along with usage of memory during the solving of the problem. This major issue is currently being dealt with by proposing a hierarchical PT POMDP model relying on grouping several machines under the same cluster which will be detailed in future works along with improving IAPTS pre-processing.

6. Conclusions

This paper explores a novel application of reinforcement learning techniques into the offensive cyber-security domain which allows penetration testing systems and frameworks to become intelligent and autonomous and thus perform most of testing and re-testing tasks with no or little human intervention. The proposed system named IAPTS can act as a module and integrate with most of the industrial PT frameworks to improve significantly the efficiency and accuracy on medium and large networks context. The proposed modelling of PT in form of RL problem allowed the coverage of the entire PT practice and thus producing a system fit for the real-world context, the current implementation of IAPTS is integrated to the most commonly used PT frameworks called Metasploit and permitted highly efficient testing in term of consumed time, allocated resources, covered tests and accuracy of the produced results. The main drawback of IAPTS is the need of high-calibre human expert supervision during early learning phases where a human trainer will perform PT along with IAPTS and adjust the learning and veto the output of the system to ensure a good quality training by acting as rewarding provider for the RL agent actions.

The major contribution of this approach is to apply RL techniques to a real-world problem of automating and optimising PT practice and resulted into a net improvements of PT framework performances notably in terms of consumed time and covered attack-vectors as well as enhancing the produced results reliability and persistence which will lead optimistically to a PT system free from human-error. The second major contribution of the system will be the capturing the expertise of human experts without instructing it as IAPTS will rely initially the expert feedback in form of rewarding values until it reach a certain maturity. Thirdly, IAPTS will on the top of saving time and reduce human labour, increase testing coverage by attempting tests that a human expert won't be able to explore because of the frequent lack of time. Finally, IAPTS permit the re-usability of the testing output by either learning and/or capturing the expertise during test and storing it with the system memory for

future use and was proved to be very efficient on re-testing scenario (very common in PT) and nearly similar cases when the testing time and accuracy of the produced results were exceptional.

References

1. J. Creasey, and I. Glover, A guide for running an effective Penetration Testing program, <http://www.crest-approved.org>. CREST Publication , 2017.
2. N. Almubairik, G. Wills, Automated penetration testing based on a threat model. 11th International Conference for Internet Technologies and Secured Transactions, ICITST, 2016.
3. A. Applebaum, D. Miller, B. Strom, C. Korban, and R. Wol, Intelligent, automated red team emulation. 32nd Annual Conference on Computer Security Applications (ACSAC '16), 2016, pp. 363-373.
4. J. Obes, G. Richarte, and C. Sarraute, Attack planning in the real world. Journal CoRR Article, 2013, abs/1306.4044.
5. M. Spaan, Partially Observable Markov Decision Processes, Reinforcement Learning: State of the Art, Springer Verlag, 2012.
6. J. Hoffmann, Simulated penetration testing: From Dijkstra to aaTuring Test++. 25th Int. Conf. on Automated Planning and Scheduling, 2015, AAAI Press.
7. C. Sarraute, Automated attack planning. Instituto Tecnológico de Buenos-aires, Ph.D. Thesis, 2012, Argentina.
8. X. Qiu, Q. Jia, S. Wang, C.Xia and L. Shuang, Automatic generation algorithm of penetration graph in penetration testing, 19th Int. Conf. on P2P, Parallel, Grid, Cloud and Internet Computing, 2014.
9. C. Sarraute, O. Buffet and J. Hoffmann, POMDPs make better hackers: Accounting for uncertainty in penetration testing. 26th AAAI Conf. on Artificial Intel. (AAAI'12), pp. 1816–1824, July 2012.
10. C. Heintz, Artificial (intelligent) agents and active cyber defence: policy implications. 6th Int. Confe. on Cyber Conflict. NATO CCD COE Publications, 2016, Tallinn.
11. M. Backes, J. Hoffmann, R. Kunemann, P. Speicher and M. Steinmetz, Simulated penetration testing and mitigation analysis. <http://arxiv.org/abs/1705.05088>, 2017.
12. S. Jimenez, T. De-la-rosa, S. Fernandez, F. Fernandez and D. Borrajo, A review of machine learning for automated planning. The Knowledge Engineering Review, Vol. 00:0, pp. 1–24. 2009.
13. Y. Andrew and M. Jordan, PEGASUS: A policy search method for large MDPs and POMDPs. 16th Conf. on Uncertainty in Artificial Intel., 2013.
14. T. Schaul, J. Quan, I. Antonoglou and D. Silver, Prioritized experience replay, Google DeepMind. ICLR 2016.
15. K. Veeramachaneni, I. Araldo, A. Cuesta-Infante, V. Korrapati, C. Bassias and K. Li, AI2: Training a big data machine to defend. CSAIL, MIT Cambridge, 2016.
16. K. Durkota, V. Lisy, B. Bosansk and C. Kiekintveld, Optimal network security hardening using attack graph games. 24th Int. Joint Conf. on Artificial Intelligence (IJCAI-2015), 2015.
17. N. Meuleau, K. Kim, L. Kaelbling and A. Cassandra, Solving POMDPs by searching the space of finite policies. 15th Conf. on Uncertainty in Artificial Intel., 2013.
18. NIST, Computer Security Resource Center - National Vulnerability Database, <https://nvd.nist.gov>, 2018.
19. E. Walraven, and M. Spaan. Accelerated Vector Pruning for Optimal POMDP Solvers, Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, 2017.
20. M. Ghanem, and T. Chen. Reinforcement Learning for Intelligent Penetration Testing. World Conference on Smart Trends in Systems, Security and Sustainability. 2018.
21. C.Dimitrakakis, and R. Ortner. Decision Making Under Uncertainty and Reinforcement Learning. Book chapter. 2019.
22. I. Osband, D.Russo, and B. Van Roy. efficient reinforcement learning via posterior sampling. In NIPS, 2013.
23. R. Grande, T. Walsh, and J. How. Sample efficient reinforcement learning with gaussian processes. In International Conference on Machine Learning, pages 1332–1340, 2014.
24. S. Agrawal, and R. Jia. Optimistic posterior sampling for reinforcement learning: worst-case regret bounds. Annual Conference on Neural Information Processing Systems, 2017, pages 1184–1194, 2017.

© 2019 by the authors. Submitted to *Information* for possible open access publication under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).