

## Estudio sobre agentes de inteligencia artificial aplicados a la ciberseguridad y pruebas de penetración en sistemas

Study on Artificial Intelligence Agents Applied to Cybersecurity and Penetration Testing in Systems

Gerardo **Cruz Espinosa**<sup>1</sup>  
Mariana **Velasco Trejo**<sup>2</sup>

Instituto Politécnico Nacional,  
<sup>1</sup> Escuela Superior de Ingeniería Mecánica y Eléctrica,  
Ciudad de México, MÉXICO  
ORCID: 0009-0005-4054-8702 / geecr07@gmail.com

Instituto Politécnico Nacional,  
<sup>2</sup> Escuela Superior de Cómputo,  
Ciudad de México, MÉXICO  
ORCID: 0009-0006-5522-9250 / mvelascot1601@alumno.ipn.mx

<https://cientifica.site>

Recibido 13/03/2026, aceptado 05/06/2026, publicado 18/06/2026.



## Resumen

La presente investigación analiza el funcionamiento técnico y la implementación de agentes de inteligencia artificial orientados a dominios específicos, particularmente en el ámbito de la ciberseguridad desde un enfoque ofensivo. Estos agentes permiten un cierto grado de automatización en ejecución de pruebas de penetración en sistemas informáticos, ya que, en la actualidad, este tipo de pruebas se realizan con herramientas operadas por profesionales de ciberseguridad.

El objetivo del estudio es examinar su arquitectura y capacidades operativas. Las pruebas realizadas evidencian su capacidad para procesar información contextual, coordinar el uso de herramientas y ejecutar acciones mediante el uso de modelos de lenguaje de gran tamaño. Los resultados sugieren que estas nuevas tecnologías pueden apoyar a los profesionales de ciberseguridad en diversas tareas operativas que actualmente se realizan de manera manual.

La metodología empleada es estructurada y se desarrolla por etapas, integrando la adquisición de información, el análisis contextual, la planificación de acciones, la ejecución controlada y la retroalimentación continua.

**Palabras clave:** ciberseguridad, agentes de inteligencia artificial, ciberseguridad ofensiva, pruebas de penetración, modelos de lenguaje grande, automatización.

## Abstract

This investigation analyzes the technical functioning and implementation of artificial intelligence agents oriented toward specific domains, particularly in the field of cybersecurity from an offensive perspective. These agents enable a certain degree of automation in the execution of penetration tests on computer systems due to this type of proofs are made through tools operated by cybersecurity professionals.

The objective of this study is to examine their architecture and operational capabilities. The conducted tests demonstrate their ability to reason and perform actions through the use of large language models. The results suggest that these emerging technologies can support cybersecurity professionals in various operation tasks, it is important to note that these tasks are already being done manually.

The methodology employed is structured and developed in stages, integrating information acquisition, contextual analysis, action planning, controlled execution, and continuous feedback.

**Index terms:** cybersecurity, artificial intelligence agents, offensive cybersecurity, penetration testing, large language models.

## I. INTRODUCCIÓN

En el mundo actual, la tecnología y el internet se han convertido en herramientas indispensables en el día a día de la mayoría de personas. En paralelo, en el campo de la inteligencia artificial (IA) se han desarrollado modelos de aprendizaje automático y modelos de lenguaje a gran escala. De acuerdo con datos de Google, la percepción positiva hacia la inteligencia artificial ha aumentado recientemente [1].

En este sentido y con respecto al ámbito profesional, el surgimiento y auge de la tecnología en la nube, la IA, los *smartphones* y las aplicaciones móviles, han ampliado la superficie de ataque, lo cual es aprovechado por los ciberdelincuentes. Ejemplo de ello, y de acuerdo con datos del NIST [2], los ataques informáticos han ido al alza durante los últimos años. Actualmente, la ciberseguridad constituye un pilar fundamental para garantizar la integridad, disponibilidad y confidencialidad de la información dentro de los sistemas digitales [3].

Desde esta perspectiva, se busca minimizar el riesgo de ciberataques mediante pruebas de penetración, ya que estos permiten identificar vulnerabilidades en sistemas informáticos al analizar información como: CVEs, registros de red, configuraciones, datos obtenidos de herramientas de escaneo y respuestas de servicios. Sin embargo, el resultado de estas pruebas depende en gran medida de la experiencia del profesional que las lleva a cabo, así como del tiempo que se tenga disponible para realizarlo y de la capacidad que se tenga para analizar grandes volúmenes de información, lo cual podría ocasionar que no se identifique la totalidad de las vulnerabilidades del sistema durante la prueba, dejando oportunidades de explotación para los ciberdelincuentes.

3

Con el surgimiento de los Modelos de Lenguaje Grande (LLM, por sus siglas en inglés), se ha observado que estos permiten analizar grandes cantidades de información y generar respuestas fundamentadas en patrones identificados durante su entrenamiento. Asimismo, cuando forman parte de arquitecturas basadas en agentes, pueden apoyar la selección de acciones a ejecutar de acuerdo con los objetivos definidos y la información disponible. Lo anterior es relevante debido a que, hasta hace poco tiempo, los modelos de lenguaje se utilizaban principalmente para tareas de generación y procesamiento de texto. En años recientes, han surgido agentes de IA capaces de interactuar con herramientas externas y ejecutar acciones dentro de entornos controlados.

En el ámbito de la ciberseguridad, los agentes de IA podrían permitir llevar a cabo tanto tareas ofensivas como defensivas, lo que suele realizarse por *Red Teams* y *Blue Teams*, respectivamente. Dentro del *Red Team*, se encuentra un especialista conocido como *Pentester*, el cual es el encargado de realizar las pruebas de penetración en sistemas, estas pruebas también son conocidas como *pentests*. Un agente de IA podría asistir en un *pentest* al realizar una serie de tareas, como lo son: reconocimiento, escaneo, evaluación de vulnerabilidades, explotación y generación de informes [4], haciendo uso de herramientas que los *Pentesters* suelen emplear. Con respecto al *Blue Team*, esta tecnología podría contribuir en la detección de Indicadores de Compromiso (IoCs), detección y correlación de eventos, así como respuesta ante incidentes [5].

Cabe destacar que esta tecnología permite generar una retroalimentación con base en los resultados obtenidos en cada prueba, lo que posibilita que el modelo itere hasta que determine que no es necesario continuar. Asimismo, puede aprender patrones a lo largo del tiempo, lo que les permite que mejoren sus resultados en la medida que sea utilizada para que pueda generar aprendizaje.

## II. REVISIÓN DE LA LITERATURA

### A. Inteligencia artificial

La Inteligencia Artificial (IA) es un área de las ciencias computacionales que tiene el objetivo de crear sistemas con capacidades suficientes para realizar tareas que tradicionalmente requieren inteligencia y razonamiento humano. Estos sistemas adquieren la capacidad de generalizar patrones a partir de datos previamente observados y de la adaptación a entornos cambiantes [6].

Los primeros desarrollos en IA estuvieron basados en sistemas expertos y modelos fundamentados en reglas lógicas explícitas, las cuales, junto con los datos, eran proporcionadas por programadores especialistas. Estos sistemas se ejecutaban en entornos controlados; sin embargo, presentaban importantes limitaciones. Además, resultaba inviable supervisar y corregir constantemente la información que el sistema generaba. Por lo que la necesidad de reducir, en la medida de lo posible, la intervención humana, así como la recopilación de grandes volúmenes de información y los avances en el estudio del lenguaje natural, dirigieron a la inteligencia artificial hacia enfoques probabilísticos y estadísticos [7].

Adicionalmente, el avance de capacidad en el ámbito computacional permitió realizar el almacenamiento y procesamiento de grandes volúmenes de datos. En este contexto, el aprendizaje profundo, aprendizaje automático, análisis masivo de información y procesamiento del lenguaje natural, se consolidaron como pilares de la IA moderna, haciendo posible que los sistemas aprendieran patrones y comportamientos a partir de grandes volúmenes de datos.

### B. Modelos de lenguaje y razonamiento (LLM)

Los LLM se fundamentan en técnicas de Procesamiento de Lenguaje Natural (PLN), disciplina de la inteligencia artificial orientada al análisis y generación de lenguaje humano considerando el contexto y significado de la información procesada [8].

Mientras se realizaban estudios referentes al PLN, se identificó un problema: los textos estaban siendo procesados de manera secuencial, lo cual causaba dificultad para recordar un contexto muy grande, lentitud en el entrenamiento y problemas con dependencias largas. Ante esta problemática, un grupo de investigadores crearon el concepto de *Transformer* [8] (Fig. 1).

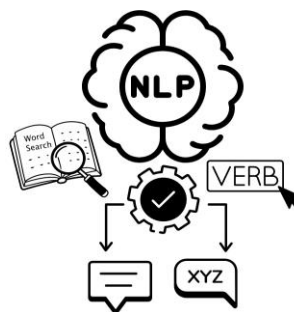


Fig. 1. Procesamiento de lenguaje natural.

Un *Transformer* se refiere a la arquitectura neuronal que permitió que los PLMs modernos funcionaran a la perfección, ya que esta arquitectura observa una oración completa al mismo tiempo y decide cuáles son las partes importantes, es decir, ve una frase como un todo y no como un conjunto ordenado de palabras individuales. Cabe mencionar que existen tres tipos de *Transformers* (Fig. 2):

- **Encoder.** Entiende texto al procesar una entrada haciendo uso de la capa de atención.
- **Decoder.** Genera texto de acuerdo con el objetivo del modelo gracias a que toma la representación realizada por el *encoder*.
- **Encoder-Decoder.** Entiende y genera texto.

