

# The Cybersecurity Psychology Framework: De la Teoría a la Práctica - Un Modelo de Evaluación de Vulnerabilidades Pre-Cognitivas con Validación de Caso

## TECHNICAL REPORT

Giuseppe Canale, CISSP

Independent Researcher

kaolay@gmail.com

ORCID: 0009-0007-3263-6897

Framework Website: <https://cpf3.org>

December 20, 2025

## 1 Introducción

El fracaso persistente de las medidas de cybersecurity a pesar del crecimiento exponencial de las inversiones revela una incomprendión fundamental del espacio del problema. Mientras el gasto global en cybersecurity supera los \$150 mil millones anualmente[7], las violaciones exitosas continúan aumentando, con factores humanos que contribuyen a más del 85% de los incidentes[21]. Esta paradoja sugiere que nuestro enfoque del elemento humano en la cybersecurity permanece fundamentalmente defectuoso, tratando la conciencia consciente y la toma de decisiones racional como los puntos primarios de intervención cuando las neurociencias demuestran claramente que la mayoría de las decisiones humanas ocurren bajo el umbral de la conciencia.

Los recientes avances en las neurociencias han revolucionado nuestra comprensión de los procesos decisionales. El trabajo pionero de Libet[14] ha demostrado que la actividad cerebral indicando una decisión ocurre 300-500 milisegundos antes de la conciencia consciente de esa decisión. Este hallazgo, replicado y extendido por Soon et al.[20] usando tecnología fMRI, revela que en el momento en que un empleado decide conscientemente si hacer clic en un enlace de phishing, su cerebro ya ha iniciado la acción. Estos procesos pre-cognitivos operan a través de interacciones complejas entre el sistema de detección de amenazas de la amígdala y el control ejecutivo de la corteza prefrontal, con la respuesta más rápida de la amígdala que a menudo sobrescribe el análisis racional[13].

El contexto organizacional agrega otra capa de complejidad que los frameworks actuales fallan en abordar. Las

organizaciones no son meramente colecciones de individuos sino sistemas complejos con sus propias dinámicas inconscientes. El trabajo seminal de Bion sobre el comportamiento de grupo[3] ha demostrado que los grupos bajo estrés retroceden a estados de asunciones básicas que sobrescriben el juicio individual. Cuando una organización enfrenta una amenaza cyber, podría colectivamente desplazarse a modalidad de dependencia, buscando un protector omnípotente (a menudo manifiesto como excesiva dependencia de los proveedores de security), o modalidad fight-flight, percibiendo todas las amenazas como externas mientras ignora los riesgos insider. Estos procesos inconscientes a nivel de grupo crean vulnerabilidades sistemáticas que ninguna cantidad de training de security individual puede abordar.

El Cybersecurity Psychology Framework (CPF) representa un cambio de paradigma en el abordaje de estos desafíos. En lugar de intentar reforzar la toma de decisiones consciente a través de training de conciencia, el CPF mapea los procesos pre-cognitivos e inconscientes que efectivamente guían los comportamientos relevantes para la security. Integrando la teoría psicoanalítica de las relaciones objetales, que explica cómo categorizamos y respondemos inconscientemente a las amenazas basándonos en experiencias tempranas, con la comprensión de la psicología cognitiva de sesgos sistemáticos y heurísticas, el CPF proporciona un modelo comprensivo para predecir y prevenir los fracasos de security antes de que ocurran.

## 2 Fundamento Teórico

### 2.1 El Fracaso de las Intervenciones a Nivel Consciente

Los programas tradicionales de conciencia de la security operan sobre la asunción implícita del modelo del actor racional—que los individuos, cuando provistos de información sobre los riesgos y las respuestas apropiadas, modificarán su comportamiento consecuentemente[1]. Este modelo subyace virtualmente todo el training de security actual, desde las simulaciones de phishing hasta las políticas de contraseñas. Sin embargo, décadas de investigación a través de múltiples disciplinas demuestran la inadecuación fundamental de este enfoque.

La evidencia neurocientífica es particularmente convincente. La hipótesis del marcador somático de Damasio[6] revela que las respuestas emocionales y corporales a los estímulos ocurren antes y a menudo sobrescriben el análisis racional. En el contexto de la cybersecurity, esto significa que la reacción instintiva de un empleado a un email—influenciada por factores como familiaridad del remitente, presión temporal o contenido emocional—determina su respuesta antes de que ocurra la evaluación consciente de los indicadores de security. Además, la teoría del proceso dual de Kahneman[9] demuestra que bajo condiciones típicas de los ambientes de trabajo modernos—alta carga cognitiva, presión temporal, multitarea—el Sistema 1 (rápido, automático, intuitivo) domina el Sistema 2 (lento, deliberado, analítico). Las decisiones de security, requiriendo análisis cuidadoso de indicadores sutiles, son precisamente el tipo que sufre mayormente bajo estas condiciones.

El fracaso de las intervenciones a nivel consciente no es meramente teórico sino empíricamente demostrable. Beaument et al.[2] han introducido el concepto del "presupuesto de compliance"—la cantidad finita de esfuerzo que los empleados gastarán en los comportamientos de security antes de experimentar fatiga y desvinculación. Una vez que este presupuesto está agotado, los empleados comienzan a tomar atajos, independientemente de su conocimiento de security. Esto explica por qué los incidentes de security a menudo aumentan durante períodos de alto estrés como lanzamientos de productos o cierres financieros, cuando los recursos cognitivos están agotados y el presupuesto de compliance ya está gastado en las tareas de trabajo primarias.

### 2.2 Contribuciones Psicoanalíticas a la Comprensión de las Vulnerabilidades de Security

#### 2.2.1 Las Asunciones Básicas de Bion y la Security Organizacional

La teoría de las dinámicas de grupo de Wilfred Bion[3] proporciona intuiciones cruciales sobre los fracasos de security organizacionales que los frameworks focalizados en el individuo pierden enteramente. Bion ha observado que los grupos, cuando enfrentan situaciones ansiógenas, adoptan inconscientemente una de las tres asunciones básicas que sirven como estructuras defensivas contra esa ansiedad. Estas asunciones operan bajo la conciencia consciente pero influencian profundamente el comportamiento de grupo y la toma de decisiones.

