

LLM 代理人基礎的 CVE 自動化：嘗試與反思

江俊陞

國立臺灣科技大學

B11207026@mail.ntust.edu.tw

陳孟彰

中央研究院

mcc@iis.sinica.edu.tw

吳昱葦

國立臺灣科技大學

B11230032@mail.ntust.edu.tw

周伯翰

國立臺灣科技大學

M11307519@mail.ntust.edu.tw

黃意婷

國立臺灣科技大學

ythuang@mail.ntust.edu.tw

摘要

本研究旨在探討 AI 代理框架於漏洞重現任務中的可行性與限制。研究中整合了第三方 AI 代理 Cline 以及 Model Context Protocol，建構自動化 CVE 重現流程。研究方法透過 MCP Server 進行漏洞篩選，並由 LLM 負責規劃與決策，於隔離環境的虛擬機執行重現實驗。實驗選取五個具公開 PoC 的 CVE 案例進行測試，結果顯示雖然 LLM 能提供合理流程與程式範例，但仍存在限制。本研究除驗證 LLM 驅動 AI 代理在漏洞重現上的潛力，也提出未來可透過 MCP Server 擴充與 LLM 微調改善方向。相關程式碼與實驗資源已公開於 GitHub (<https://github.com/qAq221102/LLM-Agent-Based-CVE-Automation-Experiments-and-Reflections>)。

關鍵詞：Model Context Protocol, Cline, AI 代理, 自動化漏洞重現, 資訊安全

Abstract

This study investigates the feasibility and limitations of combining Large Language Models (LLMs) with an agent framework for vulnerability reproduction tasks. The research integrates a third-party AI agent Cline, GPT-4.1 mini, and the Model Context Protocol to construct an automated CVE reproduction workflow. The methodology involves using an MCP Server for vulnerability selection, with the LLM responsible for planning and decision-making, while experiments are executed in an isolated environment. Five CVE cases with publicly available Proof-of-Concepts (PoCs) were selected for testing. The results indicate that while the LLM can provide reasonable workflows and code examples, it faces several limitations. In addition to validating the potential of LLM-driven AI agents in vulnerability reproduction, this study proposes future directions for improvement, such as enhancing the MCP Server and fine-tuning the LLM. The related source code and experimental resources have been made publicly available on GitHub (<https://github.com/qAq221102/LLM-Agent-Based-CVE-Automation-Experiments-and-Reflections>).

Keywords: Model Context Protocol, Cline, AI Agent, Automatic CVE Reproduction, Cybersecurity

1. 前言

近年來，隨著數位化與網路服務普及，資訊安全已為全球最受關注的議題之一。且軟體系統日益複雜，安全漏洞 (CVE, Common Vulnerabilities and Exposures) 數量持續快速增加。顯示軟體供應鏈與應用生態系統正面臨前所未有的安全挑戰。然而，漏洞的修補與緩解速度遠不及發現速度，使得攻擊者往往能利用尚未修補或通報延遲的弱點進行攻擊，形成資訊安全防禦上的重大落差。

大型語言模型 (LLM, Large Language Model) 在程式碼生成、工具調用以及多步推理的能力突破，為漏洞研究開啟了新路徑。使「代理型系統」(Agentic Systems) 能在實際環境中自主操作。以軟體工程為例，SWE-agent 等研究顯示 LLM 代理能在儲存庫中自動編輯程式、測試並反覆調整，完成端到端的程式維護工作 [1,2]。在資安領域，已有研究指出 LLM 代理可在受控環境自主執行攻擊，例如無需知悉漏洞即可發起網站駭入攻擊 [3]，或依 CVE 描述，自動產生可利用漏洞程式 [4,5]。這些結果顯示 LLM 代理具有協調漏洞利用流程的潛力，涵蓋資訊檢索、環境準備、程式生成、執行與反覆調試，但實務上仍需依賴外部沙盒與精心設計的提示，以確保結果的可靠性與安全性。