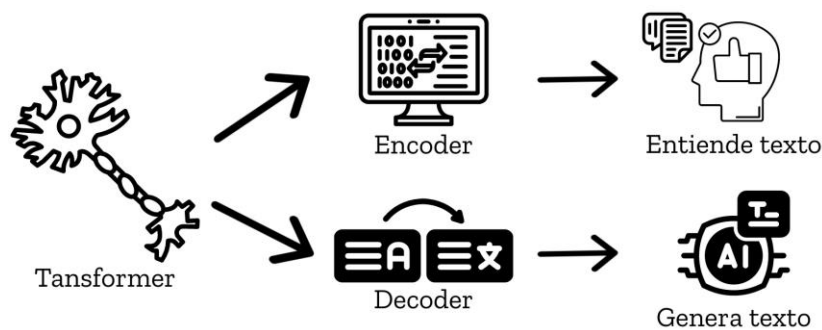


Fig. 2. Funciones de los tipos de Transformers.

Con la creación de los *Transformers*, surgieron los modelos, que son modelos matemáticos (red neuronal tipo *Transformer*) que aprenden patrones estadísticos. Aunque, un modelo no razona, sí predice, y al hacerlo millones de veces, acierta la mayor parte del tiempo.

Posteriormente y gracias a este avance, se crearon los Modelos de Lenguaje Grande (*Large Language Models*, LLM, por sus siglas en inglés), mismos que consisten en modelos de IA que tienen la finalidad de comprender contexto, generar texto y resolver tareas complejas mediante patrones aprendidos durante el entrenamiento, así como predecir la siguiente palabra en una secuencia (Fig. 3).

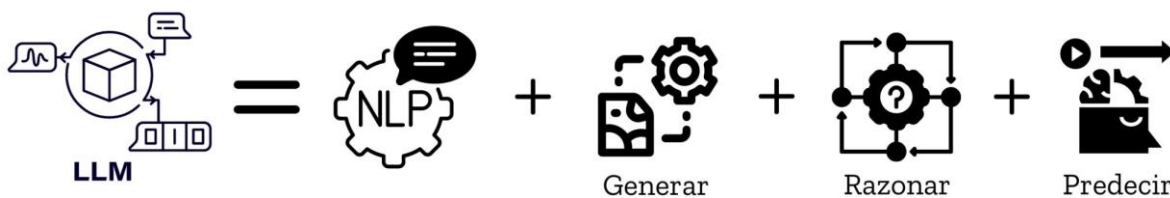


Fig. 3. Estructura de un LLM.

Una vez finalizado su proceso de entrenamiento, el modelo entra en la parte llamada inferencia, misma en la que aplica los parámetros aprendidos para generar texto a partir de una entrada, la cual es denominada como *prompt* (instrucción de entrada) en el contexto de los *chatbots*. Durante esta fase, el mecanismo de atención desempeña un papel fundamental, ya que permite comprender el contexto y generar respuestas coherentes. Es importante destacar que los LLM no son capaces de interactuar con su entorno por sí solos.

Por otra parte, además del proceso de inferencia, existen otros métodos que permiten ajustar o complementar el conocimiento de un modelo previamente entrenado. Uno de dichos métodos es el ajuste fino (*fine tuning*), técnica que consiste en continuar el entrenamiento del modelo sobre un conjunto de datos específico con el objetivo de modificar sus parámetros internos y especializarlo en una tarea u objetivo específico. Sin embargo, reentrenar modelos a gran escala requiere de una considerable capacidad computacional y recursos económicos, por lo que su desarrollo suele ser realizado por grandes empresas.

Ante dicha problemática, una alternativa es el uso de la técnica conocida como *Retrieval-Augmented Generation* (RAG), ya que, a diferencia del ajuste fino, esta no altera los parámetros del modelo, sino que incorpora información externa durante la inferencia por medio de la recuperación e incorporación de documentos relevantes. De esta forma, el modelo es capaz de generar respuestas fundamentadas en información adicional sin necesidad de modificar su entrenamiento original (Fig. 4).

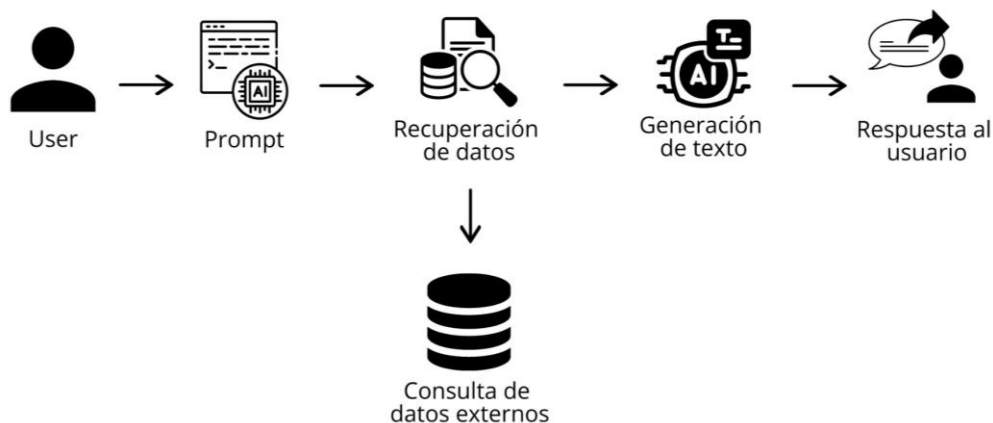


Fig. 4. Funcionamiento de la estructura RAG.

Las capacidades de razonamiento, generación de contenido y acceso a conocimiento externo han permitido que los LLM trasciendan las aplicaciones tradicionales de procesamiento de lenguaje natural. En años recientes, estos modelos han comenzado a incorporarse en dominios especializados como la ciberseguridad, donde apoyan actividades relacionadas con la detección de amenazas, análisis de vulnerabilidades, automatización de tareas operativas y asistencia en pruebas de penetración. Este avance ha impulsado el desarrollo de arquitecturas más complejas basadas en agentes inteligentes capaces de interactuar con herramientas y sistemas externos.

### C. Agentes inteligentes

Un agente inteligente es un sistema basado en modelos de IA capaz de percibir información de su entorno, procesarla y ejecutar acciones mediante herramientas externas con un grado de autonomía. En este contexto, el modelo de lenguaje actúa como componente de decisión dentro de una arquitectura orientada al cumplimiento de objetivos específicos [9].

Dentro de las arquitecturas basadas en herramientas, el modelo de lenguaje selecciona las acciones más adecuadas para cumplir un objetivo y delega su ejecución a componentes especializados (Fig. 5).

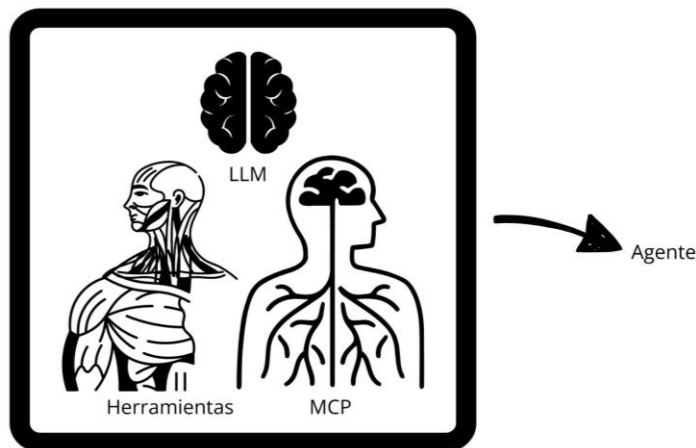


Fig. 5. Estructura de un agente.

Uno de los enfoques más utilizados para coordinar este comportamiento es ReAct (Reasoning + Acting), técnica que estructura la resolución de tareas mediante un ciclo iterativo compuesto por pensamiento, acción y observación [10] (Fig. 6).

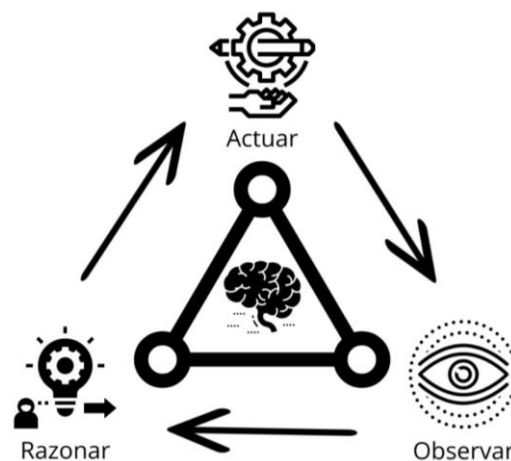


Fig. 6. Funcionamiento del enfoque ReAct.

La combinación de modelos de lenguaje, mecanismos de razonamiento y herramientas externas ha convertido a los agentes inteligentes en una de las principales líneas de investigación actuales. Particularmente en ciberseguridad, estos sistemas han demostrado potencial para automatizar procesos de análisis, correlación de eventos, generación de consultas, ejecución de herramientas especializadas y apoyo a la toma de decisiones, lo que ha motivado el desarrollo de soluciones orientadas tanto a actividades ofensivas como defensivas.