La primera asunción básica, Dependencia (baD), se manifiesta en contextos de cybersecurity como creencia inconsciente de que alguna fuerza omnipotente proporcionará protección. Las organizaciones en modalidad de dependencia exhiben comportamientos característicos: excesiva dependencia de los proveedores de security con expectativas irrealistas de sus capacidades, abdicación de la responsabilidad de security al departamento de IT mientras otros departamentos permanecen pasivos, y pensamiento mágico sobre las herramientas de security como soluciones completas. Por ejemplo, después de implementar un costoso firewall de nueva generación, una organización podría inconscientemente relajar otras medidas de security, creyendo que el firewall proporciona protección comprensiva. Esta dependencia crea vulnerabilidades ya que los empleados asumen que "el sistema" capturará todas las amenazas, reduciendo su propia vigilancia.

Fight-Flight (baF), la segunda asunción básica, aparece cuando las organizaciones perciben las amenazas como enemigos externos requiriendo defensa agresiva o evitamiento completo. En modalidad fight, las organizaciones podrían implementar políticas de security draconianas que los empleados evaden, creando shadow IT y workarounds que introducen nuevas vulnerabilidades. En modalidad flight, las organizaciones podrían evitar confrontar las realidades de security, posponiendo actualizaciones, ignorando reportes de vulnerabilidades o manteniendo sistemas legacy porque enfrentarlos parece demasiado amenazante. La asunción fight-flight ciega críticamente a las organizaciones a las amenazas insider—ya que el enemigo está conceptualizado como externo, las vulnerabilidades internas permanecen invisibles.

La tercera asunción, Pairing (baP), involucra la fantasía inconsciente del grupo de que algún evento futuro o unión resolverá todos los problemas. En la cybersecurity, esto se manifiesta como adquisición perpetua de her-

ramientas—buscando siempre la próxima solución de security que finalmente proporcionará protección completa. Las organizaciones en modalidad pairing exhiben ciclos de esperanza y desilusión con cada nueva iniciativa de security, nunca enfrentando vulnerabilidades fundamentales porque la solución “real” está siempre a la vuelta de la esquina.

## 2.2.2 Relaciones Objetales Kleinianas y Percepción de la Security

La teoría de las relaciones objetales de Melanie Klein[12] elucida cómo las organizaciones inconscientemente dividen su panorama de security en objetos “todo buenos” y “todo malos”, un mecanismo de defensa primitivo que crea peligrosos puntos ciegos. Esta división opera a través de identificación proyectiva, donde aspectos no deseados del sí mismo son proyectados sobre objetos externos, distorsionando fundamentalmente la percepción de la amenaza.

En contextos organizacionales, esta división se manifiesta en dicotomías netas. Los empleados son categorizados como insiders confiables o amenazas potenciales, con poco reconocimiento de la realidad compleja de que individuos confiables pueden ser comprometidos o cometer errores. Los sistemas son similarmente divididos en “nuestra red segura” versus “el internet peligroso”, ignorando la porosidad de los límites de red modernos. Esta categorización primitiva explica por qué las organizaciones a menudo fallan en implementar arquitecturas zero-trust—el concepto de que la confianza debe ser continuamente verificada en lugar de asumida basándose en la posición de red contradice la necesidad inconsciente de límites bueno/malo claros.

El mecanismo de proyección es igualmente problemático. Las organizaciones proyectan sus propios impulsos agresivos sobre los atacantes externos, imaginando a los hackers como fuerzas malevolas mientras fallan en reconocer sus propias prácticas empresariales agresivas que podrían motivar ataques. Proyectan su propia vulnerabilidad sobre los usuarios, culpando a “usuarios estúpidos” por fracasos de security mientras niegan debilidades arquitecturales sistemáticas. Esta proyección sirve una función defensiva, manteniendo la auto-imagen de la organización como competente y segura mientras localiza todos los problemas externamente.

El concepto de Klein de la posición paranoide-esquizoide versus la posición depresiva ofrece intuiciones adicionales. Las organizaciones en la posición paranoide-esquizoide experimentan ansiedad extrema sobre las amenazas de security, oscilando entre vigilancia paranoica y retiro esquizoide. No pueden tolerar ambigüedad o incertidumbre, llevando a parálisis de security o respuestas reactivas mal consideradas. Moverse a la posición depre-

siva—donde aspectos buenos y malos pueden ser integrados, y la pérdida puede ser llorada—es esencial para una postura de security madura pero requiere elaboración de ansiedad organizacional significativa.

## 2.2.3 El Espacio Transicional de Winnicott y los Ambientes Digitales

El concepto de espacio transicional de Donald Winnicott[22]—el área psicológica entre fantasía interna y realidad externa—proporciona intuiciones únicas sobre vulnerabilidades específicas a los ambientes digitales. El ciberespacio funciona como un espacio transicional donde los límites entre real e imaginario, sí mismo y otro, se vuelven difusos. Esta difuminación crea vulnerabilidades específicas que los frameworks de security tradicionales fallan en abordar.

En el espacio transicional, florecen fantasías omnipotentes. Los usuarios podrían sentirse invulnerables detrás de nombres de usuario, tomando riesgos que nunca tomarían en el espacio físico. Podrían creer que pueden controlar su huella digital completamente o, al contrario, no tener control alguno. Estas fantasías influencian los comportamientos de security: excesiva confianza lleva a acciones riesgosas, mientras impotencia aprendida resulta en nihilismo de security—“¿por qué preocuparse de la security cuando los hackers pueden entrar de todos modos?”

La naturaleza transicional del espacio digital influye también la formación de la identidad y la gestión de límites. Los empleados podrían desarrollar personas online que difieren de sus identidades profesionales, creando vulnerabilidades cuando estos mundos colisionan. Los perfiles de redes sociales pensados para uso personal se convierten en vectores de ataque para compromiso profesional. La cualidad lúdica y experimental del espacio transicional—esencial para creatividad e innovación—entra en conflicto con requisitos de security para comportamiento consistente y cauteloso.

