MCP Security for Enterprise Organizations

Claves para mejorar la postura seguridad en nuestros MCPs



EMPECEMOS →



Un poco sobre mí

- Trabajo en Cyber hace ya unos 7 años.
- Estoy cursando el Master de Seguridad Informática en la UBA.
- Co-fundador de ThreatX Security.
- Actualmente estoy como Sr Product Security y como Instructor de un Bootcamp.
- Tengo algunas certs de la industria.
- Soy el creador de PhiloCyber, mi proyecto personal.
- Amo los animales, sobre todos los perros.
- Me gustan los deportes, especialmente natación.
- Y disfruto mucho lo que hacemos.



Temas de hoy

01. Qué es MCP?

02. Arquitectura

03. Cómo Funciona?

04. Nivel de Adopción

05. Vulnerabilidades

06. Riesgos

07. Mitigaciones

08. Cierre y Conclusión

Milenios atrás...







FUNDAMENTOS EXPERIMENTALES

- Function calling en LLMs (OpenAI)
- JSON-RPC 2.0 como estándar
- Arquitectura cliente-servidor



CREACIÓN DE MCP

 En noviembre se comparte el protocolo, como un marco de trabajo open-source.



SE FUNDA ANTHROPIC

- Hermanitos Amodei
- Ex VP de Investigación, OpenAl
- Ex VP de Políticas y Seguridad, OpenAl.



MADURACIÓN Y AGENTES AUTÓNOMOS

- Al Agents autónomos.
- Integración con sistemas externos.
- Protocolos de comunicación LLM.



¿Qué es MCP?

Model Context Protocol

Es un protocol que define y estandariza como los LLMs se conectan con herramientas externas.

Qué busca?

Posicionado como el <u>USB-C para aplicaciones de IA</u>, establece una arquitectura cliente-servidor que permite que los modelos de IA interactúen con diversas fuentes de datos y herramientas a través de una interfaz unificada.



¿Qué es un LLM?

"Es un tipo de inteligencia artificial entrenado con una gran cantidad de texto para comprender y generar lenguaje humano".