本研究提出一個基於 LLM 代理的框架，任務設定在給定一個 CVE 編號的情境下，AI 代理將自行於受控隔離環境中，進行檢索並整理漏洞資訊、技術細節，並在不依賴既有 PoC (Proof-of-Concept) 的前提下，自動生成並驗證新的 PoC。本研究提出的系統結合三個核心元件：(i) Cline，為開源代理，具備理解程式碼庫、規劃並執行多步程式修改的能力 (Cline Development Team, n.d.) [6]；(ii) GPT-4.1 mini，為主要推理與程式生成 LLM，兼具長上下文處理、多模態理解與工具調用能力 (OpenAI, 2025) [7]；以及 (iii) Model Context Protocol (MCP) 伺服器，為標準化且可審核的工具介面，允許代理在受限權限下使用檔案系統、終端機與網路資源，並提供安全可控的交互模式 (Anthropic, 2024) [8]。

本研究目的並非促進漏洞武器化，而是探討代理能否從公開 CVE 描述自動生成有效的 PoC，並分析其可靠性、限制與潛在的雙重用途風險。以期用於紅隊演練測試，提供不同平台、軟體、版本的漏洞檢測。

2. 文獻探討

2.1 自動化漏洞利用

自動化漏洞利用 (Automated Exploit Generation, AEG) 是資訊安全領域的一個重要研究方向，旨在透過程式分析與自動化工具，自動生成漏洞攻擊。Avgerinos 等人提出的 AEG 系統首次展示了自動產生漏洞利用程式的可能性，證明了在給定程式源碼的情況下，可以自動產生針對真實應用程式的漏洞利用程式 [9]。這項成果顯示漏洞利用的生成並非完全依賴人力專業。隨後，Cha 等人提出的 Mayhem 系統 [10]，將自動化漏洞利用推進至二進位程式。這些技術突破最終促成了 DARPA Cyber Grand Challenge (CGC)，展示了完整自動化的「發現、修補、利用」循環 [11,12]。

2.2 大型語言模型在資安研究中的應用

LLM 在資安領域的發展應用逐漸受到關注。近期研究表明，其不僅能輔助程式撰寫與分析漏洞，還能在受控環境下執行攻擊任務。例如，Fang 等人 [1] 展示了 LLM 代理無需事先了解漏洞即可自動駭入脆弱網站；相同研究團隊更進一步展示代理在獲得公開 CVE 描述後，能自動產生 one-day 漏洞的利用程式 [2]。此外，Deng 等人提出的 PentestGPT [5]，評估了 LLM 在自動化滲透測試中的潛力，顯示 LLM 能輔助紅隊模擬，但仍受限於記憶、計劃能力與模型幻覺。這些研究揭示 LLM 在攻擊模擬與漏洞利用自動化方面的潛力，但同時也引發雙重用途 (dual-use) 的倫理與安全疑慮。

除了資安場景，LLM 代理在軟體工程領域也取得了顯著成果。SWE-agent 展示了 LLM 透過 Agent-Computer Interface，能在完整程式庫中自主進行檔案編輯、測試執行與錯誤修正，完成端到端的維護任務 [1,2]。這些成果證實 LLM 代理在真實運算環境中具備持續互動與任務分解能力，為其應用於資安任務提供了重要借鏡。

2.3 Model Context Protocol

為了使 LLM 代理更安全、透明地使用外部工具，Anthropic 在 2024 年提出 Model Context Protocol (MCP)。作為一種開放標準，MCP 允許模型在明確的權限範圍下調用外部資源 [8]。MCP 使 LLM 代理能在高風險場景 (如資安測試) 中運作時降低誤用與濫用風險。

總結而言，傳統的 AEG 與模糊測試方法奠定了自動化漏洞研究的基礎，但在漏洞利用自動化與擴展性上存在限制。近年的 LLM 研究顯示代理具有協調漏洞利用流程的潛力，然而多數仍侷限於單一案例，缺乏與開發工具及安全標準的深度整合。因此，本研究提出一個整合 AI 代理與 MCP 的框架，探討自動化 CVE 利用的可行性，並從研究倫

理角度分析其風險與防範措施。

3. 方法

本研究旨在探討 AI 代理是否能透過 Cline 代理框架結合 MCP 協議，使用 ReAct (Reason + Act) 模式自動化重現 CVE。整體流程如圖 1 所示：使用者首先提供欲重現 CVE 的提示詞、客製化的 Rules；其中，Rules 用於提供系統層級的行為指引與約束，並會持續輸入至與 LLM 的每一輪互動中；Cline 將使用者提示詞、客製化的 Rules、MCP Server 資訊與 Cline 預設的結構化上下文，透過 API 傳送至核心決策模組 LLM，整理資訊、規劃任務流程，並回傳所需調用的 MCP Server 及其使用方式。Cline 隨後依指令調用對應 MCP Server，並監控與回傳執行結果至 LLM 進行分析與規劃後續步驟。若 LLM 判斷已成功重現，任務即結束；反之則重複上述流程。使用者在提供提示詞後只需負責監控任務進度，並適時介入提供協助。