#### D. Ciberseguridad ofensiva y defensiva (*Red Team/Blue Team*)

La ciberseguridad en la industria moderna se estructura en torno a capacidades ofensivas y defensivas. Para describir estas funciones, se han adoptado los términos *Red Team* y *Blue Team*, los cuales se refieren a equipos que se encargan de simular ataques, así como a los profesionales responsables de la defensa y protección de sistemas (Fig. 7).

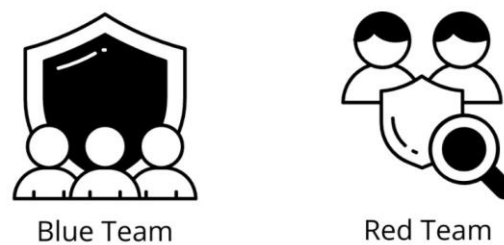


Fig. 7. Equipos de ciberseguridad ofensiva y defensiva.

8

El *Red Team* se centra en ejecutar las Tácticas, Técnicas y Procedimientos (TTPs) que suele utilizar un adversario real con el objetivo de identificar fallas de configuración, vulnerabilidades, código vulnerable y deficiencias en los controles de seguridad empleados [11]. Para estructurar y estandarizar la descripción de dichas TTPs, se emplean marcos de referencia, como lo es MITRE ATT&CK, el cual permite modelar amenazas y proporciona una taxonomía estructurada de las tácticas y técnicas observadas en ataques reales.

Adicionalmente, las pruebas de penetración permiten a las empresas priorizar vulnerabilidades de alto impacto y facilitar su mitigación. Los profesionales que las desarrollan utilizan clasificaciones de riesgos, uno de ellos es el OWASP Top 10, este es un documento de referencia que describe las principales vulnerabilidades en aplicaciones *web*, es importante destacar que este es actualizado periódicamente para adaptarlo a nuevas amenazas [13].

El estudio de estos patrones ha permitido estructurar los ataques en fases que describen la evolución del compromiso, entre las principales fases se encuentran: acceso inicial, ejecución, persistencia, elevación de privilegios, movimientos laterales y filtración de información, entre otras [12].

Por otro lado, el *Blue Team* está conformado por diferentes equipos que se encargan de realizar tareas específicas, mismas que se describen a continuación:

- Gestionar incidentes y dar respuesta a ellos, por lo que otro referente relevante es el NIST, el cual proporciona mejores prácticas y describe la capacidad de respuesta ante incidentes, la detección de anomalías y los mecanismos para reducir el impacto de los ataques. De acuerdo con el NIST, los pasos a realizar para la gestión de un incidente son: preparación, detección y análisis, contención, erradicación, recuperación y actividades posteriores al incidente [14].

- Monitorear las 24 horas del día y los 7 días de la semana la red de una empresa con el objetivo de identificar comportamiento inusual que sugiera posible actividad maliciosa. Cada vez que se detecta un comportamiento que sigue ciertos patrones, se genera un incidente. Para esta actividad, el NIST también proporciona mejores prácticas, las cuales consisten principalmente en combinar diferentes fuentes de datos, comportamientos y firmas para la detección de ataques.

Es fundamental identificar y responder de manera rápida y efectiva ante un incidente para evitar que un adversario logre comprometer un sistema [15]. Incorporar a la IA en los procesos para asegurar la integridad, confidencialidad y disponibilidad de la información representa una evolución natural que ya comienza a observarse en la actualidad: con la ayuda de LLM que verifican datos en tiempo real y responden en consecuencia, los agentes de IA aplicados a la ciberseguridad pueden apoyar tanto los procesos del *Red Team* como a los del *Blue Team*. Sin embargo, todavía existen múltiples áreas de oportunidad en esta iniciativa.

La convergencia entre agentes inteligentes, modelos de lenguaje y procesos de ciberseguridad ha dado lugar a nuevas soluciones capaces de asistir tareas tradicionalmente realizadas por especialistas. Como resultado, diversas organizaciones y grupos de investigación han desarrollado herramientas que integran capacidades de razonamiento, acceso a información contextual y ejecución de acciones, ampliando el nivel de automatización alcanzable en entornos de seguridad. En la siguiente sección se presentan algunas de las propuestas y plataformas más representativas identificadas en la literatura reciente.

## 9

### III. ESTADO DEL ARTE

En los últimos años, la automatización en el ámbito de la ciberseguridad ha evolucionado desde mecanismos basados en firmas, reglas predefinidas, *scanners* y *playbooks* estáticos hacia enfoques más avanzados que incorporan análisis conductual, correlación de eventos y modelos basados en IA.

Dentro del enfoque del *Blue Team*, existen herramientas conocidas como “Agentes en *endpoints*” (EDR/XDR). Estas herramientas son programas que se instalan en cada una de las computadoras que se desean proteger y que se comunican con un sistema central o “cerebro” [16]. Algunos de los datos que pueden recolectarse mediante estos agentes son: programas instalados con versiones vulnerables, comportamientos anómalos e indicadores de compromiso.

Algunos ejemplos claros y populares de dichos agentes son: Qualys [17], soluciones de *Endpoint Detection and Response* (EDR), Nessus que permite escanear redes, potenciando su funcionamiento mediante *plugins* para que sea utilizado en tecnologías específicas [18].

Asimismo, los agentes inteligentes se están empleando para clasificar eventos y patrones, así como correlacionarlos. Algunos ejemplos de este enfoque son los siguientes:

- Sec-PaLM, este es un modelo de lenguaje diseñado para apoyar a la seguridad en la nube de Google. Sec-PaLM brinda apoyo para la detección de indicadores de amenazas, comportamientos sospechosos, presencia de *malware* y existencia de vulnerabilidades [19].
- Microsoft Security Copilot, esta es una solución basada en IA que implementa agentes para simplificar tareas como la corrección y clasificación de vulnerabilidades, así como la detección y clasificación de ataques de *phishing*, entre muchas otras funciones [20].

- VulnBot, aplica la automatización tradicional de las pruebas de penetración apoyándose en herramientas que realizan escaneos, enumeración y, en algunos casos, explotación asistida. Una herramienta que propone utilizar agentes y modelos de lenguaje para automatizar pruebas [21].
- HexStrike, una herramienta que combina las capacidades de inferencia y planificación de un LLM con agentes capaces de ejecutar hasta 200 *software* diferentes de *pentesting*. Este proyecto puede integrarse con modelos ampliamente utilizados, como Cursor, ChatGPT o Gemini [22].

Una parte fundamental para el uso de los agentes inteligentes son los servidores MCP (*Model Context Protocol*) [23]. Este protocolo está diseñado para que los modelos de IA puedan conectarse a herramientas, datos y sistemas externos, de manera que MCP actúa como un puente estandarizado entre LLM, agentes inteligentes y aquellas herramientas disponibles para su ejecución (Fig. 8).



Fig. 8. Capas de un servidor MCP.

Además, este protocolo define la forma en que se realiza la comunicación y delimita su alcance. Cada acción que ejecuta un agente debe pasar por el servidor MCP. En otras palabras, un servidor MCP funciona como una capa intermedia entre el agente y las distintas herramientas (Fig. 9).

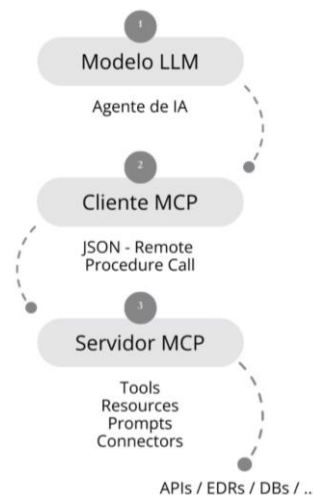


Fig. 9. Arquitectura de comunicación en MCP.

En el contexto de las pruebas de penetración, el servidor MCP habilita el acceso del agente a herramientas como: *nmap*, *sqlmap*, *Burp Suite*, *ffuf*, *amass*, *nuclei*, entre muchas otras que son ampliamente utilizadas en la industria.

Es importante destacar que, si bien estas herramientas permiten automatizar múltiples procesos, actualmente no es viable que estos reemplacen por completo al ser humano, solo deben ser considerados mecanismos de apoyo a la toma de decisiones. Los LLM, al ser modelos probabilísticos, no están exentos de generar errores, sesgos o interpretaciones no fundamentadas, por ello se requiere validar los resultados, descartar posibles falsos positivos y considerar el impacto asociado a su implementación.

Asimismo, la integración de estas tecnologías en entornos empresariales requiere considerar los riesgos de seguridad derivados de su uso. Entre estos se encuentran ataques como *prompt injection*, uso inadecuado de la herramienta, deficiencias en los controles de acceso y la necesidad de alinearse a marcos de referencia como el NIST. De manera que, su incorporación debe estar acompañada de políticas, mecanismos de supervisión y controles técnicos adecuados.

## IV. OBJETIVOS Y PREGUNTAS DE INVESTIGACIÓN

Una minuciosa revisión sobre la literatura actual acerca de la detección de ataques Evil Twin (ETA) muestra un panorama tecnológico disperso, en el cual los investigadores luchan entre la precisión algorítmica y la factibilidad operativa en redes reales. Se examinan a continuación las disparidades más importantes, los huecos tecnológicos no cubiertos y las proyecciones futuras en esta área.

11

### A. OBJETIVO GENERAL

Analizar la viabilidad del uso de agentes inteligentes basados en LLM e integrados mediante servidores MCP para llevar a cabo la ejecución de pruebas de penetración, identificando sus beneficios, limitaciones y posibles áreas de oportunidad.

