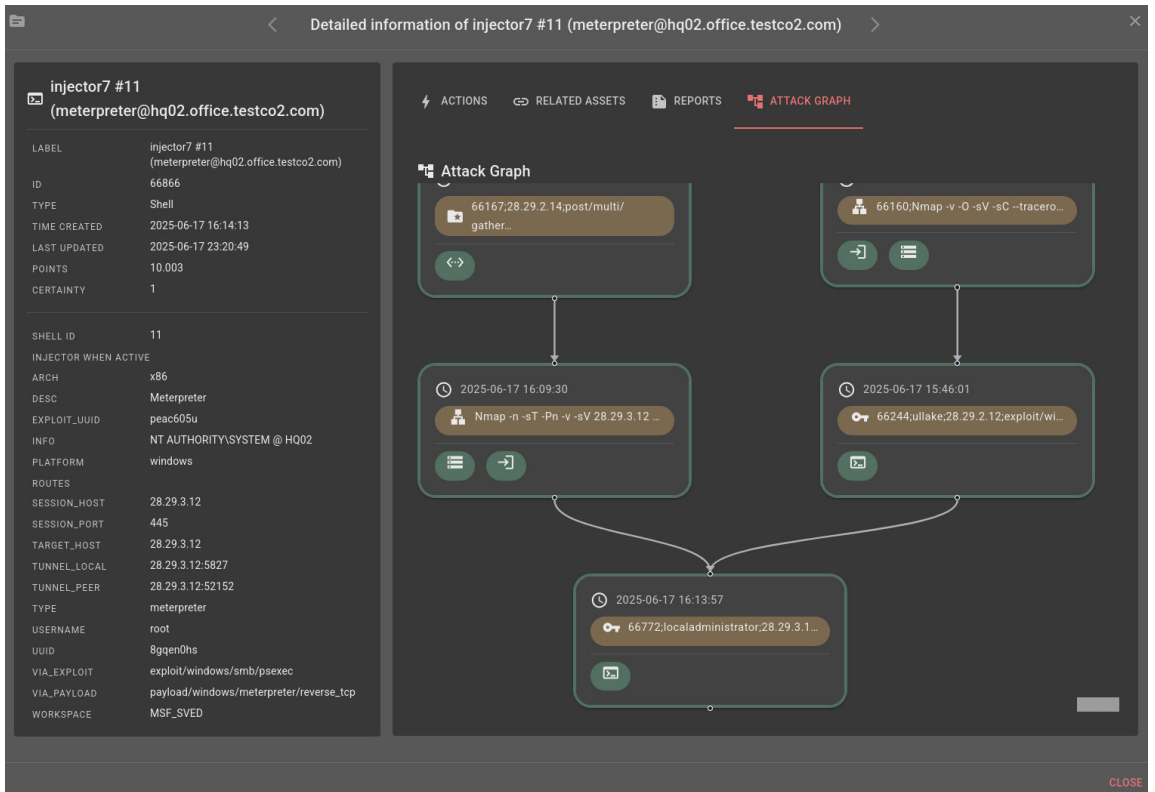


HANNES HOLM



The screenshot displays a web-based interface for analyzing a security event. The main window title is "Detailed information of injector7 #11 (meterpreter@hq02.office.testco2.com)".

injector7 #11 (meterpreter@hq02.office.testco2.com)

LABEL	injector7 #11 (meterpreter@hq02.office.testco2.com)
ID	66866
TYPE	Shell
TIME CREATED	2025-06-17 16:14:13
LAST UPDATED	2025-06-17 23:20:49
POINTS	10.003
CERTAINTY	1

SHELL ID	11
INJECTOR WHEN ACTIVE	
ARCH	x86
DESC	Meterpreter
EXPLOIT_UIID	peac605u
INFO	NT AUTHORITY\SYSTEM @ HQ02
PLATFORM	windows
ROUTES	
SESSION_HOST	28.29.3.12
SESSION_PORT	445
TARGET_HOST	28.29.3.12
TUNNEL_LOCAL	28.29.3.12:5827
TUNNEL_PEER	28.29.3.12:52152
TYPE	meterpreter
USERNAME	root
UUID	8gqen0hs
VIA_EXPLOIT	exploit/windows/smb/psexec
VIA_PAYLOAD	payload/windows/meterpreter/reverse_tcp
WORKSPACE	MSF_SVED

Attack Graph

The Attack Graph shows a sequence of actions:

- 66167;28.29.2.14;post/multi/gather... (2025-06-17 16:09:30)
- 66160;Nmap -v -O -sV -sC --tracero... (2025-06-17 15:46:01)
- 66244;ullake;28.29.2.12;exploit/wi... (2025-06-17 16:13:57)
- 66772;localadministrator;28.29.3.1...

Hannes Holm

Verktyg och teknik för CNO- övningar

Slutrapport

Titel	Verktyg och teknik för CNO-övningar: Slutrapport
Title	Tools and techniques for cyber defence exercises: Final Report
Rapportnr/Report no	FOI-R--5817--SE
Månad/Month	December
Utgivningsår/Year	2025
Antal sidor/Pages	28
ISSN	1650-1942
Uppdragsgivare/Client	Försvarsmakten
Forskningsområde	Cyberförsvar och cybersäkerhet
FoT-område	Operationer i cyberdomänen
Projektnr/Project no	E38557
Godkänd av/Approved by	Emil Hjalmarson
Ansvarig avdelning	Cyberförsvar och ledningsteknik

Bild/Cover: Hannes Holm

Detta verk är skyddat enligt lagen (1960:729) om upphovsrätt till litterära och konstnärliga verk, vilket bl.a. innebär att citering är tillåten i enlighet med vad som anges i 22 § i nämnd lag. För att använda verket på ett sätt som inte medges direkt av svensk lag krävs särskild överenskommelse.

This work is protected by the Swedish Act on Copyright in Literary and Artistic Works (1960:729). Citation is permitted in accordance with article 22 in said act. Any form of use that goes beyond what is permitted by Swedish copyright law, requires the written permission of FOI.

Sammanfattning

Utförande av cyberangrepp krävs för att uppnå en god cyberförsvarsförmåga, exempelvis för att identifiera sårbarheter, träna logganalytiker och för att anpassa intrångsdetektionssystem. Projektet *Verktyg och teknik för Computer Network Operations (CNO) övningar (VECNO)* bedrev från 2021 till 2025 forskning rörande automatisering av cyberangrepp. Forskningen operationaliserades huvudsakligen genom verktyget Lore som automatiserar cyberangreppsplaner utan krav på operatör eller förkunskap om mål.

Lore var före 2021 ett relativt enkelt verktyg som utförde ett fåtal typer av angrepp och som enbart hade använts till ett fåtal tekniska tester. Genom arbetet inom projektet har Lore vidareutvecklats på olika sätt, och idag har verktyget ett omfattande modulstöd och har använts för diverse ändamål. Lore har exempelvis använts för åtta cybersäkerhetsövningar, för forskning kring cyberförsvarsmekanismer, samt för en säkerhetstävling.

Projektet har även involverat flera utvärderingar av Lore, i synnerhet huruvida det erbjuder samma upplevelse som mänskliga inspel för logganalytiker som deltar i cybersäkerhetsövningar. Resultaten visar att den upplevda realismen var ungefär lika hög, och att inspelen var ungefär lika lärorika, oavsett om Lore eller mänskliga inspel nyttjades som källa för cyberangrepp.

Nyckelord: Cyberangrepp, cybersäkerhetsövningar, automation, artificiell intelligens

Summary

This report describes research carried out within the project *Tools and techniques for cyber defence exercises* (VECNO) from 2021 to 2025 on the topic of cyber threat automation. The research was implemented in the cyber threat automation tool Lore.

