

Programa Estratégico de Ciberseguridad en la Era de la IA de Frontera

PREPARACIÓN ORGANIZATIVA ANTE LA ACELERACIÓN
DE VULNERABILIDADES Y ATAQUES AUTÓNOMOS

GIA | GRUPO DE
INTELIGENCIA
ARTIFICIAL

isms
forum

INTERNATIONAL
INFORMATION
SECURITY
COMMUNITY

Programa Estratégico de Ciberseguridad en la Era de la IA de Frontera

Preparación organizativa ante la aceleración de vulnerabilidades y ataques autónomos

Angel Pérez, Eduardo de Prado, Francisco Lázaro, ISMS Forum¹

I. Resumen ejecutivo

La evolución reciente de la inteligencia artificial ha introducido un cambio cualitativo en el ámbito de la ciberseguridad que trasciende cualquier mejora incremental observada en los últimos años. La aparición de modelos avanzados capaces de analizar código, razonar sobre sistemas complejos y ejecutar procesos de forma autónoma empieza a alterar de manera profunda el equilibrio tradicional entre atacantes y defensores y, en un breve espacio de tiempo, veremos una aceleración significativa de este desequilibrio.

Un ejemplo especialmente relevante es el caso del modelo Mythos, desarrollado por Anthropic. El fabricante anunció que, al ver sus capacidades, decidió retirar el modelo del acceso generalizado y limitarlo inicialmente a un conjunto restringido de organizaciones y entornos controlados, es lo que se ha llamado el proyecto Glasswing, precisamente por el impacto potencial que supone su uso indiscriminado. Este hecho es indicativo de dos realidades clave: por un lado, el nivel de madurez alcanzado por este tipo de herramientas y, por otro, la preocupación sobre su posible uso ofensivo a gran escala.

La restricción en el acceso a modelos de IA de frontera como Mythos introduce una asimetría relevante: las organizaciones que tienen acceso anticipado a estas herramientas pueden reforzar sus capacidades defensivas y reducir su exposición, mientras que el resto de las organizaciones quedan en desventaja. Esta dinámica, actualmente concentrada en EE. UU. y Reino Unido a nivel gubernamental y, previsiblemente en breve a organismos europeos, refuerza la necesidad de avanzar hacia una soberanía digital europea que garantice acceso, desarrollo y control sobre

capacidades de IA aplicadas a la ciberseguridad.

Aunque inicialmente esta ventaja es limitada en el tiempo, el conocimiento y las técnicas derivadas tienden a difundirse rápidamente. En consecuencia, el riesgo no desaparece, sino que se amplifica progresivamente a medida que estas capacidades se generalizan.

Esta dinámica se ve reforzada por un hecho reciente especialmente revelador: pocos días después del lanzamiento controlado de Mythos, se reportó un acceso no autorizado al modelo a través de un proveedor externo de Anthropic. El incidente expone dos realidades incómodas. La primera, que disponer de capacidades como Mythos no garantiza protección plena: la propia herramienta defensiva se convierte en superficie de ataque. La segunda, que ni siquiera su desarrollador puede asegurar que su uso quede limitado a los participantes autorizados de Project Glasswing; la difusión hacia actores maliciosos puede producirse antes de lo previsto, comprimiendo aún más la ventana de ventaja defensiva sobre la que se ha construido el modelo de acceso restringido.

Este fenómeno se traduce en lo que ya puede denominarse una **“tormenta de vulnerabilidades”**: un incremento exponencial en la detección, explotación y encadenamiento de fallos, que supera la capacidad tradicional de respuesta de las organizaciones. El tiempo entre la aparición de una vulnerabilidad y su explotación efectiva se reduce de semanas a horas o minutos, eliminando en la práctica la ventana clásica de parcheo.

En este contexto, se consolida una asimetría estructural: los atacantes operan a velocidad de máquina, mientras que

¹ Autores: [Angel Pérez](#), [Eduardo de Prado](#), [Francisco Lázaro](#) e [ISMS Forum](#)

muchas organizaciones continúan operando a velocidad humana. Esta brecha no es únicamente tecnológica, sino organizativa, operativa y cultural.

La ciberseguridad debe evolucionar hacia un modelo centrado en la resiliencia, la

II. Cambio de paradigma: implicaciones estructurales

Este cambio viene impulsado por la denominada inteligencia artificial de frontera, caracterizada por su capacidad de razonamiento, generación de código, análisis autónomo y ejecución de secuencias complejas con mínima intervención humana. Siendo sistemas que operan en el límite del conocimiento actual, tienen capacidades generales (no narrow AI) y pueden escalar rápidamente.