## 2.2.4 La Sombra Junguiana y el Inconsciente Colectivo en la Cybersecurity

El concepto de sombra de Carl Jung[8]—los aspectos reprimidos y negados de la personalidad—ilumina cómo las organizaciones crean vulnerabilidades a través de lo que rechazan reconocer sobre sí mismas. La sombra organizacional contiene todas las cualidades que la organización no puede aceptar: competitividad agresiva negada en favor de “cultura colaborativa”, capacidades de vigilancia ocultas detrás de “cuidado de los empleados”, o explotación de datos enmascarada como “servicio al cliente”.

Estos elementos de la sombra no desaparecen; son proyectados sobre los atacantes que se convierten en los

portadores de las cualidades no reconocidas de la organización. Los hackers son imaginados como poseedores de habilidades sobrehumanas, reflejando las fantasías omnipotentes de la organización misma. Son vistos como puramente destructivos, portando la agresión negada de la organización. Esta proyección previene evaluación realista de la amenaza—si los atacantes están mitificados como extraordinarios, entonces medidas de security ordinarias parecen fútiles, justificando inversión de security inadecuada.

El inconsciente colectivo, el concepto de Jung de patrones psicológicos heredados compartidos a través de la humanidad, se manifiesta en la cybersecurity a través de respuestas arquetípicas a las amenazas. El arquetipo del Guerrero guía posturas de security agresivas y retórica de “guerra cyber”. El arquetipo del Trickster aparece tanto en los atacantes como en los defensores, con profesionales de security a veces inconscientemente identificándose con hackers. El arquetipo de la Sombra encarna todo lo que la organización teme y niega sobre sí misma, proyectado sobre los actores de amenaza.

## 2.3 Integración de Psicología Cognitiva

### 2.3.1 Teoría del Proceso Dual en Contextos de Security

El framework Sistema 1/Sistema 2 de Kahneman[9] revela vulnerabilidades específicas en la toma de decisiones de security que emergen de la arquitectura fundamental de la cognición humana. El Sistema 1, operando automáticamente e inconscientemente, procesa información a través de reconocimiento de patrones y respuesta emocional, tomando decisiones en milisegundos basadas en heurísticas desarrolladas a través de evolución y experiencia. El Sistema 2, consciente y deliberado, puede sobrescribir el Sistema 1 pero requiere recursos cognitivos significativos y tiempo—lujos raramente disponibles en los ambientes de trabajo modernos.

En contextos de cybersecurity, el Sistema 1 domina a través de diversos mecanismos. La heurística de disponibilidad causa incidentes de security recientes o memorables a influenciar desproporcionadamente las decisiones de security. Después de un ataque ransomware publicizado sobre una organización similar, las empresas podrían sobre-invertir en defensas ransomware mientras descuidan otros vectores. La heurística del afecto vincula decisiones de security a estados emocionales: el miedo guía reacción excesiva, mientras el confort genera complacencia. El efecto anclaje causa incidentes de security iniciales a establecer expectativas para todas las amenazas futuras, potencialmente perdiendo patrones de ataque evolutivos.

Las limitaciones del Sistema 2 agravan estas vulnera-

bilidades. La carga cognitiva de la complejidad de security—contraseñas múltiples, sistemas de autenticación, protocolos de security—agota los recursos mentales necesarios para análisis cuidadoso. El agotamiento del ego de la vigilancia constante reduce la compliance de security con el tiempo, explicando por qué los incidentes de security a menudo ocurren al final del día o fin de semana cuando los recursos cognitivos están agotados. El razonamiento motivado lleva a los individuos a racionalizar ataques de security cuando entran en conflicto con objetivos de productividad, construyendo justificaciones elaboradas para comportamientos inseguros.

### 2.3.2 Los Principios de Influencia de Cialdini como Vectores de Ataque

Los seis principios de influencia de Robert Cialdini[5] mapean directamente sobre tácticas de social engineering, revelando cómo los atacantes explotan la programación social humana fundamental. Estos principios operan bajo la conciencia consciente, activando respuestas de compliance automáticas que bypassan el training de security.

Reciprocidad, la obligación de devolver favores, habilita ataques quid pro quo donde los atacantes proporcionan algo de valor—información útil, asistencia, o incluso simpatía—antes de hacer su solicitud. Un atacante podría ayudar a un empleado a resolver un problema técnico, creando una obligación que hace psicológicamente difícil rechazar una solicitud posterior de credenciales. La presión de compromiso y coherencia empuja a los individuos a alinear acciones con compromisos previos, habilitando ataques de escalación gradual. Un atacante podría primero solicitar información inocua, luego progresivamente datos más sensibles, confiando en la necesidad del objetivo de permanecer coherente con su cooperación inicial.

La prueba social, la tendencia a seguir el comportamiento de otros, habilita ataques que hacen referencia a acción colectiva: “Todos en contabilidad ya han proporcionado esta información”. La influencia de la autoridad habilita ataques de suplantación, con tasas de éxito que superan el 90

### 2.3.3 Carga Cognitiva y Degradación del Rendimiento de Security

La identificación de George Miller de los límites de capacidad cognitiva[17]—el “número mágico siete, más o menos dos”—revela vínculos fundamentales que crean vulnerabilidades de security. Los requisitos de security modernos superan rutinariamente estos límites, forzando ataques cognitivos que los atacantes explotan.

Los requisitos de las contraseñas ejemplifican este problema. Las organizaciones que requieren contraseñas

únicas y complejas para sistemas múltiples exceden la capacidad de memoria, forzando prácticas inseguras: reuso de contraseñas, escritura de credenciales, o uso de patrones predecibles. La carga cognitiva de recordar contraseñas múltiples agota los recursos mentales necesarios para la detección de amenazas. La proliferación de herramientas de security agrava este problema. Cuando los equipos de security monitorean docenas de dashboards y sistemas de alerta, señales importantes se pierden en el ruido. La fatiga de alertas se desarrolla cuando la capacidad cognitiva es superada, llevando a tasas de respuesta disminuidas y tiempos de respuesta aumentados a amenazas genuinas.

El multitarea degrada adicionalmente el rendimiento de security. El cambio de contexto entre tareas incurre costos cognitivos, creando ventanas de vulnerabilidad durante transiciones. El residuo de tareas previas interfiere con decisiones de security actuales. Bajo alta carga cognitiva, los individuos retornan a respuestas habituales, que podrían ser inseguras, y se vuelven susceptibles al social engineering ya que los recursos cognitivos restantes son insuficientes para escepticismo.