In 2021, Lore was a relatively simple tool with capabilities primarily limited to server exploits and had only been used for a few technical tests. Within VECNO, the development of Lore continued in various ways to improve its capabilities. It was employed for many different purposes, including eight cyber defence exercises, research on defence mechanisms, and a capture-the-flag exercise.

There have also been several evaluations of Lore. In particular to examine if the experiences differ for log analysts participating in cyber defence exercises depending on whether they are subjected to threats generated by Lore or by human red teams. The results indicate that the perceived realism is about the same, and that the cyber threats are regarded as equally educational, regardless of their source.

Keywords: Cyber threat emulation, cyber defence exercises, automation, artificial intelligence

Innehållsförteckning

1	Inledning	7
	1.1 Bakgrund	8
	1.2 Frågeställningar	9
2	Automatisering av test- och träningsscenarion	11
	2.1 Lore.....	12
	2.2 Enkel att använda, underhålla och vidareutveckla.....	15
	2.3 Göra lämpliga val.....	15
3	Utvärderingar av automatiserade test- och träningsscenarion. 19	
	3.1 Utvärderingar under cybersäkerhetsövningar	19
	3.2 Jämförelser med andra verktyg	21
4	Slutsatser och framtida arbete	24
5	Referenser	27

1 Inledning

Denna rapport sammanfattar det arbete som utförts inom projektet *Verktyg och teknik för Computer Network Operations (CNO) övningar (VECNO)* som utfördes från 2021 till och med 2025. En översikt av de öppna rapporter som producerats inom projektet presenteras i Tabell 1.

Tabell 1. Öppna rapporter inom VECNO.

Rapportnummer	Titel	År
FOI-R--5148--SE	Automatisering av cybersäkerhetsövningar - Vidareutveckling och evaluering av Lores beslutsprocess	2021
FOI Memo 8047	Omvärldsbevakning och statusmemo 2022	2022
FOI-R--5366--SE	Verktyg som döljer skadlig kod – En systematisk granskning	2022
FOI Memo 8085	Användartester för teknisk informationshämtning i Lore GUI	2023
FOI Memo 8217	Cyber Defence Exercise - Cybon 2022	2023
FOI Memo 8354	Utvärdering av verktyg för emulering av hotaktörer	2023
FOI-S--6759--SE	Hide My Payload: An Empirical Study of Antimalware Evasion Tools	2023
FOI-S--6760--SE	Lore a Red Team Emulation Tool	2023
FOI-S--6761--SE	Evaluation of a Red Team Automation Tool in live Cyber Defence Exercises	2023
FOI Memo 8662	Utvärdering av verktyg som emulerar hotaktörer	2024
FOI-R--5533--SE	Förstärkningsinläring för cyberangrepp	2024
FOI-S--7112--SE	Realistic and Balanced Automated Threat Emulation	2025

Syftet med projektet var att bedriva forskning kring förmågan att utföra logganalys och incidenthantering, med ett särskilt fokus på cybersäkerhetsövningar. Dyliga övningar utgör en vanlig förmågehöjande verksamhet där kan logganalytiker kan observera och hantera cyberangrepp utan risk för operativa system.

Cybersäkerhetsövningar ställer dock anspråk på ett flertal tidskrävande aktiviteter, såsom design av scenarier, miljöer och cyberangrepp. Planering och exekvering av cyberangrepp är en av de dyraste aktiviteterna vid övningar. I FOI:s övningar uppgår dessa kostnader till närmare 50 % av budgeten. Under

övningen Locked Shields¹ behövs cirka 100 personer för att utföra cyberangreppen. Av denna anledning finns det stort intresse av att automatisera dem.

Projekt VECNO innefattade huvudsakligen forskning kring automation av cyberangrepp och operationaliserades genom vidareutveckling, tillämpning och utvärdering av det automatiserade cyberangreppsverktyget Lore.

1.1 Bakgrund

Forskning kring ökad förmåga att utföra effektiv logg- och incidentanalys har bedrivits av många tidigare projekt på FOI. Bakgrunden till VECNO kan spåras till projektet *Övning och experiment för operativ förmåga i cybermiljön (ÖvExCy)* [5], som pågick mellan 2015 och 2018. ÖvExCy studerade bland annat vilka verktyg som krävs för övning och experiment av logganalysförmåga. Innan ÖvExCy utfördes cyberangrepp under cybersäkerhetsövningar som FOI anordnade i bästa fall genom enkla skripthack och i värsta fall helt manuellt. Detta skapade stora kostnader och diverse problem. I synnerhet blev övningarna mycket personberoende eftersom varje person hade sina egna skripthack. Dessutom var dokumentationen bristfällig avseende vilka angrepp som utfördes – angriparna hade helt enkelt för mycket att göra för att prioritera loggning. I förlängningen medförde detta att försvarare ofta inte kunde få svar på frågor som ”när komprometterades maskin X?” eller ”när utfördes angrepp Y?”. För att lösa detta problem utvecklades verktyget *Scanning, Vulnerabilities, Exploits And Detection (SVED)* [6] inom ÖvExCy. SVED möjliggör skriptade inspel enligt standardiserade scheman, och har använts för att automatisera inspel under de allra flesta cybersäkerhetsövningar som utförts av FOI sedan 2017.

Ett problem med rena planeringsverktyg likt SVED är dock att de kräver cybersäkerhetskunskap, särskilt om de skall kunna hantera olika möjliga förändringar eller fel som kan tänkas uppstå under en cybersäkerhetsövning. Exempelvis kanske ett angrepp som används bara fungerar ibland, eller så kanske en försvarare har möjlighet att genomföra förändringar i miljön, såsom att stänga ned processer eller starta om datorer. I praktiken finns det därför alltid människor med under exekvering av SVED-inspel för att manuellt hantera eventuella oförutsedda förändringar.

Verktyget Lore [7] skapades i projektet *Övning och Experiment för cyberförsvarsförmåga (ÖvExCND)* som pågick mellan 2018 och 2021 för att möjliggöra helt automatiserade angreppsplaner utan krav på varken operatörer eller förkunskap om målen som skall angripas. Resultatet blev ett verktyg som inte kräver någon omfattande cybersäkerhetskunskap för att exekvera. Det

¹ <https://ccdcoe.org/locked-shields/>

innebär inte heller någon större kostnad att konfigurera verktyget att utföra specifika beteenden.

Lore var under 2021 ett relativt enkelt och otestat verktyg som primärt involverade nätverkskartläggning och serverangrepp i Metasploit [8]. VECNO har fortsatt arbetet med att vidareutveckla och utvärdera Lore på olika sätt. Idag har Lore ett omfattande modulstöd och kan användas för diverse tillämpningar, såsom cybersäkerhetsövningar (se kapitel 3), anpassning av intrångsdetektionssystem [1], [2] eller för att lösa capture-the-flag-utmaningar (se kapitel 4).

1.2 Frågeställningar

VECNO var ämnat att svara på följande två forskningsfrågor:

1. Hur bör verktyg för att automatisera generering av test- och träningsscenarion inom logganalys och incidenthantering konstrueras?
2. Hur påverkar val av tekniska analysverktyg och arbetsmetod logganalysförmåga och incidenthanteringsförmåga?

Det finns flera tidskrävande processer inom skapande av test- och träningsscenarion för logganalytiker som traditionellt utförs manuellt av cybersäkerhetsexperter, och därmed skulle vara värdefulla att automatisera:

- cyberangrepp
- skapande av relevanta datormiljöer
- generering av godartat beteende, såsom simulering av legitima användare.