### B. OBJETIVOS ESPECÍFICOS

1. Determinar en qué medida las herramientas que implementan agentes inteligentes cumplen con los marcos y estándares de ciberseguridad, tales como OWASP Top 10, MITRE ATT&CK y NIST.
2. Identificar y analizar las herramientas que utilizan agentes inteligentes para ejecutar pruebas de penetración automatizadas, descritas en el estado del arte.
3. Evaluar la contribución de los agentes de IA en la automatización de las pruebas de penetración, comparando sus capacidades con enfoques manuales y semi automatizados.
4. Realizar una comparativa de los resultados proporcionados por el agente mediante el uso de distintos LLM.

A partir de los objetivos previamente definidos, se plantean las siguientes preguntas de investigación:

1. ¿En qué actividades puede apoyar la automatización que ofrecen los agentes inteligentes y el razonamiento de los LLM en la ciberseguridad?
2. ¿Cuál es la implementación actual de estos agentes en la industria, dentro de soluciones comerciales?
3. ¿Cuáles son las ventajas, desventajas y limitaciones de realizar pruebas de penetración en sistemas utilizando estas nuevas herramientas?

Con base en las preguntas listadas anteriormente, se estructura un análisis desde diversas perspectivas, considerando su uso futuro en entornos empresariales.

## V. METODOLOGÍA

Una minuciosa revisión sobre la literatura actual acerca de la detección de ataques Evil Twin (ETA) muestra un panorama tecnológico disperso, en el cual los investigadores luchan entre la precisión algorítmica y la factibilidad operativa en redes reales. Se examinan a continuación las disparidades más importantes, los huecos tecnológicos no cubiertos y las proyecciones futuras en esta área.

La metodología utilizada en este trabajo se basa en un enfoque conceptual y operativo para el diseño y análisis de agentes de IA dedicados a pruebas de penetración.

Esta investigación se desarrolló en cinco etapas:

1. Adquisición de información
2. Análisis y razonamiento
3. Planificación de acciones
4. Ejecución controlada
5. Retroalimentación y aprendizaje

Estas etapas fueron definidas con base en el ciclo operativo comúnmente utilizado por agentes inteligentes orientados a la ejecución de tareas complejas, así como en procesos empleados durante pruebas de penetración. Esta estructura permite evaluar de manera sistemática la capacidad de los agentes para obtener información, procesarla, planificar acciones, interactuar con herramientas especializadas y ajustar su comportamiento a partir de los resultados obtenidos. Debido a que estas actividades forman parte de múltiples escenarios de evaluación de seguridad, la metodología puede ser reproducida y adaptada a diferentes entornos de pruebas de penetración asistidas por agentes de IA.

### A. Adquisición de información

Durante esta etapa se recopiló información sobre conceptos clave para el desarrollo de la presente investigación, como: PLN, modelos, *transformers*, LLM y agentes inteligentes, por mencionar algunos, ya que para la implementación de este agente inteligente primero se requería conocer y comprender el funcionamiento de los mismos para poder utilizarlo de manera correcta y obtener datos confiables.

Asimismo, se realizó una investigación acerca de herramientas disponibles públicamente que permiten la ejecución de *software* de *pentesting* a través de un LLM con la finalidad de conocer las funcionalidades ofrecidas por cada una de ellas.

La información recopilada en esta fase puede ser consultada en las secciones 'II. Revisión de la literatura' y III. 'Estado del arte'.

### B. Análisis y razonamiento

Luego de haber llevado a cabo la recopilación de información, se realizó una evaluación respecto a las herramientas identificadas que permiten la ejecución de *software* de *pentesting*, ya que se requería elegir una de ellas para poder implementarla y realizar las pruebas necesarias. Tras la evaluación de las mismas, se determinó que se haría uso de HexStrike, ya que, entre las herramientas identificadas, esta es la que proporciona un mayor número de *software* de *pentesting* y una arquitectura basada en agentes capaz de coordinar su ejecución. Esta característica resulta particularmente relevante para los objetivos del presente trabajo, ya que permite analizar el comportamiento de agentes de IA en escenarios que abarcan múltiples fases de una prueba de penetración. Mediante 12 agentes especializados,

HexStrike integra herramientas dedicadas a redes y reconocimiento, seguridad en aplicaciones *web*, autenticación y contraseñas, análisis de binarios, ingeniería inversa y seguridad en la nube, facilitando la evaluación de un amplio conjunto de actividades que comúnmente son realizadas por un *pentester*.

Si bien HexStrike fue utilizado como plataforma experimental debido a su amplia cobertura de herramientas de *pentesting*, la contribución de este trabajo se centra en analizar la viabilidad del uso de agentes inteligentes integrados mediante servidores MCP, así como en identificar sus beneficios, limitaciones y áreas de oportunidad dentro de escenarios de pruebas de penetración.

Para llevar a cabo las actividades anteriores, la arquitectura lógica de este agente sigue el enfoque ReAct y se encuentra conformado de la siguiente manera:

- **Módulo de adquisición de información.** Se encarga de recolectar datos sobre el sistema objetivo mediante herramientas de enumeración y reconocimiento.
- **Módulo de razonamiento y análisis.** Examina la información obtenida y define la superficie de ataque.
- **Módulo de planificación.** Selecciona las tácticas y los procedimientos que deben utilizarse.
- **Módulo de ejecución autónoma.** Se encarga de ejecutar las herramientas de *pentesting* con ayuda del servidor MCP.
- **Módulo de retroalimentación.** Su objetivo es que el modelo aprenda de las entradas que reciba y las respuestas que proporciona con el objetivo de mejorar las decisiones en futuras ejecuciones.

13

Esta arquitectura permite medir y observar el comportamiento, la eficacia y los resultados del modelo durante su ejecución. Para ello, se define un conjunto estructurado de reglas que el modelo incorpora en cada interacción con un usuario como parte del *prompt* o contexto de entrada. De esta manera, el modelo toma en consideración dichas instrucciones al generar la respuesta o al elegir la ejecución de alguna herramienta. Este enfoque no modifica los parámetros del modelo, sino que orienta su comportamiento en tiempo de inferencia mediante la incorporación dinámica de restricciones o criterios operativos.

Por otra parte, se realizó una evaluación respecto a qué LLM utilizar, luego de valorar las características de los más populares, se decidió hacer uso de Gemini, ChatGPT y Cursor.

### C. Planificación de acciones

Para la realización de las pruebas se utilizaron las siguientes herramientas y configuraciones:

- Cada agente requería ser desplegado en un entorno aislado y controlado: una máquina virtual en el *software* de virtualización VMware, asignándole 6 GB de memoria RAM. Este enfoque permite que el agente no pueda causar daños a la máquina física, lo que se alinea con los principios de diseño y contención utilizados en la industria de la ciberseguridad.
- Gemini como LLM encargado de analizar el contexto disponible, generar estrategias de acción y seleccionar los siguientes pasos a ejecutar.
- ChatGPT como LLM encargado de analizar la información disponible, planificar acciones y determinar los siguientes pasos dentro del flujo del trabajo.

- Cursor aprovechando sus capacidades de agente embebido para trabajar directamente sobre archivos del sistema.
- HexStrike como servidor MCP que expone herramientas de *hacking* ético para ser utilizadas por LLM mediante agentes.

Posteriormente, se definieron las acciones puntuales realizadas en los escenarios, las cuales consistieron en: escaneos de vulnerabilidades y servicios, búsqueda de puertos abiertos y, en algunos casos, la explotación de aplicaciones vulnerables.

Las pruebas fueron ejecutadas por dos personas:

- Un *pentester senior* con 5 años de experiencia en pruebas de penetración de sistemas *web*, aplicaciones móviles e infraestructura bancaria crítica, incluyendo plataformas SWIFT y otros sistemas de alta criticidad.
- Un analista de Inteligencia de Amenazas (Nivel 2) con 5 años de experiencia en ciberinteligencia, modelado de amenazas y análisis de actores, campañas e Indicadores de Compromiso (IoC).

Se llevaron a cabo las mismas actividades dentro del mismo alcance y bajo criterios de evaluación equivalentes, la diferencia radicó en que la analista utilizó el agente de inteligencia y el *pentester* llevó a cabo las actividades manualmente y con base en su experiencia.

Las fases consideradas para este *pentest* fueron:

1. *Reconnaissance*, el objetivo de esta es identificar superficie expuesta.
2. *Scanning*, en esta etapa se tiene el objetivo de identificar puertos y servicios.
3. *Enumeration*, durante esta etapa se extrae información sensible.
4. *Gaining access*, en esta es posible obtener RCE inicial.
5. *Lateral movement*, se busca realizar un *pivot* a usuario real del sistema.
6. *Privilege escalation*, en este punto el objetivo es escalar a *root*.
7. *Maintaining access*, la finalidad de esta fase es establecer persistencia (conceptual).
8. *Clearing tracks*, finalmente se busca borrar evidencia (teórico).

#### D. Ejecución controlada

En esta etapa, tanto la analista como el *pentester* llevaron a cabo las actividades de *pentesting* correspondientes a las fases definidas anteriormente, esto con ayuda de distintos LLM y siguiendo una secuencia estructurada de ataque. De manera paralela, el *pentester* realizó las mismas fases sin ayuda de ningún modelo de IA.