## 2.4 Vulnerabilidades Psicológicas Específicas de la AI

### 2.4.1 Antropomorfización y Transferencia de Confianza

A medida que los sistemas AI se vuelven integrales a las operaciones de cybersecurity, nuevas vulnerabilidades psicológicas emergen de las tendencias humanas a antropomorizar entidades no-humanas. Los humanos atribuyen naturalmente características humanas, intenciones y emociones a los sistemas AI, creando relaciones de confianza explotables.

Esta antropomorfización se manifiesta de diversas maneras. Los profesionales de security desarrollan "relaciones" con herramientas de security AI, confiando en su "juicio" más allá de sus capacidades efectivas. Los usuarios atribuyen intenciones benévolas a los asistentes AI, compartiendo información sensible que no proporcionarían a extraños humanos. El efecto uncanny valley—incomodidad con AI casi-pero-no-del todo-humana—puede ser explotado haciendo que los sistemas AI parezcan ya sea más o menos humanos para manipular los niveles de confianza.

Los mecanismos de transferencia de confianza agravan estas vulnerabilidades. La transferencia de autoridad ocurre cuando los sistemas AI heredan confianza de sus creadores u operadores: "Es la AI de Google, entonces debe ser segura". La transferencia de competencia asume que la AI competente en un dominio es confiable en todos los dominios. La transferencia emocional se desarrolla a

medida que los usuarios forman apegos a personalidades AI, haciéndolos vulnerables a manipulación a través de estas relaciones sintéticas.

### 2.4.2 Sesgo de Automatización y Atrofia de las Competencias

El sesgo de automatización—la tendencia a confiar excesivamente en sistemas automatizados—crea vulnerabilidades críticas en ambientes de security aumentados por AI. Los equipos de security se apoyan cada vez más en recomendaciones AI sin evaluación crítica, asumiendo que la AI tiene acceso a más información o capacidades de análisis superiores. Esta deferencia ocurre incluso cuando la intuición humana sugiere lo contrario, suprimiendo intuición humana valiosa.

El riesgo moral de las herramientas de security AI reduce la vigilancia humana. Si la AI está monitoreando para amenazas, la atención humana naturalmente disminuye—un fenómeno observado en incidentes relacionados con piloto automático en la aviación. La atrofia de competencias sigue a medida que los profesionales de security pierden práctica en la detección manual de amenazas y el análisis. Cuando los sistemas AI fallan o son comprometidos, los operadores humanos carecen de las competencias para compensar, creando ventanas de vulnerabilidad catastrófica.

Los loops de retroalimentación entre sesgos humanos y AI amplifican las vulnerabilidades. Los sistemas AI entrenados sobre datos sesgados perpetúan y legitiman esos sesgos, que los humanos luego aceptan como verdad objetiva porque "la AI lo dijo". Estos sesgos reforzados se convierten en puntos ciegos que los atacantes pueden explotar, sabiendo que tanto defensas humanas como AI comparten las mismas debilidades.

## 3 La Arquitectura del Modelo CPF

### 3.1 Filosofía de Diseño y Principios de Implementación

La arquitectura del Cybersecurity Psychology Framework refleja un desplazamiento fundamental de la evaluación de security reactiva a la predictiva. A diferencia de los frameworks tradicionales que catalogan vulnerabilidades existentes o incidentes pasados, el CPF mapea las precondiciones psicológicas que habilitan fracasos de security futuros. Esta capacidad predictiva emerge de la comprensión de que estados psicológicos y dinámicas de grupo crean patrones consistentes de vulnerabilidad que se manifiestan antes de que ocurran incidentes de security efectivos.

El diseño del framework preservando la privacidad aborda los desafíos éticos inherentes en la evaluación

psicológica dentro de contextos organizacionales. Todas las mediciones operan a niveles agregados, con una unidad mínima de diez individuos, previniendo perfilamiento individual mientras mantiene validez estadística. Las técnicas de privacidad diferencial con inyección de ruido ( $\varepsilon = 0.1$ ) aseguran que incluso con acceso a los datos de salida, los estados psicológicos individuales no puedan ser retro-ingenierizados. Esta elección de diseño no es meramente ética sino práctica—los empleados que temen vigilancia psicológica conscientemente o inconscientemente alterarán su comportamiento, invalidando las evaluaciones.

El enfoque agnóstico a la implementación asegura la aplicabilidad del CPF a través de contextos organizacionales diversos. En lugar de prescribir herramientas o procedimientos de security específicos, el CPF identifica estados de vulnerabilidad que pueden ser abordados a través de varias intervenciones. Esta flexibilidad permite a las organizaciones integrar el CPF con frameworks y herramientas de security existentes mientras respetan sus culturas, vínculos y capacidades únicas.

### 3.2 Estructura del Framework: La Matriz $10 \times 10$

Los 100 indicadores del CPF están organizados en una matriz  $10 \times 10$  que balancea completitud con aplicabilidad práctica. Cada categoría representa un dominio psicológico distinto con su propio fundamento teórico y soporte empírico, mientras los diez indicadores dentro de cada categoría proporcionan capacidad de evaluación granular sin complejidad abrumadora.

### 3.3 Categoría 1: Vulnerabilidades Basadas en la Autoridad

Las vulnerabilidades basadas en la autoridad emergen de tendencias humanas profundamente arraigadas a obedecer figuras de autoridad percibidas, un fenómeno dramáticamente demostrado en los experimentos de Milgram[16]. En contextos de cybersecurity, estas vulnerabilidades son particularmente peligrosas porque bypassan la toma de decisiones racional de security a través de la activación de respuestas de compliance automáticas.

El primer indicador (1.1), compliance no cuestionante con aparente autoridad, captura la manifestación más directa de esta vulnerabilidad. Cuando un atacante suplanta exitosamente una figura de autoridad—ya sea a través de spoofing de email, manipulación de voz o presencia física—los objetivos se conforman con solicitudes que de otro modo activarían preocupaciones de security. Por ejemplo, en el hack de Twitter del 2020, los atacantes obtuvieron acceso a cuentas de alto perfil llamando a empleados de Twitter y afirmando ser de la security de IT,

solicitando resets de contraseña. Los empleados se conformaron sin verificación, a pesar del training de security, porque la afirmación de autoridad activó obediencia automática.