Vår erfarenhet är att cyberangrepp är svårast att automatisera väl, och därmed mest spännande för forskning. Av de två frågeställningarna behandlades primärt forskningsfråga (1) eftersom:

- ett svar på fråga (1) gör det enklare att svara på fråga (2)
- det fanns stora behov av automatiserade cyberangrepp för utförande av annan verksamhet inom FOI och Försvarsmakten, såsom för forskning rörande automatiserade försvarsmekanismer (t.ex. demonstratorprojektet SAC3S [1], [2] och EDF-projektet AInception²) samt för inspel under cybersäkerhetsövningar (t.ex. Safe Cyber 2020 [3] och cybersoldaternas slutövning under 2025)
- logganalys- och incidenthanteringsförmåga studerades av andra projekt inom FOI, såsom projektet *Metod och kompetens för CNO* [4] samt projektet AInception.

² <https://www.ainception.eu>

Studier av fråga (2) utfördes inom VECNO i form av empiriska tester av hur olika anti-virus upptäcker skadlig kod obfuskerad på olika sätt [10], [15]. Dessa studier involverade 29 504 testfall skapade av 16 obfuskeringsverktyg mot 100 datorer skyddade av antingen Microsoft Windows Defender eller Symantec Endpoint Protection. Alla testfall kartlades mot en kategoriseringsmodell med 11 obfuskeringstekniker. Av testfallen passerade 54% anti-virusens statisk analys (inför kodexekvering) och 5% även den dynamiska analysen (kodexekvering).

Emedan verktyg utvecklade inom projektet har tillämpats för att generera angrepp under ett flertal cybersäkerhetsövningar där logganalytiker deltagit (se kapitel 3) så har forskningen som utförts inom projektet kring dessa övningar fokuserat på att utvärdera angreppsverktygens förmåga snarare än att studera effekten för olika arbetsmetoder eller logganalysverktyg. Av denna anledning avgränsas vidare diskussion av fråga (2) från denna rapport.

2 Automatisering av test- och träningsscenarion

Detta kapitel svarar på forskningsfrågan ”*Hur bör verktyg för att automatisera generering av test- och träningsscenarion inom logganalys och incidenthantering konstrueras?*”.

Test- och träningsscenarion i denna kontext rör som tidigare beskrivits cyberangrepp. Studierna av forskningsfrågan har varit tätt kopplade till verktyget Lore, vilket beskrivs i kapitel 2.1. Ett sammanfattande svar på hur verktyg som automatiskt utför cyberangrepp bör fungera presenteras i [20]:

1. **Angreppsemulatorer bör vara enkla att använda, underhålla och vidareutveckla:** Cyberförsvarsövningarna Lore har använts för har haft olika typer av planerade angreppsaktiviteter. Till exempel användes många serverangrepp (t.ex. MS17-010) under övningen Safe Cyber 2020, medan inga serverangrepp var tillämpliga för övningen Safe Cyber 2022. Safe Cyber 2022 involverade istället konfigurationsfel och horisontell förflyttning med hjälp av Windows-tjänster. Automatiserade cyberangreppsverktyg behöver därför stödja många olika typer av angreppsaktiviteter. Ett stöd för många angreppstyper är en del Lores kärndesign, och förmodligen en huvudsaklig anledning till att Lore kan producera lagom utmanande uppgifter för logganalytiker.
2. **Angreppsemulatorer bör göra lämpliga val.** Förutom att stödja många typer av angrepp är det viktigt att de lämpligaste utförs. Vilka angrepp som är lämpliga beror på den hotaktör som emuleras. Till exempel betar sig en datormask annorlunda än en motiverad statsaktör. Planeraren som används av Lore kombinerar boolesk logik och tränade modeller. Den booleska logiken innebär att tillämpbara angrepp väljs ut baserat på domänkunskap, och att helt orealistiska angrepp undviks. Maskininlärningsmodeller möjliggör sedan prioritering av angrepp med låg arbetsinsats (ingen specialskriven kod krävs för att prioritera nya typer av angrepp).
3. **Angreppsemulatorer bör inte ta genvägar.** Det är vanligt att mänskliga röda lag tar genvägar i övningar, t.ex. hoppar över nätverkskartläggning eftersom nätverkstopologin på förhand är känd, eller nyttjar behörigheter för att fjärrlogga in i en dator utan att dessa behörigheter samlats in (t.ex. genom mimikatz). Sådana genvägar gör logganalys svårare och mer konstgjorda. Lore tillämpar i sitt standardutförande inga genvägar, vilket även var fallet för de övningar som diskuteras i denna rapport. Detta är sannolikt en viktig anledning till varför Lores angrepp ses som realistiska och lagom svåra.

Arkitekturval och utvärderingar för Lore gällande dessa punkter beskrivs ytterligare i kapitel 2.1, 2.2 och 2.3.

2.1 Lore

Som beskrivs i [7] är Lore ett verktyg som utför samma typ av aktiviteter som (simulerade) hotaktörer. Detta kapitel redogör kort för verktygets övergripande funktion. En översikt av Lore ges av Figur 1. Läsaren hänvisas till [7], [9], [10] för mer omfattande beskrivningar av Lore. För ytterligare beskrivningar av verktyget SVED som Lore använder hänvisas läsaren till [6] och för vidare beskrivning av Crate som är den IT-infrastruktur Lore och SVED är inbyggda i hänvisas läsaren till [11].

Lore automatiserar de allra flesta typer av aktiviteter som röda lag utför, såsom:

- kartläggning av IP-rymder, portar och domäner
- mjukvaruangrepp mot servrar och klienter
- lösenordsgissning och lösenordsknäckning
- exekvering av kommandon på fjärrstyrda datorer, såsom avlyssning av nätverkstrafik, genomgång av filsystem och extrahering av lösenord.

Totalt automatiseras cirka 2000 olika typer av aktioner av standardautomationsfunktionen i verktyget. Utöver dessa kan en operatör manuellt utföra cirka 4000 andra typer av aktioner.

Lore arbetar utan krav på användarinteraktion enligt det scenario som specificerats av dess operatör. En operatör har möjlighet att ändra på ett antal parametrar för ett scenario, exempelvis:

- vilka maskiner, nätverkssegment och användarkonton som skall inkluderas eller exkluderas i testerna
- vilken förkunskap om målet som det röda laget besitter
- vilka handlingar som är tillåtna
- vilka handlingar som skall prioriteras upp eller ner
- vilka handlingar som skall utföras utöver standardbeteendet
- vilka typer av filer som skall laddas ner från fjärrstyrda maskiner och delade mappar
- vilka särskilda ordlistor som skall användas (t.ex. för lösenordsknäckning eller fuzzning)
- vilka bakdörrar som skall nyttjas (de som stöds av ramverken Metasploit och sliver)
- vilka pivoteringsmekanismer som skall användas (meterpreter och/eller chisel)
- vilka datorer som skall användas för angreppen, samt olika inställningar för dessa