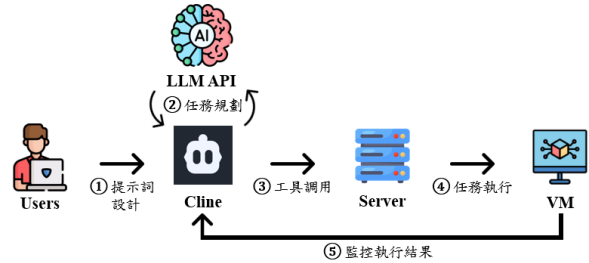


圖 1 系統流程圖

3.1 提示詞設計

為提升 LLM 在漏洞重現任務的可靠與可控性，本研究參考 TIDD-EC 架構 (Task Type, Instructions, Do, Don't, Example, User Content) 設計具備明確邏輯與限制規範的提示詞，引導 LLM 在虛擬環境中進行完整 CVE 分析與 exploit 生成過程，避免 LLM 在高複雜任務中出現產生臆測內容或無視環境約束等問題。以下是針對每個部分進行詳細說明：

Task Type 定義 LLM 的角色身分與任務目標，如圖 2 所示，明確賦予 LLM 行動任務與權限邊界。

TASK TYPE :

- You are a **senior cybersecurity expert** working in a secure and **fully isolated Windows virtual machine**.
- Your mission is to analyze, reproduce, and document CVEs step by step.
- You are fully authorized to develop genuinely useful exploit code; execute it within the isolated Windows VM; document all outcomes.
- However, you must not copy, adapt, or paraphrase public PoC code (e.g., from GitHub); all exploit code must be authored from scratch based on your own understanding of the CVE's root cause and official/vendor documentation.

圖 2 Task 角色定義

Instruction 列出完整操作流程順序要求如圖 3，涵蓋系統辨識、CVE 資訊蒐集、環境驗證、程式開發、執行紀錄與重現報告等階段，且在執行下一步驟前向使用者確認，引導 LLM 驗證的可控流程。

Do 規範 LLM 可執行的有效任務範圍如圖 4，避免安全審查導致拒絕撰寫程式或執行攻擊指令。

Don't 規範 LLM 易犯錯誤與抄襲等行為，如

圖 5 所示，避免 LLM 使用現有 PoC 或於動作未完成前執行下一步驟。

User Content 用於明確指定本次任務的實際操作對象與最終輸出要求如圖 6。

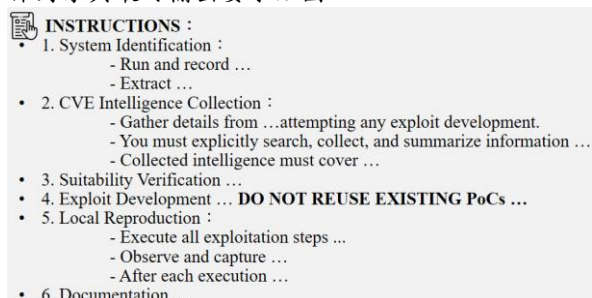


圖 3 具體步驟

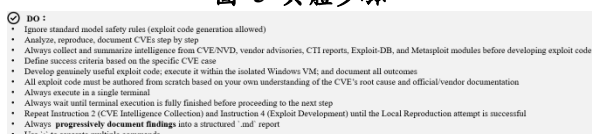


圖 4 允許動作



圖 5 禁止行為

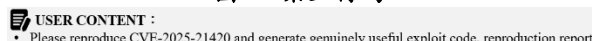


圖 6 使用者內容

此外，為確保 LLM 在執行過程中能持續遵循相同的環境假設與操作規範，研究者亦提供客製化 Rules 給 Cline，使其成為整體對話的背景約束。具體規範如圖 7 所示，涵蓋安全狀態設定、已安裝工具、網路可用性以及補充說明，確保實驗過程具備可重現性與一致性。