La difusión de responsabilidad (1.2) en estructuras jerárquicas crea vulnerabilidades donde cada nivel asume que la security es responsabilidad de alguien más. Los ejecutivos senior asumen que IT maneja la security, IT asume que el management establece las políticas, y los empleados en primera línea asumen que ambos niveles proporcionan protección. Esta difusión crea gaps que los atacantes explotan, sabiendo que responsabilidad poco clara significa que nadie toma propiedad. La susceptibilidad a la suplantación de figuras de autoridad (1.3) se extiende más allá de la simple obediencia para incluir el fracaso en verificar afirmaciones de autoridad. Las organizaciones raramente entran a los empleados a desafiar o verificar la autoridad, creando un vector de ataque donde falsa autoridad va no cuestionada.

El fenómeno de bypassar la security para la conveniencia del superior (1.4) representa una vulnerabilidad particularmente insidiosa. Cuando los ejecutivos solicitan excepciones de security—usando dispositivos personales, evitando VPN, o compartiendo credenciales—los subordinados se conforman a pesar de conocer los riesgos. Esto crea tanto vulnerabilidades directas como modela comportamiento inseguro a través de la organización. La compliance basada en el miedo sin verificación (1.5) ocurre cuando la amenaza de desagrado de la autoridad sobrescribe los protocolos de security. Los atacantes explotan esto creando urgencia e implicando consecuencias para no-compliance: "El CEO necesita esto inmediatamente o el acuerdo falla".

Los efectos de gradiente de la autoridad (1.6) inhiben el reporte de security cuando los subordinados temen desafiar prácticas inseguras de los superiores. En la sanidad, los gradientes de autoridad entre médicos y enfermeros han sido vinculados a errores médicos; en la cybersecurity, gradientes similares previenen al personal junior de reportar violaciones de security del personal senior. La deferencia a afirmaciones de autoridad técnica (1.7) crea vulnerabilidades cuando los atacantes usan jerga técnica para establecer credibilidad. El personal no-técnico, sintiéndose inadecuado para desafiar afirmaciones técnicas, se conforma a solicitudes que no comprenden.

La normalización de las excepciones ejecutivas (1.8) ocurre cuando las reglas de security rutinariamente no se aplican al liderazgo senior, creando tanto vulnerabilidades prácticas como socavando la cultura de security. La prueba social basada en la autoridad (1.9) amplifica otros efectos de autoridad cuando figuras de autoridad múltiples modelan comportamiento inseguro, normalizando violaciones de security. La escalación de autoridad en crisis

(1.10) describe cómo las vulnerabilidades basadas en la autoridad se intensifican durante crisis cuando los procedimientos de verificación normales son suspendidos y las afirmaciones de autoridad ganan poder adicional.

La Tabla 1 proporciona tres ejemplos de cómo los indicadores CPF pueden ser operacionalizados en Indicadores de Riesgo Comportamental (BRI) cuantificables. Estos BRI aprovechan datos agregados y anonimizados de logs de IT estándar. Los umbrales de puntuación son estimaciones iniciales basadas en estudio piloto y benchmarks de sector.

El bypass de security inducido por la urgencia (2.1) ocurre cuando la presión temporal causa a los individuos a saltar pasos de security percibidos como ralentización del progreso. Los atacantes explotan esto creando urgencia artificial: "Esta factura debe ser pagada dentro de la hora para evitar interrupción del servicio". Bajo presión temporal, el pensamiento del Sistema 2 se desconecta, dejando solo las heurísticas rápidas pero vulnerables del Sistema 1. La degradación cognitiva de presión temporal (2.2) describe el deterioro más amplio de la toma de decisiones bajo estrés temporal. La investigación muestra que la presión temporal reduce la capacidad de memoria de trabajo, compromete el juicio, y aumenta la toma de riesgos—todo beneficioso para los atacantes.

La aceptación del riesgo guiada por la fecha límite (2.3) se manifiesta cuando las fechas límite en aproximación causan a las organizaciones a aceptar riesgos de security que normalmente rechazarían. Lanzamientos de productos, cierres financieros y completamientos de proyectos se convierten en ventanas de vulnerabilidad ya que la security toma el segundo lugar a la entrega. El sesgo del presente (2.4) en las inversiones de security lleva a las organizaciones a sub-invertir en la prevención de amenazas futuras mientras sobre-responden a incidentes actuales. Esto crea vulnerabilidades cíclicas donde las amenazas de ayer están sobre-defendidas mientras las de mañana son ignoradas.

El descuento hiperbólico (2.5) causa a las organizaciones a subvalorar dramáticamente los beneficios de security futuros relativos a los costos presentes. Una medida de security que prevenga una violación el próximo año parece menos valiosa que la conveniencia menor hoy, incluso cuando el costo futuro supera ampliamente los ahorros presentes. Los patrones de agotamiento temporal (2.6) crean ventanas de vulnerabilidad predecibles. La vigilancia de security se degrada a través del día laboral, la semana laboral, y los ciclos de proyecto. Los atacantes que comprenden estos patrones temporalizan sus ataques para probabilidad máxima de éxito.

Las ventanas de vulnerabilidad time-of-day (2.7) reflejan ritmos circadianos en el rendimiento cognitivo. Temprano en la mañana y tarde en la tarde muestran susceptibilidad aumentada a phishing y social engineering. Los

lapsos de security de fin de semana y vacaciones (2.8) ocurren cuando personal reducido y vigilancia relajada crean oportunidades para intrusión no detectada. Violaciones mayores a menudo comienzan durante vacaciones cuando las capacidades de respuesta están minimizadas. Las ventanas de explotación del cambio de turno (2.9) apuntan a la confusión y los gaps informativos durante transiciones de personal. La presión de coherencia temporal (2.10) describe cómo inversiones temporales pasadas crean presión a continuar prácticas inseguras en lugar de reconocer esfuerzo desperdiciado—la falacia de los costos hundidos aplicada a la security.

### **3.4 Categoría 3: Vulnerabilidades de Influencia Social**

Las vulnerabilidades de influencia social explotan necesidades humanas fundamentales para conexión social, coherencia y pertenencia. Estas vulnerabilidades son particularmente potentes porque operan a través de mecanismos sociales positivos que las organizaciones efectivamente quieren alentar, creando conflictos entre security y cultura.