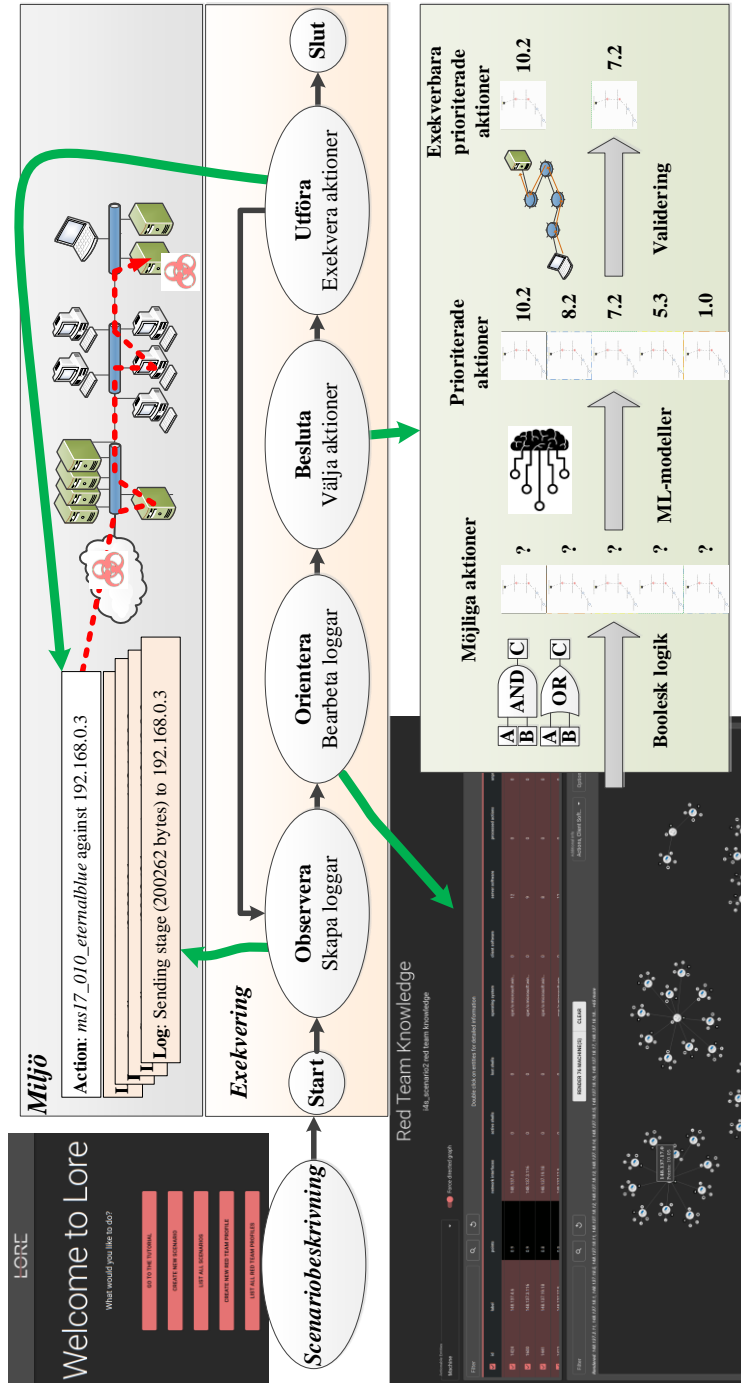
- vilka specialanpassade skript som skall användas för att drastiskt ändra beteendet (t.ex. för att få Lore att simulera hotaktören APT29 [12]).

En exekvering av ett scenario börjar med att ett antal förberedelsesteg utförs. Sådana steg involverar exempelvis att ställa in IP-konfigurationen för angripardatorerna och att starta globala lyssnare för inkommande skal (bakdörrar). Efter förberedelsefasen startar den huvudsakliga processloopen. Denna bygger på OODA-loopen [14] och innefattar:

- *observation* av loggar (eng: observe),
- *bearbetning* av loggar (eng: orient),
- *beslut* om aktioner baserat på situationen (eng: decide),
- *utförande* av beslutade aktioner (eng: act).

Observation och bearbetning handlar om att bygga upp den kunskapsdatabas som Lore sparar för ett scenario under dess utförande, såsom vilka händelser som utförts, vilka behörigheter som erhållits, vilka datorer som identifierats, och närhet mellan olika datorer. Processloopen pågår tills att antingen alla mål som finns definierade för ett scenario har uppfyllts eller tills dess att det inte längre finns handlingar kvar att utföra. Exempel på mål är att ta över särskilda datorer eller hitta filer på datorer. Under de övningar som Lore nyttjats för har scenariot avslutats efter att en viss mängd tid passerats (under cybersoldaternas slutövning 2022 var det exempelvis efter 24 timmar).

Mycket av forskningen kring Lore har rört beslutsprocessen: att identifiera vilka handlingar som är bäst att utföra givet ett visst tillstånd. Detta beskrivs vidare i kapitel 2.3. Prioriterade händelser undersöks sedan för att bedöma om de är giltiga att exekvera under ett givet tillfälle, och vid behov komplettera gällande särskilda inställningar. Ett exempel på en sådan särskild inställning är den aktuella pivoteringskedja (ofta via verktyget proxychains) som behövs för att nå ett givet mål vid ett givet tillfälle. Optimal väg för nätverkstrafik till en viss servermjukvara kan förändras beroende på vilka resurser hotaktören äger i terrängen vid ett givet tillfälle.



Figur 1: Översikt av Lore.

2.2 Enkel att använda, underhålla och vidareutveckla

Lore har en modulariserad arkitektur där det mesta av dess beteende kan anpassas via inställningar eller skript som läses in reflektivt. Detta medför kostnadseffektiv anpassning av beteende samt utveckling av helt nya beteenden, såsom att få Lore att bete sig som en särskild hotaktör.

Granstedt [12] studerade hur stor andel av stegen som hotaktören APT29 tillämpar som enkelt kunde konverteras till en angriparsprofil i Lore. Granstedt kom fram till att 90% av stegen utförda av APT29 helt kunde implementeras, och att de resterande 10% av stegen delvis kunde implementeras.

Alla cybersäkerhetsövningar som Lore har använts inom har nyttjat dess standardskript och tillämpat olika typer av inställningar för att anpassa dess beteende. De vanligast använda inställningarna är prioriteringsförändringar, exkludering av utvalda moduler, och särskilda kommandon som skall utföras på vissa övertagna maskiner.

Användartester av Lores grafiska gränssnitt visade att de flesta undersökta uppgifter gick att utföra utan någon hjälp oavsett roll och erfarenhet från Lore eller dess gränssnitt [13].

2.3 Göra lämpliga val

Före den vidareutveckling av Lore som skett inom VECNO var väntevärdet för utförande av olika handlingar mestadels hårdkodat. Exempelvis kunde handlingen "nmap -sn" ha väntevärdet 0.91 och "nmap -A" ha väntevärdet 0.88. Det huvudsakliga undantaget var serverangrepp i Metasploit, för vilka det hade tränats en maskininlärningsmodell som förutsåg sannolikheten för framgång [8], [10].

Denna lösning var problematisk eftersom att utvecklare av nya Lore-moduler var tvungna att ta ställning till prioriteten för alla andra handlingar vid skapande av en ny. Detta kunde exempelvis innefatta att bedöma hur stor vikt "4/255 datorer svarade på ping" har jämfört med "1 NTLM-hash hittades för användaren Kalle" eller "10% chans att serverangrepp fungerar". Dessutom beror vikten för en händelse på förutsättningarna. Exempelvis kan analyser gjorda av olika verktyg överlappa, där ordningen på utförande får stor effekt på väntevärdet.