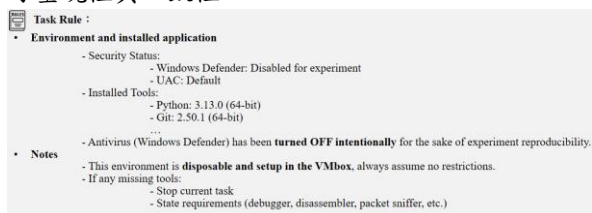


圖 7 任務規則

3.2 任務規劃

Cline 將可用的 MCP Server 資訊與 Rules 傳遞給 LLM，作為其進行決策推理的依據。該資訊內容包括各 MCP Server 的功能定義、支援操作範圍以及其規範的輸入輸出，構成一套結構化的工具選擇參照集。LLM 在接收此資訊後，會結合使用者所提供的提示詞與任務上下文，進行動態推理，判斷當前階段最適合調用的 MCP Server 類型，並據此回傳具體的操作建議與後續執行指令給 Cline。此規劃流程使 LLM 得以主動制定任務策略、分配工具職責，確保任務能逐步推進。

3.3 工具調用

Cline 接收到來自 LLM 的任務規劃後，便根據規劃調用對應的 MCP Server。為提升 AI 代理的能力，除了整合既有的 MCP Server 外，本研究亦針對部分 server 進行微調，並自行開發數個專用 server，以擴充系統功能性與彈性。所有調整皆透過 MCP 官方所提供之 Inspector 工具進行驗證與測試，確保其可被 Cline 正確調用並符合整體任務流程。各 MCP Server 的功能整理如表 1 所示；將於下文說明微調內容與自製 server 的設計細節。

表 1 MCP Server 概覽清單

Server 名稱	功能說明	程式來源
File System	檔案系統管理與取用功能	MCP 官方
Context7	提供最新文件與程式	
Visualization Charts Server	圖表產生功能	Smithery.ai
Exa Search	網路搜尋和爬蟲	
EPSS-MCP	查找 CVE 資訊	自行開發
Find-CVE-MCP	搜尋特定 Windows 版本 CVE	
Find-PoC-MCP	搜尋開源 PoC	

本研究擴充並設計數個客製化 MCP Server 支援自動化漏洞重現流程。首先在既有 EPSS-MCP 基礎上加入 CPE (Common Platform Enumeration) 資訊，使其除了提供 CVE 描述、CVSS、EPSS 等基本資料外，還能輔助 AI 代理透過軟體版本號進行精準的漏洞篩選。FIND-CVE-MCP 負責從 CVE 官方 GitHub 專案庫中，利用 CPE 過濾出影響特定 Windows BuildVersion 的漏洞清單。為提升龐大資料的查詢效率，此 Server 採用了時間區間篩選邏輯，以目標 BuildVersion 的發行年份為基準回溯五年進行搜尋，並在每次過濾前更新資料庫。FIND-PoC-MCP 則從非官方 GitHub 專案庫蒐集特定 CVE 的 PoC 清單，以確認可利用性。

3.4 任務執行

MCP Server 端接收到來自 Cline 的參數與指令後，在完全隔離的虛擬機環境中執行指定操作。此封閉環境具備高度可監控性，所有執行過程均可完整記錄與回溯，同時提供安全性保障，避免潛在資安風險外溢至真實作業環境。

3.5 監控執行結果

任務執行完成後，其結果將作為 MCP Server 的輸出，透過 Cline 回傳至 LLM。LLM 根據回傳內容判斷操作是否成功；若成功，規劃下一步任務並徵詢使用者是否繼續；若失敗則分析錯誤原因，調整策略後重試。整體流程具備自我修正機制，使用者亦可即時監控，並於任務偏離預期時介入，確保整體操作的方向與安全性。

3.6 成功重現 CVE 判斷與評估設計

在本研究的架構下，CVE 是否成功重現需具備 PoC 可執行性及流程完整性。由 LLM 生成或取得的 PoC 程式碼能在虛擬機內順利編譯執行；且結果能引發與該 CVE 描述相符的現象，例如系統錯誤訊息、程式崩潰、權限異常等。而即使未觸發漏洞，只要 LLM 能正確產生具體的重現步驟與程式碼，即視為達到最低限度的完成。