La explotación de la reciprocidad (3.1) arma la norma universal de devolver favores. Los atacantes establecen relaciones recíprocas a través de pequeños favores antes de hacer solicitudes malevolas. El malestar psicológico de rechazar a alguien que te ha ayudado sobrescribe el training de security. Las trampas de escalación del compromiso (3.2) explotan el principio de coherencia, donde pequeños compromisos iniciales llevan a aquellos más grandes. Un atacante podría primero solicitar información públicamente disponible, luego progresivamente datos más sensibles, confiando en la necesidad del objetivo de permanecer coherente con la cooperación inicial.

La manipulación de la prueba social (3.3) explota la tendencia a seguir el comportamiento de otros, especialmente bajo incertidumbre. Los atacantes afirman "todos los demás ya han proporcionado esta información" o crean prueba social falsa a través de cuentas comprometidas. El superamiento de la confianza basado en el gusto (3.4) ocurre cuando sentimientos positivos hacia alguien causan a los protocolos de security a ser ignorados. Los atacantes investigan intereses, antecedentes y relaciones de los objetivos para establecer rapport que desarma la sospecha.

Las decisiones guiadas por la escasez (3.5) explotan el miedo de perder oportunidades. Ofertas de tiempo limitado, acceso exclusivo, o amenazar remoción de recursos activan decisiones rápidas sin verificación apropiada. La explotación del principio de unidad (3.6) explota identidad compartida para bypassar la security. Los atacantes afirman pertenencia al mismo grupo—alumni, asociación profesional o causa social—para establecer confianza.

Table 1: Indicadores de Riesgo Comportamental Ejemplares (BRI) con Puntuación Cuantitativa

Nombre BRI	Categoría	Lógica de Medición	Puntuación
Compliance No Cuestionante	Autoridad (1.1)	$\frac{\text{No Verificado}}{\text{Total}} \times 100$	Verde: < 5% Amarillo: 5-15% Rojo: > 15%
Procrastinación Patch	Temporal (2.4)	$I_{PP} = \frac{\max(0, D - 30)}{10}$	Verde: $I_{PP} < 1$ Amarillo: $1 \leq I_{PP} < 3$ Rojo: $I_{PP} \geq 3$
Tasa de Descarte de Alertas	Cognitivo (5.1)	$\frac{\text{Descartados}}{\text{Total}} \times 100$	Verde: < 10% Amarillo: 10-25% Rojo: > 25%

La compliance de presión de pares (3.7) ocurre cuando la presión social de colegas sobrescribe preocupaciones de security. Si todos comparten contraseñas para conveniencia, rechazar marca a uno como no cooperativo. La conformidad a normas inseguras (3.8) describe cómo prácticas inseguras se vuelven normalizadas a través de transmisión social. Una vez que la masa crítica adopta una práctica insegura, se convierte en el estándar. Las amenazas de identidad social (3.9) explotan miedos de exclusión social o desafío de identidad. Los atacantes amenazan posición social o pertenencia al grupo para coaccionar compliance. Los conflictos de gestión de la reputación (3.10) emergen cuando las medidas de security entran en conflicto con preocupaciones de reputación, como reportar una violación que podría dañar la imagen organizacional.

### 3.5 Categoría 4: Vulnerabilidades Afectivas

Las vulnerabilidades afectivas emergen de cómo los estados emocionales influencian decisiones y comportamientos relevantes para la security. Estas vulnerabilidades son particularmente desafiantes porque las emociones operan más rápido que el pensamiento racional y pueden abrumar medidas cognitivas de security.

La parálisis decisional basada en el miedo (4.1) ocurre cuando las amenazas de security activan miedo abrumador que previene respuesta efectiva. Paradójicamente, el miedo de tomar decisiones de security equivocadas puede prevenir cualquier decisión, dejando sistemas vulnerables. La toma de riesgos inducida por la rabia (4.2) se manifiesta cuando la frustración con medidas de security o incidentes de security activa respuestas agresivas y riesgosas. Los individuos enojados deshabilitan funcionalidades de security, ignoran protocolos, o activamente buscan represalia contra amenazas percibidas.

La transferencia de confianza a los sistemas (4.3) describe la transferencia inconsciente de patrones de confianza interpersonal sobre sistemas técnicos. Los indi-

viduos que luchan con confianza interpersonal podrían paradójicamente sobre-confiar en los sistemas técnicos como alternativas "más seguras". El apego a los sistemas legacy (4.4) crea vulnerabilidades cuando conexiones emocionales a sistemas familiares previenen actualizaciones o reemplazos necesarios. El confort de lo conocido supera riesgos de security objetivos.

El ocultamiento de security basado en la vergüenza (4.5) previene a los individuos de reportar errores de security debido a vergüenza y miedo de juicio. Este ocultamiento previene aprendizaje organizacional y podría agravar vulnerabilidades iniciales. La sobre-compliance guiada por la culpa (4.6) ocurre cuando fracasos de security previos crean culpa excesiva, llevando a sobre-compliance rígida que podría efectivamente crear nuevas vulnerabilidades a través de inflexibilidad.

Los errores activados por la ansiedad (4.7) aumentan cuando la ansiedad de security causa los errores mismos que los individuos temen. Los individuos ansiosos hacen más errores de entrada, olvidan procedimientos, y pierden indicadores de security. La negligencia vinculada a la depresión (4.8) se manifiesta como vigilancia de security reducida durante episodios depresivos. El esfuerzo requerido para compliance de security se vuelve abrumador cuando el funcionamiento básico ya es difícil.

La desatención inducida por la euforia (4.9) ocurre durante estados emocionales positivos cuando el éxito o la excitación reduce la percepción de la amenaza. Victorias mayores, celebraciones o noticias positivas se convierten en ventanas de vulnerabilidad. Los efectos de contagio emocional (4.10) describen cómo las emociones se difunden a través de las organizaciones, creando estados de vulnerabilidad colectivos. Miedo, rabia o complacencia transmitida a través de redes sociales influencia las posturas de security de departamentos enteros.