GPT-4



Qwen



Gro



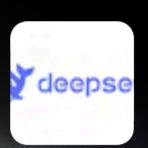
Mistral Al



Gemini



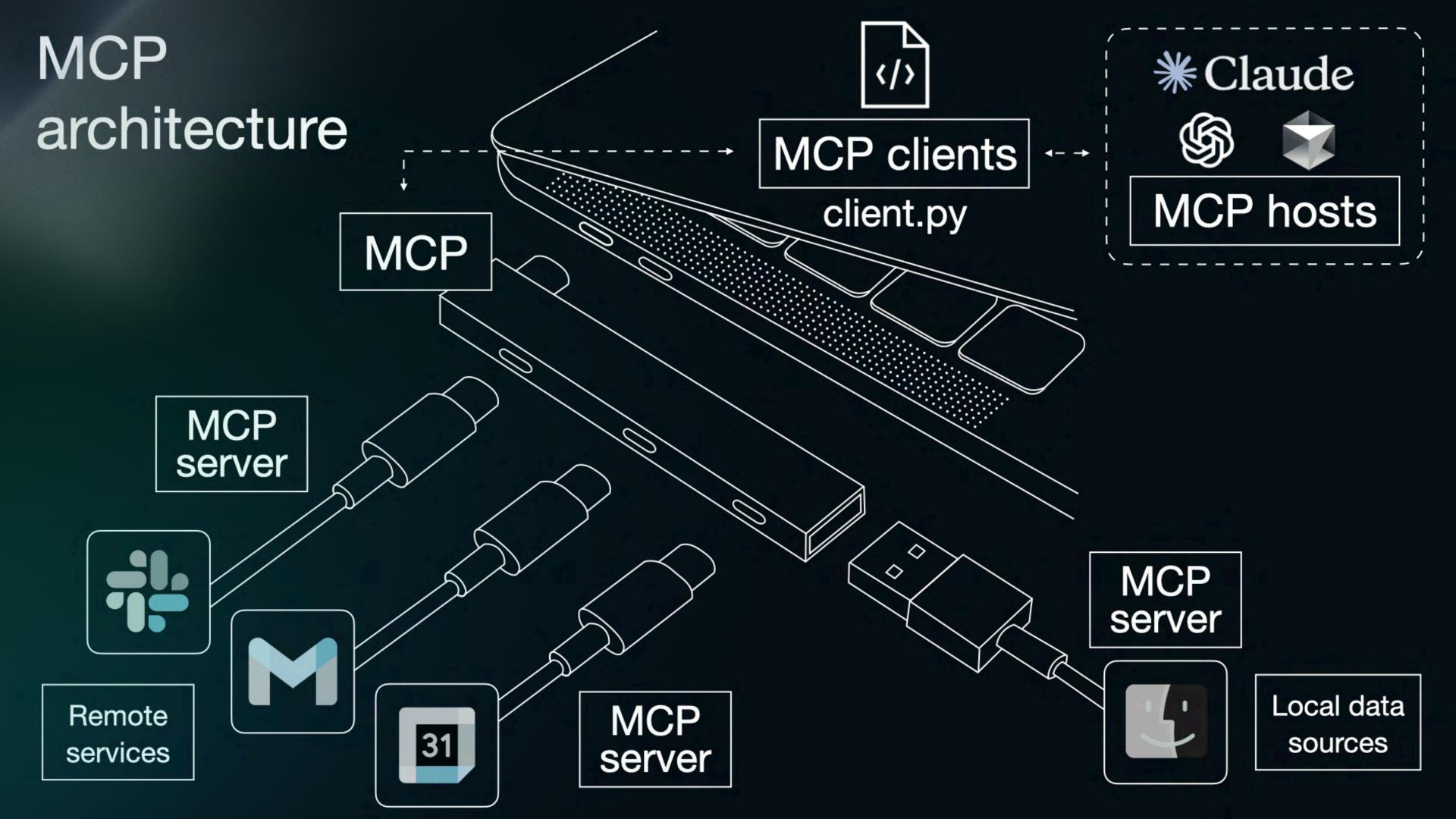
Falcon



DeepSeek-R1



Claude



Mecanismos de Transporte MCP

stdio (Local)

- Pipes STDIN/STDOUT
- Latencia µsegundos
- Solo local, seguro
- Uso: CLI, apps escritorio

Streamable HTTP (Remoto)

☑ Oficial

- Endpoint HTTP único
- 290-300 req/seg
- OAuth 2.1 + TLS
- Uso: Apps web producción

WebSocket (Remoto)

77 Comunidad

- Full-duplex persistente
- Bidireccional tiempo real
- Implementación no oficial
- Uso: Apps tiempo real

Formato de Mensajes:

JSON-RPC 2.0 (todos los transportes)

Request (con ID)

Notification (sin ID)

Response (ID coincidente)

MCP Client

Invokes **Tools**Queries for **Resources**Interpolates **Prompts**

MCP Server

Exposes **Tools**

Exposes Resources

Exposes Prompts

Tools

Model-controlled

Functions invoked by the model

Retrieve / search

Send a message

Update DB records

Resources

Application-controlled

Data exposed to the application

Files

Database Records

API Responses

Prompts

User-controlled

Pre-defined templates for AI interactions

Document Q&A

Transcript Summary

Output as JSON

¿Cómo funciona?

Flujo MCP para Agendar Charla Ekoparty Usuario Host (Claude/LLM) Cliente MCP Servidor MCP (Ekoparty) (0) "Agendar 'El poder de la automatizacion en la era de AI' mié. 11:35 - 12:20" (1) "Conectar y listar herramientas" (2) Handshake (initialize/result/initialized) (3) tools/list() (4) tools[] (schema 'agendarCharla') (5) "Herramientas disponibles: [schema]" (6) LLM (en Host) procesa (0) y (5). Decide llamar a 'agendarCharla'. (7) "Ejecuta: tools/call('agendarCharla', ...)" (8) tools/call('agendarCharla', args: ...) (9) Servidor ejecuta lógica interna (API/DB) y confirma. (10) Result(content: "Éxito: Charla agendada") (11) "Resultado de la herramienta: 'Éxito...'" (12) Respuesta final (LLM): "¡Listo! Tu charla fue agendada." Host (Claude/LLM) Usuario Cliente MCP Servidor MCP (Ekoparty)