El agente de IA se implementó en un *host* con el sistema operativo Kali Linux, desde el cual se realizaron las actividades de evaluación. Además, tanto las pruebas realizadas por el agente como por el *pentester*, se realizaron en escenarios experimentales diseñados específicamente para ejercicios de *pentesting*, desplegados en entornos aislados con el fin de garantizar condiciones controladas durante la experimentación.

A continuación, se describen los distintos escenarios en los que compitieron el analista y el *pentester*, así como los resultados y conclusiones obtenidos.

Cabe destacar que las Tácticas, Técnicas y Procedimientos (TTP) utilizadas en este estudio están alineadas con marcos de referencia ampliamente utilizados en la industria, entre estos se encuentra MITRE ATT&CK y la metodología de

Certified Ethical Hacker (CEH) de EC-Council. Estas referencias proporcionan un conjunto estructurado de prácticas comúnmente utilizadas para la evaluación de la seguridad de sistemas y redes.

En todos los escenarios descritos a continuación, el laboratorio simula un servidor en una red empresarial donde una atacante obtuvo acceso no autorizado, por lo que se busca obtener control de este hasta lograr tener los privilegios elevados, mismos que en sistemas Linux se denominan como privilegios de *root*.

### 1) Escenario 1: PENTESTER

Este escenario describe las actividades realizadas por el *pentester* de manera manual, sin uso de agentes de IA.

#### a) Caso A: RECONOCIMIENTO

El *pentester* ejecutó un escaneo con la herramienta RustScan, de este se obtuvo como resultado la identificación del servidor vulnerable y los puertos TCP abiertos en la máquina remota.

#### b) Caso B: EXPLOTACIÓN

El *pentester* buscó y logró explotar la vulnerabilidad, las actividades que realizó para lograr este resultado fueron:

- Comprobar versiones mediante el *output* de la herramienta *nmap*.
- Buscar *exploits* en internet, principalmente en plataformas como GitHub.
- Modificar el *exploit*.
- Ejecutar y obtener acceso inicial.

#### c) Caso C: ESCALACIÓN DE PRIVILEGIOS

El *pentester* realizó las siguientes actividades para poder escalar privilegios dentro del entorno de prueba:

- Obtener los grupos a los que pertenece el usuario con sesión activa.
- Verificar el grupo de *sudoers* con el objetivo de averiguar si es posible ejecutar acciones con privilegios elevados.
- Buscar archivos de configuración.
- Enlistar los procesos del sistema.
- Identificar las tareas programadas corriendo en segundo plano.
- Identificar los procesos con permisos elevados.
- Realizar búsquedas en los *logs* del sistema con la finalidad de identificar información relevante.

En la Tabla 1, se muestra a manera de resumen, las actividades que el *pentester* llevó a cabo, los resultados obtenidos en cada una de ellas y el tiempo en minutos que requirió para realizarlas.

**TABLA 1. PENTESTER**

Fases	Herramienta	Técnicas usadas	Evidencias obtenidas	Resultado	Tiempo (min)
□ <i>Reconnaissance</i>	<i>RustScan</i>	Escaneo de puertos y servicios	Se obtuvieron versiones de protocolos HTTP, SSH.	Identificación de puertos abiertos y versiones de servicios.	8
□ <i>Scanning</i>	<i>BurpSuite Scammer</i>	Detección de vulnerabilidades en páginas web.	Se obtuvieron rutas escondidas versiones exactas del CMS.	Confirmación de tecnologías que la aplicación vulnerable implementa	5
□ <i>Enumeration</i>	<i>FFUF</i>	Búsqueda de directorios ocultos mediante “fuzzing”	Se obtuvieron credenciales del archivo settings.php.	Credenciales reutilizables encontradas.	15
□ <i>Gaining Access</i>	<i>Metasploit</i>	Inicio de sesión en en gestor de contenido + Subir archivos.	Se ejecutaron comandos mediante interpretación de archivos php.	<i>Ejecución remota de comandos. Usuario www-data.</i>	7
□ <i>Lateral Movement</i>	<i>SSH</i>	Probar las credenciales obtenidas en los servicios encontrados.	Se obtuvo acceso mediante SSH.	<i>Se conectó por ssh a la máquina remota, usuario obtenido: johncusack.</i>	6
□ <i>Privilege Escalation</i>	<i>ZSH</i>	Revisar permisos del comando <i>sudo -l</i> .	Se reconoció que el usuario <i>johncusack</i> puede ejecutar <i>/usr/local/bin/bee</i> como <i>root</i> .	RCE como <i>root</i> .	12
□ <i>Maintaining Access</i>	<i>penelope</i>	Obtener <i>shell</i> reversa estable	Se obtuvo una shell reversa y estable con la herramienta <i>penelope</i> .	<i>Shell como usuario privilegiado (root)</i>	5
□ <i>Clearing Tracks</i>	N/A	Remover archivos, limpiar <i>logs</i> .	No es necesario.	—	2

## 2) Escenario 1: LLM GEMINI

En este escenario fue realizado por la analista utilizando un agente de IA con Gemini implementado como LLM.

### a) Caso A: RECONOCIMIENTO

Para comenzar el ejercicio el analista escribió el siguiente *prompt* en el agente de IA: “Realiza una enumeración de puertos y servicios, recopila la información más relevante y enumera configuraciones débiles. Si se encuentra la manera de hacer una ejecución remota de comandos y posteriormente escalar privilegios no lo hagas, es solo enumeración”.

Como resultado del *prompt*, se identificó un servidor *web* que ejecutaba un gestor de contenido *Backdrop CMS* y un repositorio *.git* accesible públicamente [T1593.003].

### b) Caso B: EXPLOTACIÓN

Luego del escaneo, el analista utilizó el siguiente *prompt*: “Si se encuentra la manera de hacer una ejecución remota de comandos y, posteriormente, escalar privilegios, hazlo, solo intenta no tirar los sistemas, no realices ataques agresivos”. El agente basándose en la información obtenida en la fase anterior y en las instrucciones del analista, realizó los siguientes pasos:

- Analizar el *output* de las herramientas y clasificar los hallazgos. Intentar explotar las vulnerabilidades identificadas en el gestor de contenido.
- Alinear configuraciones con técnicas, tácticas y procedimientos que se describen en marcos de referencia como MITRE ATT&CK e intentar descubrir cómo escalar de los privilegios.
- Revisar otros posibles vectores de ataque y generar un reporte sencillo.

Durante estas actividades, se observó que el agente siguió las reglas que se le indicaron para evitar ejecutar *exploits* y dañar sistemas. Asimismo, se identificó que el agente sí determina cuáles son los caminos más prometedores y los prioriza.

### c) Caso C: ESCALACIÓN DE PRIVILEGIOS

Posterior a la explotación, el analista utilizó el siguiente *prompt*: “Revisa si puedes obtener los máximos privilegios dentro del sistema”, ante lo que el agente realizó los siguientes pasos:

- Validación de los usuarios existentes dentro del sistema.
- Revisión de los servicios identificados en la etapa anterior.
- Reutilización de credenciales.
- Revisión de paquetes con configuraciones incorrectas.
- Reconocimiento de paquetes con permisos que permiten escalar.
- Escalación de privilegios.

Finalmente, durante esta fase, se identificó que el usuario existente contaba con permisos para ejecutar la herramienta administrativa *bee* con privilegios elevados. El abuso de su función *eval*, que permite ejecutar código PHP malicioso tras inicializar el entorno de *Backdrop*, derivó en ejecución de comandos como *root* y el compromiso total del sistema.

En la Tabla 2, se muestra a manera de resumen, las actividades que el agente llevó a cabo y los resultados obtenidos en cada una de ellas. Cabe destacar que el agente de IA obtuvo algunos falsos positivos que requirieron la validación del analista. Asimismo, se observó que el comportamiento del agente se regula de acuerdo con la información que se genera en tiempo real.

**TABLA 2. GEMINI**

Fases	Herramienta	Técnicas usadas	Evidencias obtenidas	Resultado	Tiempo (min)
□ <i>Reconnaissance</i>	<i>Nmap</i>	Escaneo de puertos y servicios	Se descubrieron puertos, tecnologías y versiones: 22 (SSH) y 80 (HTTP), Apache 2.4.41, Backdrop CMS 1.27.1.	Identificación de puertos abiertos y versiones de servicios.	3
□ <i>Scanning</i>	<i>Nikto</i>	Navegación web detección de tecnologías usadas en ese puerto.	Se obtuvo el archivo “robots.txt”, el cual revela rutas en el gestor de contenido <i>Backdrop</i> .	Confirmación de tecnologías <i>Linux + Apache + PHP CMS</i> .	2
□ <i>Enumeration</i>	<i>Feroxbuster</i>	Listar directorios y obtener archivos encontrados en la web <i>/.git</i> .	Se obtuvo el archivo “settings.php” con credenciales de una base de datos detectada al momento de la intrusión, usuario “tiffany”.	Credenciales reutilizables encontradas.	6
□ <i>Gaining Access</i>	<i>Python</i>	Inicio de sesión en gestor de contenido + Subir archivos.	Ejecución de comandos con el archivo “shell.php”.	<i>Shell</i> como <i>www-data</i> .	5
□ <i>Lateral Movement</i>	<i>SSH</i>	Probar la reutilización de contraseñas contra el servicio de SSH.	Se obtuvo acceso mediante SSH.	<i>Shell</i> como <i>johncusack</i> .	5
□ <i>Privilege Escalation</i>	<i>Bash</i>	Revisar permisos del comando <i>sudo -l</i> , abuso de del usuario <i>bee eval</i> .	Se reconoció que el usuario <i>johncusack</i> puede ejecutar <i>/usr/local/bin/bee</i> como <i>root</i> .	RCE como <i>root</i> .	7

<input type="checkbox"/> <i>Maintaining Access</i>	<i>Netcat</i>	Obtener <i>shell</i> reversa estable + <i>TTY upgrade</i> .	Se obtuvo una terminal <i>script /dev/null -c bash</i> .	<i>Shell</i> interactiva estable.	6
<input type="checkbox"/> <i>Clearing Tracks</i>	N/A	Remover archivos, limpiar <i>logs</i> .	No es necesario.	—	2

### 3) Escenario 3: LLM CHATGPT

En este escenario fue realizado por la analista con ayuda de un agente de IA en el que se implementó como LLM a ChatGPT.