Som lösning på detta problem utvecklades ett poängsystem för de olika egenskaperna i kunskapsdatabasen, där olika typer av information gavs poäng i relation till en "perfekt bakdörr", vilken ges 10 poäng [9]. Maskininlärningsmodeller som predikterar den förväntade poängen för olika handlingar tränas genom att exekvera handlingar under olika förutsättningar och mäta hur många poäng de generar. Totalt utfördes ett par miljoner olika

handlingar mellan 2023 och 2025³ för att möjliggöra övervakad inläring av maskininlärningsmodeller för elva olika typer av handlingar:

1. *lokala angrepp* (t.ex. eskaleringsangrepp i Metasploit eller skriptet `gtfonow.py`⁴)
2. *serverangrepp* (t.ex. serverangrepp i Metasploit)
3. *insamling av behörigheter* (t.ex. mimikatz, LaZagne eller `impacket-secretsdump`)
4. *nätverkskartläggning* (t.ex. via `nmap`)
5. *webb* (t.ex. `wget`, `curl`, `ffuf` eller `nuclei`)
6. *active directory* (t.ex. `bloodhound`)
7. *lösenordsknäckning* (t.ex. `hashcat`)
8. *lösenordsgissning* (t.ex. `kerbrute`, `netexec` eller `hydra`)
9. *generella skalkommandon* (t.ex. ladda upp och exekvera `chisel` på ett målsystem)
10. *generell informationsinhämtning* (t.ex. nedladdning av en fil från en delad mapp)
11. *övriga handlingar* (t.ex. `impacket-ticketConverter`).

Data insamlades genom experiment i cyberanläggningen Crate i Linköping [9]. Det genomfördes utvärderingar med diverse olika typer av modeller, såsom tre varianter av djupa neurala nätverk, random forests, XGBoost och LightGBM. En översikt av resultatet för den sista tränings-sessionen (mellan september och oktober 2025), där cirka 470 000 datapunkter samlades in, beskrivs av Tabell 2 (beskrivande statistik) och Tabell 3 (förklarad varians för två upplärda modeller). Förklarad varians beskrivs av måttet r^2 , där 1 innebär att modellen helt kan förutsäga utfallet medan 0 innebär att den är helt oförmögen att förutspå utfallet. Hur stort r^2 bör vara för att anses som bra beror på tillämpning, men ett r^2 som är större än eller lika med 0,3 ses av vissa som en godtagbar förklaringsgrad [16].

Som kan ses i Tabell 3 förklarar XGBoost (och ensemble-metoder överlag) en större del av den observerade variansen än neurala nätverk (oavsett testad variant). Detta beror förmodligen på att djupa neurala nätverk behöver betydligt fler datapunkter för att hitta mönster [17]. Modellerna har en relativt god förståelse kring vad som är bäst att göra för att maximera poäng på kort sikt: den förklarade variansen för XGBoost-modellen inom varje modell är mellan 0,33 och 1,0, med ett snitt på 0,63. Noterbart är att modellen som predikterar lokala angrepp, på grund av dess få datapunkter, förmodligen inte är särskilt tillförlitlig trots dess relativt höga r^2 -värde (0,65).

³ Datainsamlingen itererades årligen för att lära Lore vikten av nya typer av handlingar som tillkommit sedan senaste inläringstillfället.

⁴ <https://github.com/Frissi0n/GTFONow>

Tabell 2. Resultat från övervakad inläring.

Kategori	Antal tester	Poäng (medel)	Poäng (avvikelse)	Antal särdrag
Active directory	1764	2,88	7,73	199
Generella skalkommandon	52790	0,67	2,78	237
Insamling av behörigheter	10560	1,03	2,51	234
Lokala angrepp	130	4,34	8,17	311
Lösenordsgissning	153599	0,04	0,71	204
Lösenordsknäckning	8944	-0,02	0,01	190
Nätverkskartläggning	42275	0,26	0,67	197
Serverangrepp	108732	-0,01	0,25	318
Webb	44816	0	0,01	199
Övriga	46849	0,05	0,45	197

Tabell 3. Korrekthet för modeller tränade genom övervakad inläring.

Kategori	Modell	r ² (medel)	r ² (avvikelse)
Active directory	Neurala nätverk	0,63	0,02
	XGBoost	0,7	0
Generella skalkommandon	Neurala nätverk	0,48	0,25
	XGBoost	0,74	0
Insamling av behörigheter	Neurala nätverk	0,13	0,03
	XGBoost	0,33	0
Lokala angrepp	Neurala nätverk	0,42	0,28
	XGBoost	0,65	0
Lösenordsgissning	Neurala nätverk	0,38	0,05
	XGBoost	0,52	0
Lösenordsknäckning	Neurala nätverk	1	0
	XGBoost	1	0
Nätverkskartläggning	Neurala nätverk	0,68	0,02
	XGBoost	0,71	0
Serverangrepp	Neurala nätverk	0,3	0,04
	XGBoost	0,49	0
Webb	Neurala nätverk	0,15	0,02
	XGBoost	0,49	0
Övriga	Neurala nätverk	0,43	0,1
	XGBoost	0,7	0

Emedan de upptränade modellerna är av relativt god kvalitet har systemet dock en bristande förståelse kring kopplingarna mellan olika utförda händelser. Detta

eftersom prioriteten för varje händelse bedöms utan att ta hänsyn till andra händelser, mer än gällande mycket abstrakta egenskaper såsom ”hur många händelser har utförts mot komponent X?” och ”hur många poäng är komponent X för tillfället värd?”. Systemet har därmed ingen god kunskap om huruvida en viss händelse är den bästa givet ändamålet att nå en hög poäng på sikt.

Av denna anledning utvecklades en ny metod som predikterade vilken komponent som Lore bör fokusera på, såsom en viss dator eller ett visst användarkonto [18]. Med denna metod sker beslutsprocessen i två steg – först utvärderas vilken komponent som bör angripas, och sedan vilken händelse som bör utföras mot komponenten. Upplärning av modellen som predikterar vilken komponent som skall angripas sker genom oövervakad inläring med djupa neurala nätverk. Det är också möjligt att applicera övervakad inläring som ett steg inför den oövervakade inläringen.

Det gjordes en preliminär utvärdering av denna metod under 2023 [18]. Resultaten från denna studie visade att det krävdes avsevärt mer data än vad som var möjligt att generera genom tester i Crate. Av denna anledning har det i projektet påbörjats utveckling av en simulator som kan prediktera utfallet för olika handlingar baserat på utfallet för handlingar med samma förutsättningar vid tidigare tillfällen. På ett sådant sätt kan handlingar som enbart samlar in information (till skillnad från att exempelvis editera en routingtabell eller generera en baddörr) undvikas för att snabbare generera data. Arbetet med att färdigställa och tillämpa denna simulator lämnas till framtida arbete.

3 Utvärderingar av automatiserade test- och träningsscenarion

Lore har hittills använts under åtta cybersäkerhetsövningar:

- **Safe Cyber 2020 & 2022**
Övningar som involverade myndigheter och företag inom totalförsvaret
- **Övningar med försvarsmakten 2022, 2023, 2024 (2) & 2025**
Huvudsakligen cybersoldaternas slutövningar
- **Nordic American CERT exercise 2023**
En övning som involverade CERT-personal (Computer Security Incident Response Team) från USA och Europa

Av dessa övningar har enbart en organiserats av VECNO (en övning med Försvarsmakten under 2022). De andra övningarna har organiserats av andra projekt inom FOI. Att det finns behov av att nyttja Lore för cybersäkerhetsövningar kan i sig ses som en signal på att verktyget fungerar för ändamålet.

Inga deltagare under dessa övningar har fått någon förhandsinformation om att de skulle utsättas för en autonom hotaktör – från deras synvinkel har det varit helt vanliga övningar.

Det har under och efter varje övning genomförts kvalitativ återkoppling från deltagarna för att bedöma deras upplevelse av inspelen utförda av Lore. Exempelvis hur lärorika, realistiska och intressanta de ansåg dem vara. Denna återkoppling har varit överväldigande positiv, där många deltagare anmärkt hur engagerande och lärorika övningarna varit. Några övningar har förutom mer kvalitativ återkoppling dessutom erbjudit möjligheter att kvantitativt jämföra Lore med inspel som planerats av människor. Dessa utvärderingar beskrivs i kapitel 3.1.