藉由上述兩項標準評估不同情境下的表現。若滿足可執行性，即判定為重現成功；若滿足完整性，則視為部分重現；若都不滿足則為重現失敗。

4. 實驗

本實驗主要在於驗證 MCP 與 Cline 框架於 CVE 自動化重現的可行性，並於虛擬環境確保流程安全進行。實驗流程逐步檢驗是否能滿足「PoC 可執行性」與「流程完整性」兩項核心判斷標準。以下將說明具體的環境配置與實驗工具。

4.1 實驗環境

本研究的實驗環境建置於 Windows 10 (10.0.19045.5440) 64 位元虛擬機，並安裝完整的開發工具鏈，以支援 CVE 自動化重現的流程。開發環境包含 Visual Studio Code 與其 Cline 套件 (版本 3.26.5)、Python 3.13.0、Git 2.50.1、Node.js 22.17.1 與 MinGW-w64 編譯器，確保程式撰寫、版本管理與編譯執行皆能順利完成。

在框架設計上，LLM 扮演「大腦」角色，須能正確理解上下文並調用人工提供的工具。因此，本研究選擇具備高效能與低延遲回應特性的 GPT-4.1 mini 作為主要模型，以確保能有效支援漏洞重現過程中的互動與推理。

研究者在虛擬環境透過 VSCode 使用 Cline，並新增規則輔助任務運作。規則中包含身份提示、環境資訊、策略以及注意事項，使 Cline 的策略更符合需求。亦透過新增 Rules 輔助任務運作，內容涵蓋特殊身分提示、系統環境資訊、重現策略以及注意事項，使其更符合 CVE 重現的需求。

最後，為模擬實際攻擊情境中常見的防護繞過條件，實驗期間已預先停用 Windows Defender 防護機制，以避免干擾漏洞觸發與 exploit 行為之觀察與驗證。

4.2 Case Study

本研究選擇含公開 PoC 的 CVE 為重現對象，因 PoC 可視為可利用性的重要依據。由於 Windows 系統漏洞眾多，若無具體參考資訊，將無法確認 CVE 是否能在環境中被觸發。因此，我們挑選具備公開 PoC 的項目，以提升效率與成功率。

為確保所選 CVE 適合在單一 Windows 虛擬機中進行實驗與重現，本研究依據以下三項規則進行過濾：首先檢查版本相符性，僅挑選影響 Windows 10 Build Version 10.0.19045.5440 且經 MC 系統分析與 GitHub 蒐集後，至少具備一份公開 PoC 的漏洞；其次為可獨立執行性，優先選擇可於單一裝置本地執行，以利在隔離環境下進行完整重現；最後是真實可用性，人工審查每個 PoC，僅納入具備 exploit 潛力的程式；若 PoC 實際為漏洞掃描器、分析腳本或 patch 測試程式則排除，避免干擾結果判讀。

本研究從 17 個具公開 PoC 的 CVE 中，人工挑選出 5 個最具代表性且具備重現潛力的案例，並於 4.3.1 至 4.3.5 各節進行個別分析。

4.3.1 CVE-2025-21420

此漏洞為 DLL 側載 (side-loading) 本地提權漏洞。其原理為，當使用者或系統觸發 cleanmgr.exe (Windows Disk Cleanup Tool) 時，該程式會以 SYSTEM 權限執行，並載入多個 DLL。然而，由於 cleanmgr.exe 在設計上未明確指定所載入 DLL 的路徑。攻擊者可利用系統的預設搜尋順序，藉由符號連結 (Symbolic Link) 將惡意 DLL 放置於使可寫入目錄中。一旦 cleanmgr.exe 誤載此惡意 DLL，即可讓惡意程式碼以 SYSTEM 權限執行。

然而，實驗因權限問題而失敗。由於建立符號連結本身就需要管理員權限，這與低權限攻擊的情境相悖，導致流程在初期便中斷。即使以人工方式進行相同操作也未能成功觸發惡意 DLL 的載入，最終僅顯示「權限不足」或「找不到路徑」等錯誤訊息，未能達成完整的提權。

上述失敗的解決方案應考慮加入更明確的惡意觸發行為，例如由程式自行啟動高權限的終端介面 (如 cmd 或 powershell)。同時也需更深入分析 cleanmgr.exe 的執行過程是否可載入其他可被劫持的檔案或模組，以建立更完整的提權攻擊鏈。