#### a) Caso A: RECONOCIMIENTO AUTOMATIZADO

El *prompt* utilizado por la analista en este escenario fue el mismo que se utilizó en el escenario 2: “Realiza una enumeración de puertos y servicios, recopila la información más relevante y enumera configuraciones débiles. Si se encuentra la manera de hacer una ejecución remota de comandos y posteriormente escalar privilegios no lo hagas, es solo enumeración”.

Del que también se obtuvo como resultado la identificación de un servidor web y un repositorio con acceso público, la diferencia del escenario consistió en el tiempo que llevó la ejecución, mismo que puede ser consultado en la Tabla 3.

#### b) Caso B: EXPLOTACIÓN

Para este caso, se hizo uso del mismo *prompt* del escenario 2: “Si se encuentra la manera de hacer una ejecución remota de comandos y, posteriormente, escalar privilegios, hazlo, solo intenta no tirar los sistemas, no realices ataques agresivos”.

Los resultados también fueron los mismos del escenario 2 con una variación en el tiempo en el que fueron finalizadas las actividades necesarias, tal como se observa en la Tabla 3.

#### c) Caso C: ESCALACIÓN DE PRIVILEGIOS

En esta fase, se hizo uso del mismo *prompt* utilizado en el escenario 2: “Revisa si puedes obtener los máximos privilegios dentro del sistema”. El agente realizó pasos similares a los del escenario anterior, nuevamente se observaron diferencias en el tiempo de respuesta.

En la Tabla 3, se muestran los detalles sobre las actividades realizadas por el agente y el tiempo que le tomó proporcionarle los detalles a la analista.

TABLA 3. CHATGPT

Fases	Herramienta	Técnicas usadas	Evidencias obtenidas	Resultado	Tiempo (min)
□ <i>Reconnaissance</i>	<i>Nmap</i>	Escaneo de puertos y servicios	Se descubrieron puertos, tecnologías y versiones: 22 (SSH) y 80 (HTTP), Apache 2.4.41, Backdrop CMS 1.27.1.	Identificación de puertos abiertos y versiones de servicios.	4
□ <i>Scanning</i>	<i>Nikto</i>	Navegación web detección de tecnologías usadas en ese puerto.	Se obtuvo el archivo “robots.txt”, el cual revela rutas en el gestor de contenido <i>Backdrop</i> .	Confirmación de tecnologías <i>Linux</i> + <i>Apache</i> + <i>PHP</i> <i>CMS</i> .	3
□ <i>Enumeration</i>	<i>Feroxbuster</i>	Listar directorios y obtener archivos encontrados en la <i>web /.git</i> .	Se obtuvo el archivo “settings.php” con credenciales de una base de datos detectada al momento de la intrusión, usuario “tiffany”.	Credenciales reutilizables encontradas.	10
□ <i>Gaining Access</i>	<i>Python</i>	Inicio de sesión en en gestor de contenido + Subir archivos.	Ejecución de comandos con el archivo “shell.php”.	<i>Shell</i> como <i>www-data</i> .	6
□ <i>Lateral Movement</i>	<i>SSH</i>	Probar la reutilización de contraseñas contra el servicio de SSH.	Se obtuvo acceso mediante SSH.	<i>Shell</i> como <i>johnCUSACK</i> .	8
□ <i>Privilege Escalation</i>	<i>Bash</i>	Revisar permisos del comando <i>sudo -l</i> , abuso de del usuario <i>bee eval</i> .	Se reconoció que el usuario <i>johnCUSACK</i> puede ejecutar <i>/usr/local/bin/bee</i> como <i>root</i> .	RCE como <i>root</i> .	9
□ <i>Maintaining Access</i>	<i>Netcat</i>	Obtener <i>shell</i> reversa estable + <i>TTY upgrade</i> .	Se obtuvo una terminal <i>script /dev/null -c bash</i> .	<i>Shell</i> interactiva estable.	7

<input type="checkbox"/> <i>Clearing Tracks</i>	N/A	Remover archivos, limpiar logs.	No es necesario.	—	3
---	-----	---------------------------------	------------------	---	---

#### 4) Escenario 4: LLM CURSOR

En este escenario se llevó a cabo por la analista mediante el uso de un agente de IA con Cursor implementado como LLM.

Dado que los prompts utilizados fueron los mismos que los de los escenarios anteriores y que en las actividades realizadas por el agente, así como en los resultados finales obtenidos no se observaron diferencias, en este escenario únicamente se muestra la Tabla 4 en la que se encuentran los detalles y tiempos de las fases desarrolladas.

TABLA 4. CURSOR

Fases	Herramienta	Técnicas usadas	Evidencias obtenidas	Resultado	Tiempo (min)
<input type="checkbox"/> <i>Reconnaissance</i>	<i>Nmap</i>	Escaneo de puertos y servicios	Se descubrieron puertos, tecnologías y versiones: 22 (SSH) y 80 (HTTP), Apache 2.4.41, Backdrop CMS 1.27.1.	Identificación de puertos abiertos y versiones de servicios.	3
<input type="checkbox"/> <i>Scanning</i>	<i>Nikto</i>	Navegación web detección de tecnologías usadas en ese puerto.	Se obtuvo el archivo “robots.txt”, el cual revela rutas en el gestor de contenido <i>Backdrop</i> .	Confirmación de tecnologías <i>Linux + Apache + PHP CMS</i> .	3
<input type="checkbox"/> <i>Enumeration</i>	<i>Feroxbuster</i>	Listar directorios y obtener archivos encontrados en la web <i>/.git</i> .	Se obtuvo el archivo “settings.php” con credenciales de una base de datos detectada al momento de la intrusión, usuario “tiffany”.	Credenciales reutilizables encontradas.	6
<input type="checkbox"/> <i>Gaining Access</i>	<i>Python</i>	Inicio de sesión en en gestor de contenido + Subir archivos.	Ejecución de comandos con el archivo “shell.php”.	<i>Shell</i> como <i>www-data</i> .	4

<input type="checkbox"/> <i>Lateral Movement</i>	<i>SSH</i>	Probar la reutilización de contraseñas contra el servicio de SSH.	Se obtuvo acceso mediante SSH.	<i>Shell</i> como <i>johncusack</i> .	6
<input type="checkbox"/> <i>Privilege Escalation</i>	<i>Bash</i>	Revisar permisos del comando <i>sudo -l</i> , abuso de del usuario <i>bee eval</i> .	Se reconoció que el usuario <i>johncusack</i> puede ejecutar <i>/usr/local/bin/bee</i> como <i>root</i> .	RCE como <i>root</i> .	6
<input type="checkbox"/> <i>Maintaining Access</i>	<i>Netcat</i>	Obtener <i>shell</i> reversa estable + <i>TTY upgrade</i> .	Se obtuvo una terminal <i>script /dev/null -c bash</i> .	<i>Shell</i> interactiva estable.	4
<input type="checkbox"/> <i>Clearing Tracks</i>	N/A	Remover archivos, limpiar <i>logs</i> .	No es necesario.	—	2

## 22

### E. Retroalimentación y aprendizaje

Una vez que las herramientas terminan su ejecución, el agente se encarga de revisar los resultados generados y evalúa si se requiere una nueva iteración, en caso de que así sea, ejecuta herramientas adicionales; de lo contrario, el proceso se detiene. Esta retroalimentación permite:

- Refinar el proceso de planificación.
- Incrementar la efectividad.
- Ajustar las reglas y los criterios de análisis.

Cabe mencionar que el servidor MCP es capaz de generar registros (*logs*), lo que permitió el monitoreo de las acciones realizadas y de las herramientas ejecutadas.

## VI. RESULTADOS Y DISCUSIÓN

En esta sección se analizan los resultados integrando una comparación cuantitativa y una evaluación cualitativa del desempeño y comportamiento del agente de inteligencia artificial. El objetivo es comprender su capacidad de autonomía, razonamiento y los mecanismos de control cuando se ejecutan herramientas que hacen uso de agentes de IA.

Se observó que los resultados obtenidos de los dos casos descritos anteriormente muestran que la explotación con ayuda de agentes de IA, bajo controles bien definidos, puede ayudar a los *pentesters* a tomar decisiones.

### A. Análisis cuantitativo de resultados

Se revisaron métricas utilizadas en pruebas de penetración, tales como:

- Tiempo de ejecución.
- Cantidad de hallazgos y calidad de estos.
- Falsos positivos detectados.
- Evaluación de la superficie de ataque que se siguió.
- Intervención humana requerida.