Utöver att jämföra Lore med mänskliga inspel kan det tänkas finnas andra verktyg som fyller samma behov. Kapitel 3.2 sammanfattar de studier som genomförts för att jämföra Lore med andra liknande verktyg.

3.1 Utvärderingar under cybersäkerhetsövningar

Det har genomförts två studier inom VECNO som undersökt hur väl Lore kan ersätta mänskligt genererade inspel under övningar [19], [20]. Dessa studerade bland annat:

- Om Lore producerar lika realistiska inspel som de inspel som utförs av människor
- Om Lore producerar lika lärorika inspel som de inspel som utförs av människor
- Om det är lika svårt att utföra incidentanalyser för inspel som utförs av Lore som inspel utföra av människor

Studierna genomfördes baserat på data erhållen från tre övningar: Safe Cyber 2020 och 2022, samt en övning för Försvarmakten under 2024. Totalt involverade de tre övningarna 132 logganalytiker. Under Safe Cyber-övningarna nyttjades 2x2 försöksplaner som utsatte logganalytikerna för Lore och mänskliga angrepp om vartannat. För övningen för Försvarmakten under 2024 var det av övningstekniska skäl inte möjligt att tillämpa en formell försöksplan. Denna övning hade dock inga överlappande inspel mellan Lore och människor.

Analysen gjordes på två sätt:

- En central uppgift för logganalytikerna under alla tre övningarna var att rapportera alla identifierade säkerhetsincidenter. Dessa rapporter granskades med avseende på korrekthet.
- Under Safe Cyber-övningarna samlades deltagarnas åsikter in genom enkäter vid lunch samt efter övningens utförande (det vill säga, efter varje steg i försöksplanen).

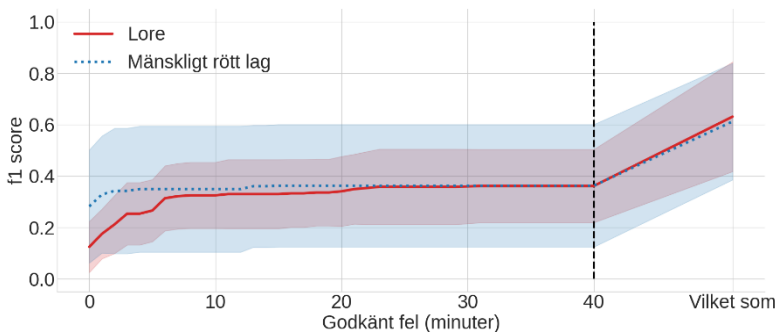
Resultatet från enkäterna visade på att det inte fanns några statistiskt signifikanta skillnader mellan inspel utförda av Lore och människor. Med andra ord, den upplevda realismen var ungefär lika hög, och inspelen ungefär lika lärorika, oavsett om Lore användes eller inte [19], [20] (se tabell 4). Det kan noteras att medelvärdet för variabeln *Upplevd svårighetsgrad* är lägre än minvärdet för skalan (1). Anledningen är att det antogs att en upplevd svårighetsgrad på 4,5 var optimalt, och medelvärdet är därför den genomsnittliga avvikelser från 4,5 (en högre avvikelse är sämre).

Tabell 4. Upplevda skillnader mellan Lore och mänskliga inspel under övningar (SD: standardavvikelse). Alla frågor besvarades med en Likertskala (1-7).

Variabel	Hotaktör	Antal	Medel	SD
<i>Upplevd lärdom</i> (1–7)	Mänskligt rött lag	61	4,72	1,29
	Lore	63	4,89	1,42
<i>Upplevd realism</i> (1–7)	Mänskligt rött lag	57	4,72	1,44
	Lore	57	4,31	1,53
<i>Upplevd svårighetsgrad</i> (1-7, mätt i form avvikelse från 4,5)	Mänskligt rött lag	59	0,76	0,67
	Lore	63	1,03	0,83

En sammanfattning av korrektheten för incidentrapporterna presenteras i figur 2. I figuren presenteras korrekthet genom måttet f1 score på den vertikala axeln och godkänt fel (i minuter) på den horisontella axeln. På ett övergripande sätt kan ett f1 score på 1 tolkas som att logganalytikerna upptäckt alla angrepp utan att felaktigt rapportera in någon godartad händelse som ett angrepp. Vidare tilläts logganalytiker ha ett fel på ett visst antal minuter i sin incidentrapport. Exempelvis innebär ett godkänt fel på 30 minuter att en rapport räknas som godkänd om det i den står att en dator angreps 11:00, när den i praktiken angreps 10:30, men inte om datorn rapporteras som angripen 11:01.

Resultatet från incidentrapporterna visade på samma mönster som enkätstudien: det var inga statistiskt signifikanta skillnader för korrektheten av incidentrapporter om inspelen skapades av Lore eller människor. Däremot kunde det noteras att spridningen var högre för de mänskliga inspelen: de olika grupperna gjorde ett jämnare resultat för inspel skapade av Lore än människorna [20].



Figur 2. Korrekthet för incidentrapporter för angrepp skapade av Lore och mänskliga röda lag.

3.2 Jämförelser med andra verktyg

Det finns många verktyg som kan automatisera de cyberangrepp som test- och träningsscenarier för logganalys och incidenthantering ställer krav på. Inom VECNO har det genomförts flera tekniska tester för att utvärdera sådana verktyg tillhandahållna som öppen källkod [21], [22]. Detta arbete fortsatte i samma spår som det arbete som tidigare utfördes inom ÖvExCND [23], [24]. Testade verktyg har involverat allt ifrån enkla skript som automatiserar serverangrepp i Metasploit till system som tillämpar stora språkmodeller för att i teorin lösa komplexa utmaningar. Angående stora språkmodeller testades verktyget PentestGPT [25] i [21]. Denna studie visade att en användare av PentestGPT behövde avsevärd kunskap kring utförande av penetrationstester då verktyget i bästa fall gav vaga beskrivningar av de steg som behövde utföras. Dessa studier

fann att de största konkurrenterna till Lore var Caldera [26] och Infection Monkey⁵ [22], [27].

Caldera har utvecklats av MITRE och är troligen det mest omnämnda (och tillämpade) verktyget i akademien. Infection Monkey har utvecklats av företaget Akamai. Jämförelser av dessa verktyg samt verktyget Deep Exploit⁶ mot Lore genom tekniska tester i cyberanläggningen Crate presenteras i [22], [27]. En översikt av utfallet presenteras i tabell 5.

Tabell 5. Jämförelser av Caldera, Deep Exploit, Infection Monkey och Lore (medelvärden och standardavvikelser).