En este escenario de la utilización de este tipo de IA para identificar y explotar autónomamente vulnerabilidades, el Riesgo no se limita a incidentes puntuales: continuidad en el tiempo, impactos transversales (economía, seguridad, democracia), efectos en cascada y de difícil contención una vez desplegados (sin parches en muchos casos) y todo ello se traduce en un Riesgo Sistémico.

A esto viene a añadirse la incertidumbre radical en las capacidades futuras, en la realidad de evaluaciones de riesgos incompletas y la necesidad del enfoque de precaución.

Nos encontramos ante una transición desde un modelo de ciberseguridad basado en eventos discretos hacia un modelo continuo, en el que la materialización del riesgo es permanente y dinámica. El gap entre el descubrimiento de una vulnerabilidad y su explotación no solo se estrecha, sino que tiende a desaparecer.

Estas capacidades atacantes, como ya hemos dicho, no solo permiten identificar vulnerabilidades, sino también explotarlas, encadenarlas y adaptarse dinámicamente a los controles existentes. Como consecuencia, se rompe el modelo clásico de defensa basado en visibilidad previa y reacción.

automatización y la contención. La prevención deja de ser suficiente como mecanismo principal de control. La organización debe asumir que determinados compromisos serán inevitables y que la clave pasa a ser limitar su impacto y garantizar la continuidad operativa.

La organización debe asumir un nuevo supuesto operativo:

las vulnerabilidades existen en una cantidad mucho mayor de las que hasta ahora se conocían, serán explotadas y pueden estar siendo utilizadas antes de ser conocidas por las empresas

Esto trae consigo riesgos emergentes como el referido de automatización de ciberataques complejos con la escalada en capacidades ofensivas y la pérdida de control humano en procesos críticos y riesgos extremos donde se manifiesten los conflictos de interés de las organizaciones (objetivos desalineados), pérdida de control e impactos irreversibles a escala global.

A este negro panorama, se añaden viejos conocidos: la gobernanza va por detrás de la tecnología (regulación fragmentada, autoridades diversas, falta de estándares homogéneos en la UE y escasa capacidad de supervisión efectiva), falta de transparencia estructural (modelos cerrados, acceso limitado a la información, dependencias de declaraciones voluntarias) y un déficit de evaluación rigurosa tanto interno, como externo.

Esto implica una ruptura definitiva con los modelos tradicionales de:

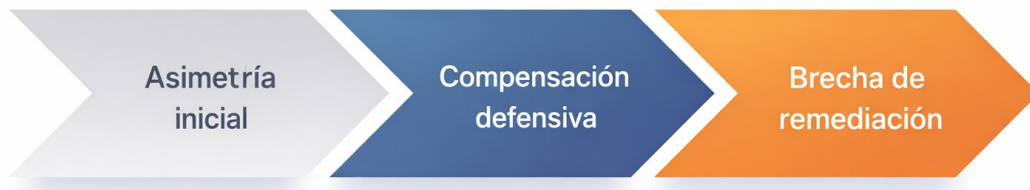
- priorización basada únicamente en CVSS (sistema estándar de la industria para evaluar la gravedad de una vulnerabilidad)
- inteligencia basada en CVE (catálogo internacional de vulnerabilidades que publica la organización MITRE)
- ciclos de parcheo periódicos

y obliga a adoptar modelos basados en:

- detección continua.
- explotación asumida.
- contención como objetivo.
- capacidad de auto-mejora o planificación estratégica avanzada.

III. Las tres ventanas de riesgo

La transición hacia un escenario en el que la IA de frontera está plenamente disponible no es instantánea ni homogénea. Conviene analizarla como una secuencia de tres ventanas temporales con perfiles de riesgo claramente diferenciados, cada una de las cuales exige una respuesta organizativa distinta.



Ventana 1. Asimetría inicial

Cuando un modelo con capacidades equivalentes a Mythos se generaliza en el mercado, los atacantes lo incorporan de forma inmediata. No median jerarquías de aprobación, comités de inversión, validaciones legales, ni procesos de adopción cultural. Las organizaciones defensoras, en cambio, deben atravesar todos esos ciclos antes de poder hacer un uso equivalente. Durante este periodo, el atacante opera con una herramienta avanzada mientras el defensor opera con la generación anterior. Es la ventana de máxima asimetría y de mayor exposición agregada.