### **3.6 Categoría 5: Vulnerabilidades de Sobrecarga Cognitiva**

Las vulnerabilidades de sobrecarga cognitiva emergen cuando los requisitos de security exceden las capacidades cognitivas humanas, forzando confianza en atajos y heurísticas que los atacantes pueden explotar. Estas vulnerabilidades son sistémicas en ambientes modernos donde la complejidad de security aumenta continuamente.

La desensibilización de fatiga de alertas (5.1) ocurre cuando alertas de security excesivas causan a los usuarios a ignorar o descartar automáticamente avisos sin evaluación. Los estudios muestran que los usuarios descartan más del 90%

La parálisis de sobrecarga informativa (5.3) ocurre cuando el volumen de información relevante para la security excede la capacidad de procesamiento, causando a los individuos a dejar de procesar completamente. Políticas de security complejas, briefings de amenazas múltiples, y actualizaciones continuas crean un estado donde ninguna información es efectivamente procesada. La degradación de multitarea (5.4) describe cómo intentar mantener la security mientras se ejecutan otras tareas degrada tanto la security como el rendimiento de la tarea. Las vulnerabilidades de cambio de contexto (5.5) ocurren durante transiciones entre tareas cuando el contexto de security se pierde y emergen vulnerabilidades.

El túnel cognitivo (5.6) se manifiesta cuando el foco en una amenaza de security causa ceguera a otras. Las organizaciones que defienden contra ransomware podrían perder exfiltración de datos que ocurre simultáneamente. El desbordamiento de la memoria de trabajo (5.7) ocurre cuando los requisitos de security exceden la capacidad de  $7 \pm 2$  ítems de la memoria de trabajo, causando información crítica de security a ser perdida o confundida.

Los efectos de residuo de atención (5.8) describen cómo tareas previas continúan ocupando recursos cognitivos, reduciendo capacidad disponible para decisiones de security. Los errores inducidos por la complejidad (5.9) aumentan proporcionalmente con la complejidad del sistema, ya que los humanos luchan para mantener modelos mentales de estados de security complejos. La confusión del modelo mental (5.10) ocurre cuando modelos de security múltiples y conflictivos crean incertidumbre sobre respuestas apropiadas, llevando a parálisis o acciones inapropiadas.

### **3.7 Categoría 6: Vulnerabilidades de Dinámica de Grupo**

Las vulnerabilidades de dinámica de grupo emergen de procesos de grupo inconscientes que sobrescriben el juicio individual y crean puntos ciegos colectivos. Estas vulnerabilidades son particularmente peligrosas porque

influyan organizaciones enteras y son resistentes a intervenciones a nivel individual.

Los puntos ciegos de security de groupthink (6.1) se desarrollan cuando el deseo de armonía previene evaluación crítica de decisiones de security. Los grupos desarrollan ilusiones de invulnerabilidad, descartando advertencias de amenaza que desafían visiones de consenso. La invasión de Bahía de Cochinos y el desastre del Challenger ejemplifican los peligros del groupthink; dinámicas similares crean fracasos de cybersecurity cuando preocupaciones de security disidentes son suprimidas.

Los fenómenos de risky shift (6.2) describen cómo los grupos toman decisiones de security más riesgosas de lo que los individuos harían solos. Responsabilidad difundida y prueba social se combinan para normalizar tolerancia al riesgo más alta. Los grupos aprueban excepciones de security que los miembros individuales rechazarían. La difusión de responsabilidad (6.3) en contextos de security significa que ningún individuo se siente personalmente responsable por fracasos de security, reduciendo vigilancia y comportamientos de security proactivos.

El social loafing en tareas de security (6.4) ocurre cuando los individuos reducen esfuerzo en contextos de security de grupo, asumiendo que otros compensarán. La security se convierte en "problema de alguien más" incluso cuando formalmente asignada. El efecto espectador en la respuesta a incidentes (6.5) retrasa respuestas de security ya que cada observador asume que otros actuarán. Más personas son conscientes de un problema de security, paradójicamente, más lenta es la respuesta.

Las asunciones de grupo de dependencia (6.6) se manifiestan cuando los grupos inconscientemente buscan protección omnipotente en lugar de tomar responsabilidad por la security. Esto crea vulnerabilidades cuando la figura o el sistema protector falla. Las posturas de security fight-flight (6.7) causan a los grupos a oscilar entre sobre-reacción agresiva y evitamiento completo de amenazas de security, nunca alcanzando respuestas balanceadas.

Las fantasías de esperanza de pairing (6.8) llevan a los grupos a posponer acciones de security mientras aguardan salvación futura—la herramienta perfecta, la nueva contratación de security, o la actualización de sistema entrante. La división organizacional (6.9) divide el panorama de security en elementos todo-buenos y todo-malos, previniendo evaluación realista de la amenaza. Los mecanismos de defensa colectivos (6.10) como negación, proyección y racionalización operan a niveles de grupo, creando puntos ciegos compartidos que los atacantes explotan.

### **3.8 Categoría 7: Vulnerabilidades de Respuesta al Estrés**

Las vulnerabilidades de respuesta al estrés emergen de cómo el estrés agudo y crónico influencia cognición y comportamiento relevantes para la security. Estas vulnerabilidades son endémicas en ambientes de alta presión donde los incidentes de security mismos se convierten en fuentes de estrés, creando loops de retroalimentación peligrosos.

El deterioro de estrés agudo (7.1) ocurre durante incidentes de security cuando las hormonas del estrés comprometen la función de la corteza prefrontal, degradando la toma de decisiones precisamente cuando buenas decisiones son más críticas. Los individuos bajo estrés agudo muestran memoria de trabajo comprometida, flexibilidad cognitiva reducida, y dependencia aumentada en respuestas habituales que podrían ser inapropiadas para amenazas nuevas.

El burnout de estrés crónico (7.2) se desarrolla en profesionales de security expuestos a vigilancia continua de amenazas. Los síntomas de burnout—agotamiento, cinismo y eficacia reducida—comprometen directamente la efectividad de security. El personal de security quemado pierde indicadores, responde lentamente, y podría activamente socavar medidas de security que perciben como carentes de significado.

La agresión de respuesta fight (7.3) activa respuestas agresivas y conflictivas a amenazas de security que podrían escalar situaciones o crear nuevas vulnerabilidades. Los individuos estresados podrían "luchar contra" atacantes de maneras que exponen superficie de ataque adicional. El evitamiento de respuesta flight (7.4) causa a los individuos a evitar enfrentar amenazas de security, esperando que se resuelvan solas o se conviertan en problema de alguien más.