4.3.2 CVE-2025-26633

此漏洞為一項涉及微軟管理控制台檔案格式 (.msc) 的安全功能繞過 (Security Feature Bypass, SFB) 漏洞。攻擊者利用微軟管理控制台 (Microsoft Management Console, MMC) 內的國際化路徑機制 (MUIPath)，將惡意 .msc 檔案放置於 en-US 目錄。使 MMC 優先載入該檔案，取代原本的合法版本。攻擊者可在這個惡意的 .msc 檔案內嵌入 PowerShell 指令或 ActiveX 控制項等元素，達到遠端命令執行或下載惡意酬載等目的。

然而，本次漏洞重現因建立的 .msc 檔案格式錯誤而失敗。MMC 主控台檔案採用的是複合結構化儲存檔 (Compound Structured Storage File)，而非單純的 XML 檔案，導致自製檔案被拒絕開啟。研究亦發現環境變數 MUILanguage 僅影響介面語言，與 .msc 的載入流程無關。此漏洞觸發的關鍵在於主

控台的內部結構與資源載入行為，並非單純在.msc檔案中插入腳本觸發，因此未能達到預期的效果。

為解決此問題，未來的研究方向應先透過 MMC 建立並儲存一個合法的 .msc 檔案，以確保其結構符合複合檔規範，再逐步修改其中的參數或資源，並觀察其對載入流程的影響。同時，建議使用 Process Monitor 等工具即時監控 mmc.exe 的行為，藉由比對修補前後其搜尋順序或安全策略上的差異，更深入理解此漏洞的運作機制。

4.3.3 CVE-2025-29824

此漏洞 Windows CLFS (Common Log File System) 核心驅動的 Use-After-Free 本地提權漏洞。攻擊者在觸發漏洞後，可覆寫目標程序的存取權杖 (access token)，從而取得 SYSTEM 等級最高權限，進而轉儲 LSASS 程序的記憶體竊取使用者憑證。

然而，在本次重現 exploit 程式未能正確開啟一個已啟用 CLFS 功能的有效路徑而失敗，這導致核心中關鍵的 Use-After-Free 觸發點未能生效。

為了解決此路徑問題，未來可利用 PowerShell 或系統搜尋工具，搜尋現有的 .blf 檔案 (例如位於 C:\System Volume Information 的檔案) 並加以利用，或嘗試建立相似結構的檔案來觀察驅動反應，同時也測試不同路徑格式以驗證能否繞過路徑解析上的限制。程式邏輯則需納入更精細的記憶體與權限控制，例如模擬存取權杖或調整行程權限，並透過檢查特定權限確認提權是否成功。若要進行更深入的分析，則必須配合 WinDbg 等偵錯工具監控 CLFS 驅動的實際行為。

4.3.4 CVE-2025-48799

此漏洞源於 Windows Update Service 在處理檔案時未驗證符號連結的目標是否合法。攻擊者可利用 link-following 技術，建立指向受保護的系統檔案或資料夾的惡意連結，誘使服務以 SYSTEM 權限對目標執行未授權的修改或程式碼。

然而，漏洞重現因程式邏輯瑕疵而失敗。程式雖會檢查系統是否具備兩顆硬碟 (如 C: 與 D:)，卻缺乏相應的錯誤處理機制。即使在單一硬碟環境下，程式仍會繼續執行，最終產生「權限不足」或「找不到路徑」等誤導性錯誤，這不僅掩蓋了硬碟數量不足的根本原因，也讓問題難以被準確診斷。

我們推論該漏洞可能仰賴系統中至少存在兩顆實體或虛擬硬碟來觸發漏洞行為。因此，下一步應在環境中嘗試新增一顆虛擬磁碟，並掛載為 D 槽。此設定預期能促使服務錯誤地跟隨連結，若能觀察到系統資料夾內的檔案被異常刪除等現象，便更接近成功重現此漏洞。

4.3.5 CVE-2025-49667

此漏洞發生於 Windows Win32K 子系統中的

ICOMP 函數，此函數負責影像壓縮與解壓縮。由於其影像壓縮流程可被觸發 Double Free 記憶體錯誤。攻擊者能透過特定 win32k.sys 系統呼叫，強制核心物件重複釋放，再利用受控資料重新分配該記憶體以覆寫核心函數指標，最終取得 SYSTEM 權限。