Variabel	Facit	Caldera	Deep Exploit	Infection Monkey	Lore
Hittade konton	134	1 (1)	0 (0)	0 (0)	124 (1)
Hittade grupper	8	0 (0)	0 (0)	0 (0)	8 (0)
Hittade domäner	2	0 (0)	0 (0)	0 (0)	2 (0)
Hittade datorer	56	15 (16)	0 (0)	13 (0)	56 (0)
Hittade nätverk	8	4 (1)	0 (0)	1 (0)	8 (0)
Övertagna nätverk (admin)	7	1 (1)	0 (0)	0 (0)	7 (0)
Övertagna nätverk (användare)	7	0 (0)	0 (0)	0 (0)	4 (0)
Övertagna nätverk (totalt)	7	1 (1)	0 (0)	0 (0)	7 (0)
Övertagna datorer (admin)	55	4 (7)	0 (0)	0 (0)	52 (1)
Övertagna datorer (användare)	55	0 (0)	0 (0)	0 (0)	12 (4)
Övertagna datorer (totalt)	55	4 (7)	0 (0)	0 (0)	53 (1)
Erhållna lösenord	134	0 (1)	0 (0)	0 (0)	13 (2)
Erhållna lösenordshashar	134	1 (1)	0 (0)	0 (0)	106 (25)
Erhållna lösenord eller lösenordshashar	134	1 (1)	0 (0)	0 (0)	106 (25)
Wazuh-larm (>= nivå 6)	-	8604 (4609)	9718 (371)	8736 (484)	9983 (950)

Infection Monkey lyckades inte ta över datorer i något studerat scenario. Caldera omfattar rent teoretiskt många av teknikerna i MITRE ATT&CK [22], men är i praktiken inte särskilt bra på att automatisera dem. Medan Caldera helt

⁵ <https://github.com/guardicore/monkey>

⁶ https://github.com/13o-bbr-bbq/machine_learning_security/blob/master/DeepExploit/

misslyckas med att ta över datorer i tre av fyra studerade scenarier⁷, och i det fjärde scenariot tar över cirka 17 datorer, tar Lore över i princip alla datorer oavsett scenario. Den största anledningen till detta är att övertagande av fler datorer ställer krav på moduler som inte stöds av Caldera, såsom pivotering via fjärrstyrda datorer, avlyssning av nätverkstrafik, serverangrepp, lokala angrepp, lösenordsknäckning, lösenordsgissning, lösenordsextrahering, och horisontell förflyttning genom exempelvis SSH, WinRM och PsExec.

Vi förmodar att anledningen till detta är att Caldera är tänkt att anpassas för varje scenario, till skillnad från Lore som är tänkt att kunna fungera väl oavsett scenario. När Lore använts i cybersäkerhetsövningar har det istället satts begränsningar på hur verktyget skall agera. Baserat på vår erfarenhet krävs det avsevärt mindre ansträngning och kunskap att begränsa ett beteende än att tillföra nya beteenden.

⁷ Varje scenario innefattade olika typer av ursprungliga bakdörrar: Windows eller Linux, samt med eller utan system/root-behörighet.

4 Slutsatser och framtida arbete

VECNO har vidareutvecklat verktyget Lore som utför automatiska cyberangrepp långt bättre än vad något annat liknande verktyg som vi testat klarar. Lore har dessutom, till skillnad från dess konkurrenter, utvärderats med goda resultat under ett flertal cybersäkerhetsövningar och användartester av olika slag (tex. [3], [12], [13]). Svaret på forskningsfrågan ”Hur bör verktyg för att automatisera generering av test- och träningsscenarion inom logganalys och incidenthantering konstrueras” ligger dels inbäddat i de olika delarna av Lore, såsom vilka handlingar som är relevanta och hur de bäst bör exekveras, och dels i de artiklar som publicerats inom projektet. Ett sammanfattande svar på hur verktyg som automatiskt utför cyberangrepp bör fungera är [20]:

1. angreppsemulatorer bör vara enkla att använda, underhålla och vidareutveckla
2. angreppsemulatorer bör göra lämpliga val
3. angreppsemulatorer bör inte ta genvägar.

Gällande forskningsfråga (1) finns det många aspekter som kräver vidare studier. Dessa beskrivs i resterande del av detta kapitel.

Lore har enbart jämförts med verktyg skrivna i (mestadels) öppen källkod. Det är möjligt att tänka sig jämförelser mot kommersiella verktyg, såsom Pentera⁸ och vPenTest⁹. Att studera kommersiella verktyg är dock av naturliga skäl dyrt. Exempelvis genomfördes två telefonmöten med utvecklarna av Pentera under 2024, och det visade sig att det hade krävt en stor del av projektets budget för att över huvud taget testa den enklaste versionen av verktyget. Det är också möjligt att stora språkmodeller blivit bättre på att utföra penetrationstester sedan vi testade PentestGPT [21]. Bland annat påstås verktyget CAI och PentestGPT numera kunna lösa många utmaningar i Hack The Box, om än med ”viss användarinput”:

“While CAI explores autonomous capabilities, our results clearly demonstrate that effective security operations still require human teleoperation providing expertise, judgment, and oversight in the security process. The Human-In-The-Loop (HITL) module is therefore not merely a feature but a critical cornerstone of CAI’s design philosophy.” [28].

Vi ämnar utföra fler utvärderingar av sådana verktyg. Vi ämnar också integrera stora språkmodeller i Lore för att kunna lösa utmaningar som ställer stora krav på kreativitet, och därmed är svåra att skapa generella lösningar för. Exempelvis att

⁸ <https://pentera.io/penetration-testing/>

⁹ <https://www.vonahi.io/>

hitta behörigheter i en fil, att hitta en okänd sårbarhet i ett skript som möjliggör eskalering av behörigheter, eller att få förslag på nya handlingar när tillståndsrymden är helt uttömd.

Det finns också planer på att bättre anpassa Lore för penetrationstester och capture-the-flag-utmaningar. Sådana tillämpningar ställer andra krav än cybersäkerhetsövningar, som fokuserar på utmaningar för försvarande deltagare. Framförallt krävs ett långt större modulstöd eftersom det till skillnad från övningsnät oftast inte är möjligt att modifiera målmiljön, såsom att lägga till en sårbarhet. Under 2025 genomfördes en förstudie av detta genom att utveckla automatiserade lösningar för ett antal av maskinerna och nätverken¹⁰ i plattformen Hack The Box. I skrivande stund har Lore testats mot 14 maskiner: för två av dessa maskiner kan Lore ge fullständiga lösningar och för 12 partiella lösningar (t.ex. eskalera användarbehörighet till root). Dessutom finns partiella lösningar för de tre nätverken "Dante", "P.O.O." och "RastaLabs". De steg som inte kan automatiseras av Lore rör huvudsakligen angrepp genom för utmaningarna specialskrivna skript, såsom eskalering via särskilda argument till ett bash-skript som en vanlig användare har rättighet att köra som root. Våra preliminära tester visar att denna typ av utmaning är en god kandidat för lösning genom stora språkmodeller.

På liknande sätt har även moduler utvecklats för att möjliggöra många av angreppsmetoderna i Game of Active Directory¹¹ (GOAD), vilken är en öppet tillgänglig labbmiljö med ett stort antal Windows Active Directory-sårbarheter. Det kvarstår dock arbete med att lära Lore hur värdefulla olika angrepp är i GOAD. Som en del i detta arbete är dess labbmiljö planerad att inkorporeras i Crate för att möjliggöra mer omfattande datainsamling.

Lore tillämpades också för cybersäkerhetsutmaningen *Stealth Cup*¹² som anordnades av Austrian Institute Of Technology (AIT). Under Stealth Cup tävlade 12 lag bestående av professionella penetrationstestare. Lore lyckades nå bättre behörigheter än 4 av dessa lag och samma behörighet som 7 av lagen. Endast ett lag lyckades lösa alla utmaningar i Stealth Cup (och därmed uppnå ett bättre resultat än Lore).

Att tillämpa Lore för att lösa capture-the-flag-utmaningar har varit mycket givande och utmanande, och utfallet har kunnat återanvändas för andra ändamål. Exempelvis nyttjades moduler som skapades för att lösa utmaningar i Hack The Box under en cybersäkerhetsövning med Försvarmakten under 2025.

¹⁰ Dessa kallas för "Pro Labs" i Hack The Box.