Nivel de Adopción Actual

Adopción de IA

90% de los desarrolladores profesionales uitilizan IA en su día a día

Productividad

80% reporto que la IA incrementó su productividad

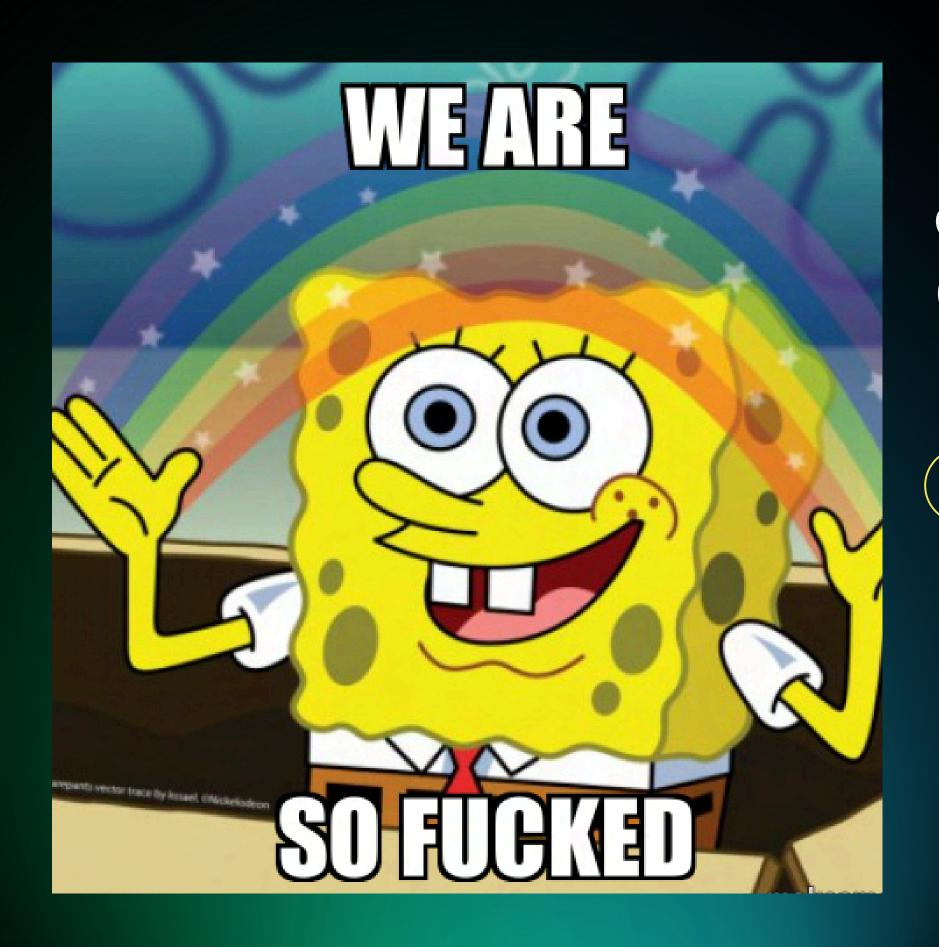
Dependencia

El 65% depende "moderadamente" de la IA en el trabajo

Calidad

59% afirmaron que la IA impactó positivamente en la calidad de código

Fuente: State of Al-assisted Software Development - DORA Research 2025



¿Cuál es el problema entonces?

API Calls





LLMs tooling



Introduce una nueva capa semántica, donde los Agentes de IA interpretan lenguaje natural para usar herramientas y recursos externos.



Amenazas e Implicaciones

Tool Poisoning Attacks

Una inyección de prompts indirecta, donde el prompt malicioso está escondido en la descripción de la herramienta.

Tool Spoofing

Un atacante crea un servidor falso que parece legítimo, y los agentes Al se confunden, enviando datos confidenciales al servidor incorrecto.

The Rug Pull

Un servidor MCP opera bien y parece confiable, pero más tarde pasa a comportarse de forma maliciosa, atacando cuando menos lo esperan.

Prompt Injection

Alguien esconde instrucciones maliciosas dentro de la descripción o entrada para la herramienta... la IA las sigue y ejecuta acciones arbitrarias.

Unauthenticated access

No hay controles de identidad adecuados, seria como dejar la puerta abierta, cualquiera puede entrar y operar como si fuera legítimo.

Excessive Permissions

El servidor MCP tiene demasiados permisos o acceso a recursos. Si lo hackean, el daño puede ser enorme porque pueden acceder a todo.



¿Cómo mitigamos estos riesgos?



No hay secretos, ir siempre de lo más básico a lo más específico. Desde diseño y desarrollo seguro, hasta seguridad en profundidad, logs y monitoreo continuo y testeos.