En la Tabla 5, se muestra una comparativa del tiempo que le tomó a cada uno de los agentes llevar a cabo las actividades correspondientes, en esta es posible observar que el agente que tenía implementado Cursor como LLM fue el más rápido.

**TABLA 5. COMPARATIVA DE AGENTES**

Fase	Tiempo (min)	Tiempo (min)	Tiempo (min)	¿A quién le toma menos tiempo?
	Gemini	ChatGPT	Cursor	
<input type="checkbox"/> Reconnaissance	3	4	3	Gemini y Cursor
<input type="checkbox"/> Scanning	2	3	3	Gemini
<input type="checkbox"/> Enumeration	6	10	6	Gemini y Cursor
<input type="checkbox"/> Gaining Access	5	6	4	Cursor
<input type="checkbox"/> Lateral Movement	5	8	6	Gemini
<input type="checkbox"/> Privilege Escalation	7	9	6	Cursor
<input type="checkbox"/> Maintaining Access	6	7	4	Cursor
<input type="checkbox"/> Clearing Tracks	2	3	2	Gemini y Cursor

En la Tabla 5, también es posible observar que, con una pequeña diferencia, el agente en el que se implementó Cursor como LLM tiene una pequeña ventaja en tiempo a los anteriores. Por lo que, en la Tabla 6, es posible observar los tiempos de este agente (escenario 3) en comparación con el escenario realizado manualmente por el *pentester*.

**TABLA 6. COMPARATIVA AGENTE VS. PENTESTER**

Fase	Tiempo (min)	Tiempo (min)	¿A quién le toma menos tiempo?	Porcentaje de ventaja
	Gemini	Pentester		
<input type="checkbox"/> <i>Reconnaissance</i>	3	8	Agente	62.5%
<input type="checkbox"/> <i>Scanning</i>	2	5	Agente	40%
<input type="checkbox"/> <i>Enumeration</i>	6	15	Agente	60%
<input type="checkbox"/> <i>Gaining Access</i>	5	7	Agente	28.5%
<input type="checkbox"/> <i>Lateral Movement</i>	5	6	Agente	16.6%
<input type="checkbox"/> <i>Privilege Escalation</i>	7	12	Agente	41.6%
<input type="checkbox"/> <i>Maintaining Access</i>	6	5	Pentester	16.6%
<input type="checkbox"/> <i>Clearing Tracks</i>	2	2	-	-

24

Dado que se observó que el agente con respuestas más rápidas también supera al tiempo del *pentester*, se realizó una comparativa del agente al que le llevó más tiempo obtener los resultados con el *pentester*, misma que se visualiza en la Tabla 7.

**TABLA 7. COMPARATIVA AGENTE VS. PENTESTER**

Fase	Tiempo (min)	Tiempo (min)	¿A quién le toma menos tiempo?	Porcentaje de ventaja
	ChatGPT	Pentester		
<input type="checkbox"/> <i>Reconnaissance</i>	4	8	Agente	50%
<input type="checkbox"/> <i>Scanning</i>	3	5	Agente	40%
<input type="checkbox"/> <i>Enumeration</i>	10	15	Agente	60%

<input type="checkbox"/> <i>Gaining Access</i>	6	7	Agente	33.3%
<input type="checkbox"/> <i>Lateral Movement</i>	8	6	<i>Pentester</i>	25%
<input type="checkbox"/> <i>Privilege Escalation</i>	9	12	Agente	25%
<input type="checkbox"/> <i>Maintaining Access</i>	7	5	<i>Pentester</i>	28.5%
<input type="checkbox"/> <i>Clearing Tracks</i>	3	2	<i>Pentester</i>	33.3%

Aunque en esta comparación el *pentester* sí tiene ventaja en algunas fases, la mayoría de ellas siguen siendo realizadas en un menor tiempo por el agente más tarde entre los evaluados.

### B. Análisis cualitativo del comportamiento del agente

El análisis cualitativo sugiere que esta solución puede complementar el trabajo de los profesionales de pruebas de penetración mediante la automatización de tareas específicas de análisis y reconocimiento. Además, la integración de información proveniente de múltiples herramientas permite adaptar la estrategia de ejecución conforme se generan nuevos datos durante la evaluación.

Si bien el agente resultó ser una herramienta de gran utilidad, también generó falsos positivos, y es en este punto donde se destaca la importancia de la validación y el juicio humano.

Adicionalmente, se observó un comportamiento alineado con las reglas que fueron indicadas al agente, evitando acciones destructivas o intentos de explotación no justificados.

Por otra parte, en la figura 10 se muestra una comparativa del total de tiempo total en minutos que se requirió en cada escenario para finalizar las tareas del ejercicio de *pentesting*.

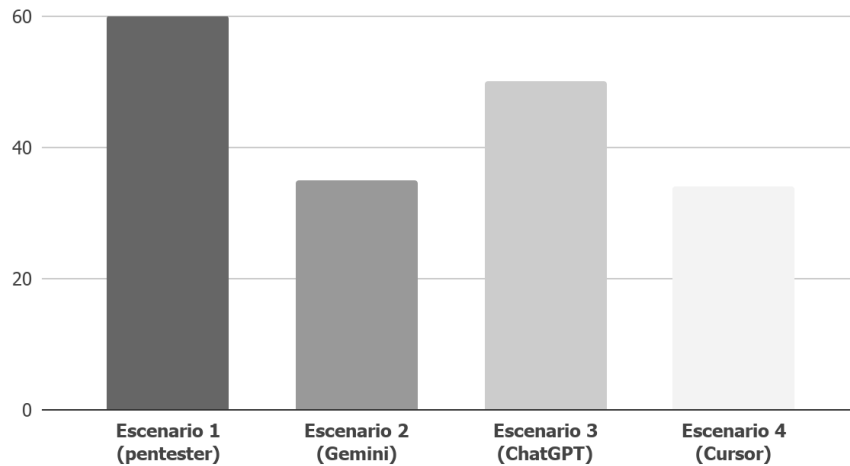


Fig. 10. Gráfica de comparación de eficiencia.

### C. Comparación con enfoques tradicionales

26

Si comparamos el proceso que realizó el agente con los métodos y herramientas tradicionales utilizados en las pruebas de penetración y el análisis de vulnerabilidades, existen diferencias relevantes, descritas a continuación:

- Aunque los escáneres semiautomáticos clásicos, como *Nessus*, permiten escanear sistemas rápidamente, carecen de capacidades de razonamiento, lo que imposibilita modificar su comportamiento ante variables que pueden presentarse durante el proceso.
- Las validaciones manuales son necesarias, pero consumen mucho tiempo y esfuerzo, además, están limitadas por la experiencia y el juicio del *pentester*.

Este enfoque demuestra que el camino óptimo es que las herramientas basadas en agentes inteligentes tengan función de apoyo, más que como un sustituto del personal humano.

Finalmente, se concluye que el uso de estas herramientas es, y seguirá siendo, un apoyo invaluable que continuará evolucionando. Actualmente, no existen datos claros que indiquen que puedan reemplazar por completo la experiencia y el juicio de un profesional de la ciberseguridad.

### B. Beneficios y métricas de rendimiento

Con base en los estudios analizados, se determinaron y organizaron las ventajas competitivas de cada propuesta técnica. Los beneficios clave que se han encontrado en la muestra documental se resaltan en la figura 2. El estudio del corpus mostró que el 73.33% ( $n=22$  de 30 estudios) de las investigaciones consideraron la Alta Precisión y la Tasa de Detección Exacta como su mayor beneficio, logrando frecuentemente niveles por encima del 95%, e incluso llegando al 100% en situaciones controladas de laboratorio. En segundo lugar, el 13.33% ( $n=4$  de 30 estudios) destacó la independencia de infraestructura, lo que posibilitó la detección sin hardware adicional en la red (enfoque del lado del cliente). El resto de la muestra priorizó el bajo costo y latencia (10.00%,  $n=3$ ) y la escalabilidad (3.33%,  $n=1$ ). Estos hallazgos evidencian que, aunque el objetivo principal es la precisión algorítmica, se está mostrando un interés cada vez mayor en soluciones autónomas y descentralizadas.

## VII. LIMITACIONES Y CONSIDERACIONES ÉTICAS

Aun con los beneficios que puede aportar el uso de agentes inteligentes en la ciberseguridad, específicamente en las pruebas de penetración, los resultados indican una serie de riesgos operativos y limitaciones.

Un área de oportunidad importante son los falsos positivos. Este problema puede presentarse en entornos que han sido personalizados, por ejemplo, cuando existen servicios desplegados en puertos distintos a los estándares. En estos casos, el sistema puede intentar escanear, enumerar e incluso explotar un puerto, asumiendo que corresponde a un servicio determinado, cuando en realidad se trata de otro, o bien de sistemas con configuraciones atípicas. Esta situación hace necesaria la validación manual de los resultados generados.

27

Otro riesgo crítico es el potencial uso indebido de estas herramientas por parte de personas con conocimientos mínimos, las cuales podrían intentar atacar sistemas sin con autorización mediante la ejecución de un solo *prompt*. Actualmente, existen evidencias de que los ciberdelincuentes ya están haciendo uso de modelos de IA para llevar a cabo actividad maliciosa, se ha observado que suelen utilizarla para la falsificación de documentos, creación de correos electrónicos de *phishing*, escaneo y explotación de vulnerabilidades, desarrollo de *malware*, entre otros. En otras palabras, estas herramientas amplifican las capacidades de estos actores maliciosos.