Ventana 2. Compensación defensiva

Progresivamente, las mismas capacidades ofensivas se incorporan a los procesos defensivos: revisión continua de código, descubrimiento automatizado de vulnerabilidades en DevOps CI/CD (Integración e Implementación continua), parcheo asistido por IA y endurecimiento sistemático en el ciclo de desarrollo. Los grandes proveedores de software comienzan a entregar productos con menos vulnerabilidades nativas y tiempos de corrección significativamente reducidos. El riesgo agregado disminuye respecto a la ventana anterior, sin llegar a desaparecer.

Ventana 3. Brecha de remediación

Esta es la ventana más larga, menos visible y previsiblemente la más peligrosa. Aunque

los proveedores entreguen software más seguro y publiquen parches con mayor frecuencia, las organizaciones consumidoras siguen operando bajo procesos de despliegue manuales, ventanas de mantenimiento heredadas, dependencias funcionales, validaciones internas y restricciones operativas que ralentizan la aplicación efectiva de actualizaciones. El software de origen mejora. El tiempo medio entre la publicación de un parche y su despliegue en producción no.

El fenómeno es observable hoy en entornos OT, sistemas industriales y plataformas de infraestructura crítica, donde el despliegue de un parche implica coordinación con fabricantes, validaciones de compatibilidad y ventanas operativas restringidas. En la era de la IA de frontera, estas dinámicas adquieren una dimensión nueva: cada parche publicado se convierte en una señal pública de explotabilidad para todos los activos que aún no lo han aplicado. La capacidad de identificar el delta entre el código corregido y el código vulnerable es precisamente uno de los usos más eficientes de los modelos de frontera.

La consecuencia estratégica es directa. La ventaja competitiva en seguridad no la determinará quién dispone del mejor modelo, sino quién es capaz de remediar más rápido. La automatización del proceso de parcheo y la reducción del tiempo de despliegue dejan de ser mejoras operativas

V. Evolución del modelo de amenazas

La evolución tecnológica se traduce en un cambio profundo del modelo de amenazas.

Los ataques se aceleran, automatizan y escalan. El movimiento lateral se optimiza mediante IA. La superficie de ataque crece debido al software reutilizado y al desarrollo asistido por IA.

Pero el impacto más relevante no es técnico, sino operativo.

Este cambio afecta directamente a las áreas de sistemas. El parcheo, tradicionalmente percibido como una carga impuesta por seguridad, deja de ser una actividad puntual para convertirse en un proceso continuo. La

organización debe pasar de un modelo de mantenimiento a un modelo de respuesta permanente.

La incapacidad de adaptación en sistemas se convierte en un riesgo crítico.

La identidad se consolida como vector principal de ataque. La ingeniería social se vuelve más eficaz. La inteligencia tradicional pierde relevancia frente a un entorno de vulnerabilidades no catalogadas.

En conjunto, el modelo de amenazas deja de ser observable y pasa a ser **emergente y continuo**.

VI. Modelo objetivo de seguridad

1. VulnOps: operación continua de vulnerabilidades

La gestión de vulnerabilidades debe evolucionar hacia un modelo continuo. VulnOps no es una mejora del modelo existente, sino su sustitución. Implica:

- detección continua
- priorización basada en explotabilidad real
- integración en desarrollo
- automatización de corrección

Este modelo exige una transformación profunda de sistemas y operaciones.

El parcheo pasa a ser una capacidad crítica del negocio.

Además, VulnOps debe apoyarse en:

- análisis basado en LLM
- integración de código y despliegue continuo (CI/CD)
- revisión obligatoria de código

Sin esta evolución, la organización no puede operar en un entorno de descubrimiento continuo de vulnerabilidades.

2. Zero Trust real: contención efectiva del impacto

El Zero Trust debe evolucionar desde el plano conceptual a su implementación real.

La microsegmentación se convierte en el principal mecanismo de contención frente a

ataques automatizados. La identidad debe gestionarse dinámicamente y bajo riesgo.

El objetivo ya no es evitar accesos, sino limitar el impacto.