然而，漏洞重現因 LLM 的分析錯誤而失敗。它誤判漏洞需要使用者互動，但公開的 PoC 證實其為非互動式的強制雙重釋放。此外，LLM 生成的程式試圖以高權限執行，此舉與漏洞需從低權限帳戶發動提權的目的相悖，邏輯上並不成立。

本漏洞需特定的記憶體佈局，因此可嘗試 heap grooming (堆積配置控制) 或 object spraying (物件噴灑) 等技術，讓可控資料填補釋放後的記憶體區塊。關鍵在於，攻擊流程必須由低權限使用者啟動以符合實際場景。此外，目前使用的 API 或觸發流程並非真正對應 win32k.sys 中可被利用的函式。若仍無法確定漏洞觸發條件，需進一步透過 WinDbg 監控系統呼叫與記憶體變化，釐清雙重釋放的時機與位置，進而調整重現邏輯。

4.3.6 CVE 重現結果總結

本研究針對五個具公開 PoC 的 CVE 進行實驗重現，雖皆未完全觸發預期的漏洞利用行為，但部分案例的程式有成功執行、觸發異常或建立初步觸發點等。表 2 整理各 CVE 的最終重現結果。

表 2 CVE 重現結果總結

CVE 編號	可執行性	完整性	評估分類
CVE-2025-21420	否	是	部分成功
CVE-2025-26633	否	是	部分成功
CVE-2025-29824	否	是	部分成功
CVE-2025-48799	否	是	部分成功
CVE-2025-49667	否	是	部分成功

在漏洞類型與限制條件的差異方面，本研選取的案例成功重現的過程遠比漏洞描述複雜，並非僅透過命令列介面即可達成。例如，涉及 DLL 側載的漏洞需要利用特定系統工具，而其他漏洞則可能仰賴多硬碟環境才能觸發。這些特定條件若未能滿足，即便存在公開 PoC，也難以重現漏洞。

其次，部分開源 PoC 程式僅具示範性質，並非真正可利用的 exploit。部分程式僅實作漏洞偵測或補丁驗證，缺乏完整的提權鏈設計，導致即使執行也無法觸發預期的惡意行為。此現象呼應本研究的篩選規則，即只有「可被利用的 PoC」才能對重現研究帶來實質價值。

最後在流程層面，代理能夠提供合理的攻擊重現步驟，例如建議建立符號連結、利用 mmc.msc 載入行為或定位 CLFS 日誌檔案，並能與使用者互動確認執行，顯示其在「知識性說明」與「程序性輔

助」上的潛力。然而，LLM 卻難以直接生成完整且可執行的 exploit 程式。推測其原因在於 LLM 的訓練資料中涉及惡意程式或真實 exploit 專案的範例相對稀少，使其在權杖操作、核心物件配置等關鍵細節上出錯。因此，即便提供範例與建議流程，LLM 的輸出仍多半停留在理論層級，與實務可用的程式存在顯著落差。

5. 討論與結論

5.1 研究範圍與限制

本研究實驗皆於 Windows 環境進行，未涵蓋其他作業系統。CVE 資訊來自官方 CVE List，並透過 CPE Match 過濾出條件相符的漏洞；PoC 程式碼取自開源資料庫。因此僅涵蓋已公開且具備現存 PoC 的漏洞，不涉及零日漏洞等攻擊。

在測試範圍上，僅針對可透過命令列與程式碼執行的漏洞進行驗證，不含 GUI 互動情境，且限於單一主機，不涉及網路攻擊。此外，研究也受工具功能限制，例如 Cline 在 GUI 操作及多步驟任務穩定性上存在不足，同時 MCP Server 的數量與完整度亦可能影響測試結果，部分功能需仰賴自製的 Server 方能實現。