¹¹ <https://github.com/Orange-Cyberdefense/GOAD>

¹² <https://stealth.ait.ac.at/>

Det finns också planer på att utöka möjligheten för användarinteraktion med verktyget. Under cybersäkerhetsövningar finns det ett behov av att kunna justera Lore's beteende, primärt för att ändra svårighetsgraden i analysarbetet för deltagande logganalytiker; under utmaningar såsom Stealth Cup finns det ett behov för cybersäkerhetsexperter att kunna delge ny information och starta specifika händelser.

Det kvarstår även mycket arbete med att förfina beslutsprocessen för Lore. Framförallt är en mer omfattande studie av öövervakad inläring planerad som en fortsättning på studien som beskrivs i [18].

Slutligen bör det poängteras att Lore inte är ämnat ersätta mänskliga operatörer – enbart underlätta för dessa. För att citera en av dess användare:

”Lore är bra på att göra allt tråkigt som man inte lär sig något av att göra. Sedan kan man göra det roliga kreativa arbetet själv.”

5 Referenser

- [1] T. Sommestad, ”Sammanfattning av FOI:s arbete inom projektet SAC3S/Tyr under 2023”, Totalförsvarets forskningsinstitut (FOI), FOI Memo 8460, mar. 2024.
- [2] T. Sommestad, ”Sammanfattning av FOI:s arbete inom projektet SAC3S/Tyr under 2024”, Totalförsvarets forskningsinstitut (FOI), FOI Memo 8812, feb. 2025.
- [3] Sommestad, Teodor, ”SAFE Cyber 2020, genomförande av distribuerad övning (FOI Memo 7522)”, Totalförsvarets forskningsinstitut (FOI), apr. 2021.
- [4] P. Lif, ”Leveransplan 2025 för projektet Metod och kompetens för CNO”, Totalförsvarets forskningsinstitut (FOI), FOI Memo 8781, jan. 2025.
- [5] T. Sommestad, ”Övning och Experiment för operativ förmåga i cybermiljön: Slutrapport”, Totalförsvarets forskningsinstitut (FOI), FOI-R--4498--SE, dec. 2017.
- [6] H. Holm och T. Sommestad, ”Sved: Scanning, vulnerabilities, exploits and detection”, i *MILCOM 2016-2016 IEEE Military Communications Conference*, IEEE, 2016, s. 976–981.
- [7] H. Holm, ”Lore A Red Team Emulation Tool”, *IEEE Transactions on Dependable and Secure Computing*, 2022.
- [8] H. Holm, T. Sommestad, I. Rodhe, M. Persson, och P. Lif, ”Lore: Ett verktyg för automatisering av cyberangrepp under IT-säkerhetsövningar”, Totalförsvarets forskningsinstitut (FOI), FOI-R--4661--SE, nov. 2018.
- [9] H. Holm, ”Automatisering av cybersäkerhetsövningar: Vidareutveckling och evaluering av Lores beslutsprocess”, Totalförsvarets forskningsinstitut (FOI), FOI-R--5148--SE, juli 2021.
- [10] H. Holm, E. Hyllienmark, och M. Persson, ”Verktyg som döljer skadlig kod - En systematisk granskning”, Totalförsvarets forskningsinstitut (FOI), FOI-R--5366--SE, dec. 2022.
- [11] J. Almroth och T. Gustafsson, ”CRATE Exercise Control--A cyber defense exercise management and support tool”, i *2020 IEEE European Symposium on Security and Privacy Workshops (EuroS&PW)*, 2020, s. 37–45. doi: 10.1109/EuroSPW51379.2020.00014.
- [12] E. Granstedt, ”Evaluating the Realism and Effectiveness of Automated APT Emulation in Cybersecurity Training : A Lore Case Study”, Master’s Thesis, KTH, School of Electrical Engineering and Computer Science (EECS), 2024.
- [13] O. Johansson, ”Användartester för teknisk informationshämtning i Lore GUI”, Totalförsvarets forskningsinstitut (FOI), FOI Memo 8085, okt. 2023.
- [14] B. Brehmer, ”The dynamic OODA loop: Amalgamating Boyd’s OODA loop and the cybernetic approach to command and control”, i *Proceedings of the*

- Second International Conference on Military Technology*, Stockholm, Sweden, okt. 2005.
- [15] H. Holm och E. Hyllienmark, "Hide My Payload: An Empirical Study of Antimalware Evasion Tools", i *2023 IEEE International Conference on Big Data (BigData)*, 2023, s. 2989–2998. doi: 10.1109/BigData59044.2023.10386838.
- [16] G. E. Gignac och E. T. Szodorai, "Effect size guidelines for individual differences researchers", *Personality and Individual Differences*, vol. 102, s. 74–78, 2016, doi: <https://doi.org/10.1016/j.paid.2016.06.069>.
- [17] N. Hollmann *m.fl.*, "Accurate predictions on small data with a tabular foundation model", *Nature*, vol. 637, nr 8045, s. 319–326, 2025.
- [18] J. Karlsson och H. Holm, "Förstärkningsinläring för cyberangrepp", Totalförsvarets forskningsinstitut (FOI), FOI-R--5533--SE, jan. 2024.
- [19] H. Holm och J. Reuben, "Evaluation of a Red Team Automation Tool in live Cyber Defence Exercises", i *38th IFIP TC 11 International Conference, SEC 2023*,
- [20] H. Holm och T. Sommestad, "Realistic and balanced automated threat emulation", *Computers & Security*, vol. 151, s. 104351, 2025, doi: <https://doi.org/10.1016/j.cose.2025.104351>.
- [21] L. Helgeson, "Utvärdering av verktyg för emulering av hotaktör", Totalförsvarets forskningsinstitut (FOI), FOI Memo 8354, okt. 2023.
- [22] H. Holm och L. Helgeson, "Utvärdering av verktyg som emulerar hotaktörer", Totalförsvarets forskningsinstitut (FOI), FOI Memo 8662, nov. 2024.
- [23] H. Holm och T. Sommestad, "Omvärldsbevakning inom ÖvExCND under 2019", Totalförsvarets forskningsinstitut (FOI), FOI Memo 6941, dec. 2019.
- [24] H. Holm, "Omvärldsbevakning inom ÖvExCND under 2020", Totalförsvarets forskningsinstitut (FOI), FOI Memo 7365, nov. 2020.
- [25] G. Deng *m.fl.*, "PentestGPT: Evaluating and Harnessing Large Language Models for Automated Penetration Testing", i *33rd USENIX Security Symposium (USENIX Security 24)*, Philadelphia, PA: USENIX Association, aug. 2024, s. 847–864. [Online]. Tillgänglig vid: <https://www.usenix.org/conference/usenixsecurity24/presentation/deng>
- [26] A. Applebaum, D. Miller, B. Strom, C. Korban, och R. Wolf, "Intelligent, automated red team emulation", i *Proceedings of the 32nd Annual Conference on Computer Security Applications*, i ACSAC '16. New York, NY, USA: Association for Computing Machinery, 2016, s. 363–373. doi: 10.1145/2991079.2991111.
- [27] H. Holm och L. Helgeson, "An Empirical Study of Automated Adversary Emulators", presenterad vid 59th Hawaii International Conference on System Sciences (HICSS), Lahaina, Hawaii, jan. 2026.
- [28] V. Mayoral-Vilches *m.fl.*, "CAI: An Open, Bug Bounty-Ready Cybersecurity AI", *arXiv preprint arXiv:2504.06017*, 2025.



FOI
Totalförsvarets forskningsinstitut
164 90 Stockholm

Tel: 08-55 50 30 00
Fax: 08-55 50 31 00

www.foi.se