3. Automatización y defensa a velocidad de máquina

La defensa debe operar a la misma velocidad que el atacante. Esto implica:

- Nada sin telemetría
- Todo recolectado.
- Correlación en tiempo real y evaluado
- automatización de detección
- Colaboración CSIRT-SOC
- respuesta automatizada preautorizada por CSIRT

En modelos donde el SIEM y la capacidad de correlación, así como la gestión de incidentes dependen del CSIRT, la respuesta preautorizada debe entenderse como una extensión operativa de decisiones previamente definidas bajo su gobierno. La detección, los casos de uso y los playbooks se diseñan y validan en el ámbito del CSIRT, mientras que el SOC ejecuta de forma automatizada acciones de contención limitadas y trazables. Este modelo permite alcanzar una respuesta a velocidad de máquina sin diluir la responsabilidad ni el control estratégico del incidente, manteniendo la separación entre operación y dirección.

4. Capacidad ofensiva interna asistida por IA

La defensa pasiva frente a un atacante con capacidades de IA es estructuralmente insuficiente. La organización debe incorporar de forma sistemática las mismas herramientas que un atacante puede adquirir hoy mediante una suscripción mensual al alcance de cualquier individuo, y aplicarlas de forma continua sobre su propia superficie expuesta.

Esto implica desplegar procesos de pentesting y descubrimiento automatizado de vulnerabilidades, asistidos por modelos de IA generales y agenticos, que operen sobre activos públicos e internos críticos. El objetivo no es únicamente identificar fallos antes que el atacante, sino garantizar

paridad temporal: asegurar que cualquier capacidad ofensiva accesible en el mercado se está aplicando primero en contexto defensivo.

Este enfoque tiene además una implicación estratégica relevante. Cuando capacidades como las descritas en el caso Mythos se generalicen, las organizaciones que ya hayan industrializado el proceso podrán incorporar el nuevo modelo sustituyendo únicamente la pieza tecnológica. Las que no lo hayan hecho deberán construir el proceso completo bajo presión operativa. La preparación no se mide en herramientas adoptadas, sino en procesos preparados para incorporar nuevas herramientas sin rediseño.

VII. Plan estratégico de transformación

La transformación descrita en este documento no puede afrontarse mediante el reajuste interno de las partidas presupuestarias existentes. Las capacidades que exige el nuevo escenario, automatización a velocidad de máquina, microsegmentación efectiva, VulnOps continuo y uso defensivo de IA, no estaban contempladas en los modelos de inversión tradicionales en ciberseguridad.

Muchos CISOs se encontrarán con un presupuesto plenamente comprometido y sin margen real de maniobra. Pretender abordar este escenario reasignando partidas ya destinadas a controles existentes equivale a desproteger la operación actual para cubrir, de forma incompleta, la amenaza emergente. La organización debe asumir que la respuesta a este cambio de paradigma requiere una partida específica, identificada y trazable, dedicada a la construcción de las capacidades descritas. La alternativa no es un menor nivel de protección, sino una

desprotección estructural frente a un riesgo que la organización ya está asumiendo en la práctica.

Sobre esa base presupuestaria, la transformación requiere un enfoque estructurado:

Corto plazo

- revisión del modelo de riesgo
- uso de IA en análisis de código
- simulación de incidentes múltiples
- Descubrimiento de vulnerabilidades asistido por IA

Medio plazo

- inventario continuo
- refuerzo de identidad
- segmentación
- automatización inicial

Largo plazo

- VulnOps completo
- Zero Trust maduro
- CSIRT/SOC automatizado

VIII. Comunicación a la dirección y al consejo

La IA acelera negocio y riesgo simultáneamente. El mensaje clave es claro:

- el riesgo ya ha cambiado
- no actuar implica asumirlo

Las métricas deben evolucionar hacia:

- tiempo de detección

- tiempo de contención
- nivel de automatización


La dirección debe entender que: *la resiliencia es ahora una capacidad estratégica.*

IX. Conclusiones

La ciberseguridad en la era de la IA de frontera no evoluciona: cambia de naturaleza.

La diferencia entre organizaciones no será quién evita incidentes, sino quién es capaz de absorberlos, contenerlos y continuar operando.

La resiliencia, la automatización y la velocidad dejan de ser ventajas competitivas para convertirse en condiciones de supervivencia.



Programa Estratégico de Ciberseguridad en la Era de la IA de Frontera

PREPARACIÓN ORGANIZATIVA ANTE LA ACELERACIÓN
DE VULNERABILIDADES Y ATAQUES AUTÓNOMOS

GIA | GRUPO DE
INTELIGENCIA
ARTIFICIAL

isms
forum | INTERNATIONAL
INFORMATION
SECURITY
COMMUNITY