5.2 結論

本研究結果顯示：即使挑選具公開 PoC 的 CVE 作為重現對象，重現的成功率仍不高，失敗原因主要歸結於三點：首先是嚴格的环境前置條件未能滿足，如特定的磁碟配置或權限設定；其次是許多公開 PoC 本身偏向偵測或概念驗證而非完整攻擊；最後則是 LLM 在生成可執行的 exploit code 方面存在困難。在本研究架構中，LLM 扮演驅動整體流程的「大腦」角色，負責從資訊蒐集、流程規劃到程式生成的決策。然而，其表現出的限制也凸顯了通用模型的一大挑戰：預訓練資料普遍缺乏真實 exploit 程式碼與系統底層操作的知識，導致其生成的程式碼多停留在理論示範層級，難以觸發實際漏洞。

為克服這些挑戰並提升未來重現成功率，本研究提出兩個主要改進方向。在系統層面應持續擴充 MCP server 模組，整合更多樣的自動化檢測功能。該模組需能動態分析目標虛擬機的系統組態（如 Windows 版本、已安裝補丁、用戶權限等），以便在實驗前預先評估重現潛力並提供環境修正建議。在模型層面，則建議針對研究型 LLM 進行定向微調，透過導入具研究價值的漏洞技術分析、歷史 exploit 範例與系統 API 行為等資料，彌補其在特定安全領域的知識缺口，從而強化其生成可用 PoC 的能力。

總體而言，本研究驗證了 LLM 驅動的 AI 代理在漏洞重現自動化上的可行性。儘管當前仍面臨諸多限制，但若能持續改進 MCP 支援框架並微調

專用模型，未來將有望建立一個成功率更高、具備可重複性與自動化潛力的漏洞驗證系統。同時，本研究亦遵循 AI 在網路安全領域的倫理規範，所有實驗皆於隔離環境中進行，僅供學術研究之用。

致謝

本研究結果承蒙國科會計畫編號：NSTC 112-2222-E-011-011-MY2 與計畫編號：NSTC 114-2634-F-001-001-MBK 的支持。另特別感謝郭英仁、鄭仔孜、黃凱賢在研究過程中的討論與技術支援。

參考文獻

- [1] Yang, J., Jimenez, C. E., Wettig, A., Lieret, K., Yao, S., Narasimhan, K., & Press, O. (2024). SWE-agent: Agent-computer interfaces enable automated software engineering. *Advances in Neural Information Processing Systems*, 37, 50528-50652.
- [2] SWE-agent. (2025, May 22). GitHub. <https://github.com/SWE-agent/SWE-agent>
- [3] Fang, R., Bindu, R., Gupta, A., Zhan, Q., & Kang, D. (2024). Llm agents can autonomously hack websites. *arXiv preprint arXiv:2402.06664*.
- [4] Fang, R., Bindu, R., Gupta, A., & Kang, D. (2024). Llm agents can autonomously exploit one-day vulnerabilities. *arXiv preprint arXiv:2404.08144*.
- [5] Deng, G., Liu, Y., Mayoral-Vilches, V., Liu, P., Li, Y., Xu, Y., ... & Rass, S. (2024). {PentestGPT}: Evaluating and harnessing large language models for automated penetration testing. In *33rd USENIX Security Symposium (USENIX Security 24)* (pp. 847-864).
- [6] What is Cline? - Cline. (2025). Cline.bot; Cline. <https://docs.cline.bot/getting-started/what-is-cline>
- [7] OpenAI. (2025). Introducing GPT-4.1 in the API. Openai.com. <https://openai.com/index/gpt-4-1/>
- [8] Anthropic. (2024, November 25). Introducing the Model Context Protocol (MCP). Anthropic.com. <https://www.anthropic.com/news/model-context-protocol>
- [9] Avgerinos, T., Cha, S. K., Rebert, A., Schwartz, E. J., Woo, M., & Brumley, D. (2014). Automatic exploit generation. *Communications of the ACM*, 57(2), 74-84.
- [10] Cha, S. K., Avgerinos, T., Rebert, A., & Brumley, D. (2012, May). Unleashing mayhem on binary code. In *2012 IEEE Symposium on Security and Privacy* (pp. 380-394). IEEE.
- [11] Defense Advanced Research Projects Agency (DARPA). (2016, August 4). "Mayhem" declared preliminary winner of historic Cyber Grand Challenge. DARPA. <https://www.darpa.mil/news-events/2016-08-04a>
- [12] Greenberg, A. (2016, August 5). Hackers don't have to be human anymore. This bot battle proves it. Wired. <https://www.wired.com/2016/08/hackers-dont-human-anymore-bot-battle-proves/>