Con base en los resultados obtenidos en esta investigación, se sugiere adoptar un enfoque de “autonomía supervisada” para la gestión de la autonomía de los agentes, así como implementar controles estrictos sobre quién puede utilizar estas herramientas dentro de las redes de una empresa.

Adicionalmente, la adopción de agentes basados en inteligencia artificial para actividades de ciberseguridad plantea desafíos relacionados con la responsabilidad de las acciones ejecutadas, la trazabilidad de las decisiones tomadas por los modelos y la rendición de cuentas ante posibles errores. En consecuencia, se considera indispensable mantener mecanismos de supervisión humana, auditoría y control que permitan garantizar un uso ético y seguro de estas tecnologías.

Desde una perspectiva ética, el uso de agentes inteligentes en pruebas de penetración debe apegarse a principios de legalidad, autorización explícita y minimización de impacto. Las capacidades de automatización analizadas en este trabajo fueron evaluadas en entornos controlados y con fines de investigación, evitando su aplicación sobre sistemas sin consentimiento. Asimismo, los resultados presentados tienen como propósito fortalecer las capacidades defensivas de las organizaciones y contribuir a la comprensión de las amenazas, no facilitar actividades ofensivas no autorizadas.

## VIII. CONCLUSIONES Y TRABAJO FUTURO

Este trabajo estudió el funcionamiento e implementación de LLM combinados con servidores MCP con el objetivo de crear un agente que brindara apoyo en pruebas de penetración a sistemas. La propuesta integra capacidades asociadas a la ciberseguridad ofensiva y defensiva, permitiendo asistir actividades relacionadas con los equipos Red Team y Blue Team, tal como se muestra en la Figura 7. Asimismo, la infraestructura incluyó herramientas ofensivas ampliamente utilizadas en la industria, las cuales se comunican mediante una arquitectura basada en MCP. Esta arquitectura, representada en las Figuras 8 y 9, proporciona una capa estandarizada de comunicación entre el modelo de lenguaje, el agente y las herramientas externas, facilitando la ejecución coordinada de tareas especializadas.

Los resultados obtenidos muestran que la integración de agentes basados en LLM con herramientas de pentesting permite automatizar parcialmente actividades de reconocimiento, análisis y ejecución de pruebas de seguridad, reduciendo la intervención manual requerida en determinadas etapas del proceso. De igual forma, se observó que la capacidad de integrar información proveniente de múltiples fuentes y herramientas favorece la generación de estrategias de acción más consistentes y adaptables al contexto evaluado.

No obstante las pruebas también evidenciaron limitaciones relacionadas con la validación de resultados, la posibilidad de generar falsos positivos y la necesidad de supervisión humana durante la ejecución de actividades críticas. Por ello, estas tecnologías aún no pueden considerarse un sustituto directo de los *pentesters*; más bien, constituyen una herramienta complementaria capaz de incrementar las capacidades operativas del profesional y optimizar ciertas tareas especificadas.

28

En conjunto, los hallazgos sugieren que los agentes inteligentes integrados mediante servidores MCP representan una alternativa prometedora para incrementar el nivel de automatización alcanzable en ciberseguridad. Su evolución podría favorecer el desarrollo de soluciones capaces de coordinar herramientas especializadas, adaptarse a distintos escenarios de evaluación y asistir de manera más efectiva a los equipos de seguridad en entornos reales.

Algunas líneas de investigación que se abrieron son las siguientes:

- Desarrollar lineamientos y políticas para agentes de inteligencia artificial ofensivos.
- Entrenamiento especializado para ramas específicas de la ciberseguridad, como, por ejemplo, seguridad en la nube, aplicaciones móviles o Active Directory.
- Reducción de falsos positivos y sesgos mediante validación híbrida.

Estas líneas de investigación podrían permitir el desarrollo de agentes experimentales que produzcan mejores resultados, alineados con los estándares utilizados en la industria de la ciberseguridad.

### **CRediT** (*Contributor Roles Taxonomy*)

**Contribuciones de los autores:** Conceptualización: GCE; Metodología: GCE; Software: GCE; Investigación: GCE, MVT Redacción y preparación del borrador original: GCE, MVT; Redacción, revisión y edición: GCE, MVT; Supervisión: GCE; Análisis formal: GCE; Administración del proyecto: GCE; Adquisición de fondos: GCE.

**Financiamiento:** Los autores declaran que el presente estudio no recibió financiamiento externo. El desarrollo experimental se realizó en un entorno controlado de laboratorio virtual.

**Declaración de disponibilidad de datos:** Los datos generados durante los experimentos corresponden a entornos simulados en laboratorio virtual. No se utilizaron datos personales ni información de sistemas reales sin autorización.

**Conflicto de interés:** El autor declara no tener conflicto de interés relacionado con el contenido de este estudio.

## REFERENCIAS

- [1] Google, “El optimismo global sobre la IA aumenta a medida que crece su uso,” Google Blog, 2023. [Online]. Available: <https://blog.google/intl/es-419/noticias-de-la-empresa/el-optimismo-global-sobre-la-ia-aumenta-a-medida-que-crece-su-uso/>
- [2] National Institute of Standards and Technology, “National Vulnerability Database (NVD),” NIST, 2024. [Online]. Available: <https://nvd.nist.gov/vuln/search/#nvd/home?resultType=statistics>
- [3] National Institute of Standards and Technology, NIST Cybersecurity White Paper 29, NIST CSWP 29, p. 18, 2023. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/CSWP/NIST.CSWP.29.pdf>
- [4] EC-Council, “Penetration Testing Phases,” EC-Council Cybersecurity Exchange, 2024. [Online]. Available: <https://www.eccouncil.org/cybersecurity-exchange/penetration-testing/penetration-testing-phases/>
- [5] P. Cichonski, T. Millar, T. Grance, K. Scarfone, *Computer Security Incident Handling Guide*, NIST SP 800-61 Rev. 2, National Institute of Standards and Technology, 2012. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-61r2.pdf>
- [6] National Institute of Standards and Technology, Secure Software Development Framework (SSDF) Version 1.1: Recommendations for Mitigating the Risk of Software Vulnerabilities, NIST SP 800-218A, 2023. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.800-218A.pdf>
- [7] Google, Artificial Intelligence at Google: Our Principles, 2018. [Online]. Available: <https://ai.google/principles/>
- [8] A. Vaswani et al., “Attention Is All You Need,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017. [Online]. Available: <https://arxiv.org/abs/1706.03762>
- [9] Microsoft, Microsoft Responsible AI Standard, 2022. [Online]. Available: <https://www.microsoft.com/en-us/ai/responsible-ai>
- [10] Hugging Face, “AI Agents Course – What Are Agents?” 2025. [Online]. Available: <https://huggingface.co/learn/agents-course>
- [11] MITRE Corporation, “Key Concepts – MITRE ATT & CK,” 2024. [Online]. Available: <https://attack.mitre.org/resources/>
- [12] MITRE Corporation, “MITRE ATT & CK®,” 2024. [Online]. Available: <https://attack.mitre.org/>
- [13] The OWASP Foundation, OWASP Top 10:2021, 2021. [Online]. Available: <https://owasp.org/Top10/2021/>
- [14] National Institute of Standards and Technology, “About NIST,” U.S. Department of Commerce, 2024. [Online]. Available: <https://www.nist.gov/about-nist>
- [15] National Institute of Standards and Technology, Incident Response Recommendations and Considerations for Cybersecurity Risk Management, NIST SP 800-61 Rev. 3, 2025. [Online]. Available: <https://doi.org/10.6028/NIST.SP.800-61r3>
- [16] Microsoft, “What Is Endpoint Detection and Response (EDR),” Microsoft Security, 2024. [Online]. Available: <https://www.microsoft.com/en-us/security/business/security-101/what-is-edr-endpoint-detection-response>
- [17] National Institute of Standards and Technology, Technical Guide to Information Security Testing and Assessment, NIST SP 800-115, 2008. [Online]. Available: <https://nvlpubs.nist.gov/nistpubs/Legacy/SP/nistspecialpublication800-115.pdf>
- [18] Qualys, Inc., “Vulnerability Management, Detection and Response (VMDR),” 2024. [Online]. Available: <https://www.qualys.com/apps/vulnerability-management-detection-response/>
- [19] E. Nijkamp, J. Hayase, C. Xiong, et al., “Sec-PaLM: Aligning Large Language Models with Security Expertise,” arXiv, arXiv:2309.06106, Sep. 2023. [Online]. Available: <https://arxiv.org/abs/2309.06106>
- [20] Microsoft, Microsoft Security Copilot Documentation – Overview, Jul. 15, 2025. [Online]. Available: <https://learn.microsoft.com/en-us/copilot/security/microsoft-security-copilot>
- [21] B. Biggio, F. Roli, “Wild Patterns: Ten Years After the Rise of Adversarial Machine Learning,” *Pattern Recognition*, vol. 84, pp. 317–331, 2018, doi: <https://doi.org/10.1016/j.patcog.2018.07.023>
- [22] A. Pautov, “AI-Driven Pentesting at Home Using HexStrike AI for Full Network Discovery and Exploitation,” Medium, 2024. [Online]. Available: <https://medium.com/@1200km/ai-driven-pentesting-at-home-using-hexstrike-ai-for-full-network-discovery-and-exploitation-00a9e88b3bde>
- [23] Model Context Protocol, “Model Context Protocol Specification,” 2024. [Online]. Available: <https://github.com/modelcontextprotocol/specification>