ZERO TRUST

Nunca confiar en nada ni nadie, privilegios mínimos para el rol del usuario o la funcionalidad de la herramienta.

DEPENDENCIAS Y PARCHEOS

Revisar y monitorear dependencias constantemente,

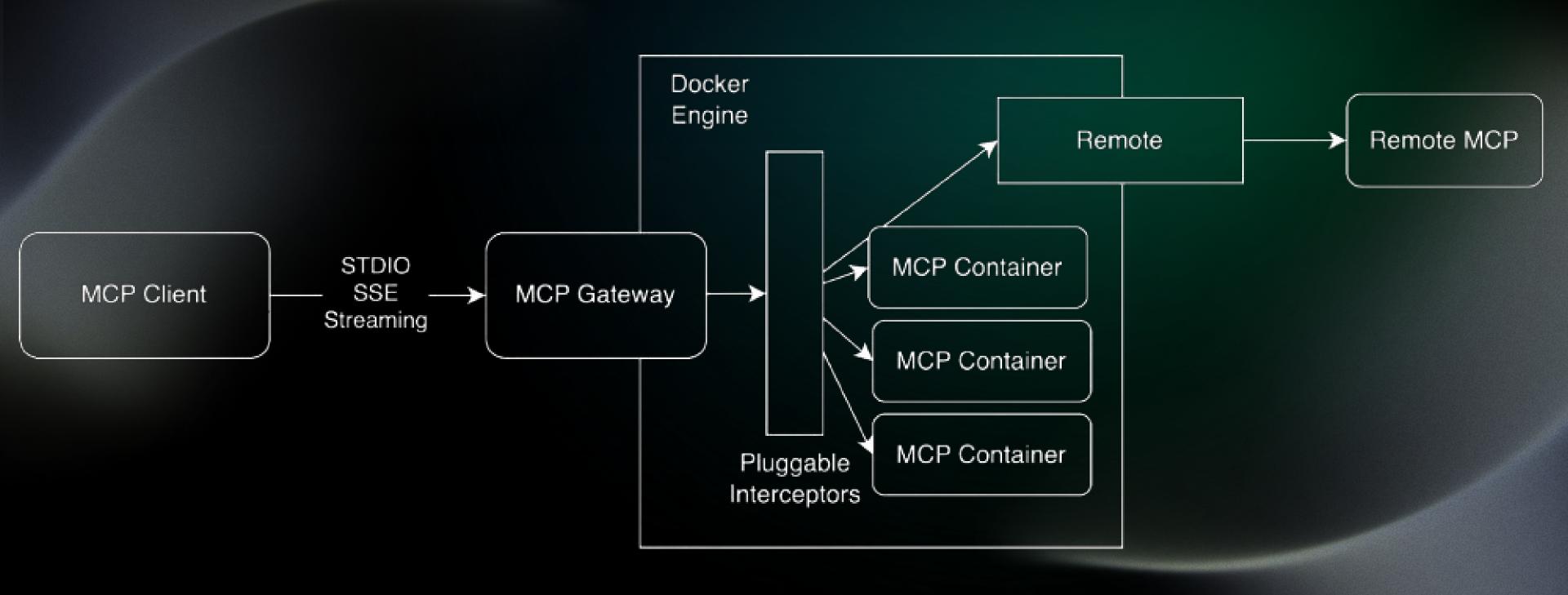
PLAN DE REPUESTA ANTE INCIDENTES

Armar un plan claro ante un eventual incidente, como aislar el sistema y controlar el impacto, roles y responsabilidades, monitoreo, y recuperación.

NO CONFIAR EN SERVIDORES DE TERCEROS

Siempre validar MCP servers, modelos, APIs de forma específica, escanear, testear, ver código y analizar si vale la pena tomar el riesgo intrínseco de usarlos.

MCP Gateway



Conclusiones

- 1. EL PROTOCOLO REVOLUCIONÓ LA INTEGRACIÓN DE IA CON SISTEMAS EXTERNOS, PERO INTRODUJO NUEVOS RIESGOS DE SEGURIDAD.
- 2. LA ADOPCIÓN REQUIERE UNA ESTRATEGIA ZERO TRUST Y MONITOREO CONSTANTE.
- 3. LA VALIDACIÓN DE SERVIDORES MCP Y CONTROL DE PERMISOS SON CRÍTICOS PARA MITIGAR AMENAZAS.

Muchas gracias!



philocyber.com