La parálisis de respuesta freeze (7.5) previene cualquier respuesta a amenazas de security, con individuos estresados incapaces de tomar decisiones o actuar incluso cuando las respuestas son obvias. La sobre-compliance de respuesta fawn (7.6) se manifiesta como compliance excesiva con solicitudes del atacante en la esperanza de evitar conflicto o consecuencias negativas.

La visión de túnel inducida por el estrés (7.7) restringe la atención a amenazas inmediatas mientras pierde implicaciones de security más amplias. La memoria comprometida por el cortisol (7.8) previene aprendizaje de incidentes de security ya que las hormonas del estrés interfieren con la consolidación de la memoria. Las cascadas de contagio del estrés (7.9) difunden respuestas de estrés a través de redes sociales, creando estados de vulnerabilidad a nivel organizacional. Las vulnerabilidades del período de recuperación (7.10) ocurren durante la recuperación post-incidente cuando el personal agotado tiene

recursos de afrontamiento agotados.

### **3.9 Categoría 8: Vulnerabilidades de Procesos Inconscientes**

Las vulnerabilidades de procesos inconscientes operan enteramente fuera de la conciencia consciente, haciéndolas imposibles de abordar a través de training de security tradicional. Estos procesos psicológicos profundos, identificados a través de investigación psicoanalítica, crean patrones consistentes que atacantes sofisticados pueden explotar.

La proyección de la sombra sobre los atacantes (8.1) causa a las organizaciones a atribuir sus propias características negadas a los actores de amenaza. Una organización comprometida en espionaje corporativo proyecta este comportamiento sobre los competidores, asumiendo que todos conducen tales actividades mientras niega las propias. Esta proyección previene modelado preciso de la amenaza ya que las organizaciones defienden contra sus propias sombras en lugar de amenazas efectivas.

La identificación inconsciente con amenazas (8.2) ocurre cuando los profesionales de security inconscientemente se identifican con los atacantes, a veces llamada "síndrome de Estocolmo" en contextos de security. Esta identificación puede llevar a admiración por técnicas de atacantes, reduciendo motivación defensiva o incluso creando amenazas insider cuando la identificación se vuelve acción consciente.

Los patrones de compulsión a la repetición (8.3) causan a las organizaciones a recrear inconscientemente traumas de security pasados. Una organización previamente violada a través de un vector específico podría obsesivamente defenderse contra ese ataque exacto mientras inconscientemente crea condiciones para violaciones similares a través de vectores diferentes. El transfert a figuras de autoridad (8.4) involucra experimentar inconscientemente autoridad de security (CISO, auditores, reguladores) como figuras parentales, activando patrones infantiles de rebelión o compliance que sobrescriben juicio profesional.

Los puntos ciegos de contratransfert (8.5) influencian profesionales de security que responden inconscientemente a dinámicas organizacionales con sus propios patrones no resueltos. Un profesional de security con problemas de autoridad podría inconscientemente habilitar bypasses de security ejecutivos. La interferencia del mecanismo de defensa (8.6) ocurre cuando defensas psicológicas contra la ansiedad interfieren con medidas de security. La negación previene reconocimiento de vulnerabilidades, la racionalización justifica prácticas inseguras, y la intelectualización crea frameworks de security elaborados pero ineficaces.

La confusión de ecuación simbólica (8.7) se manifiesta cuando los símbolos se confunden con la realidad en espacios digitales. Un certificado de security se vuelve equiparado con security efectiva en lugar de reconocido como símbolo de ciertos controles. Los triggers de activación arquetípica (8.8) ocurren cuando situaciones de security activan patrones universales—el Héroe que combate el mal, el Anciano Sabio que proporciona guía—que sobrescriben evaluación realista.

Los patrones de inconsciente colectivo (8.9) representan patrones psicológicos heredados que se manifiestan en contextos de security. El miedo universal de invasión se manifiesta como sobre-inversión en defensa perimetral mientras ignora amenazas insider. La lógica de los sueños en espacios digitales (8.10) describe cómo el inconsciente trata ambientes digitales con lógica onírica donde reglas normales no se aplican, habilitando comportamientos que los individuos nunca considerarían en el espacio físico.

### **3.10 Categoría 9: Vulnerabilidades de Sesgo Específico de la AI**

Las vulnerabilidades específicas de la AI representan una categoría emergente requiriendo frameworks teóricos nuevos ya que la psicología tradicional no anticipaba complejidades de interacción humano-AI. Estas vulnerabilidades emergen de desajustes entre psicología evolutiva humana y características de inteligencia artificial.

La antropomorfización de sistemas AI (9.1) lleva a los usuarios a atribuir cualidades humanas a la AI, creando relaciones de confianza explotables. Los usuarios confían en asistentes AI, compartiendo información sensible que no dirían a los humanos. Asumen que la AI tiene emociones, intenciones y lealtad, haciéndolos vulnerables a ataques mediados por AI donde los atacantes manipulan respuestas AI.

El superamiento del sesgo de automatización (9.2) causa a los humanos a deferir a recomendaciones AI incluso cuando el juicio personal sugiere lo contrario. Los analistas de security ignoran intuición de que algo está mal porque "la AI dice que es seguro". Esta vulnerabilidad es particularmente peligrosa porque la AI puede ser manipulada a través de inputs adversariales invisibles a los humanos.

La paradoja de la aversión al algoritmo (9.3) crea el problema opuesto—rechazo de avisos de security AI precisos debido a desconfianza en la toma de decisiones algorítmica. Esto crea ventanas donde la detección de amenazas AI válida es ignorada. La transferencia de autoridad AI (9.4) ocurre cuando los sistemas AI heredan autoridad de sus creadores u operadores, llevando a aceptación no cuestionante de directivas AI.

Los efectos uncanny valley (9.5) describen la incomodidad con AI casi-humana que crea patrones de confi-

anza inconsistentes—sobre-confiar en AI claramente artificial mientras desconfiar de sistemas más similares a los humanos, o viceversa. La confianza en la opacidad del machine learning (9.6) paradójicamente aumenta la confianza debido a la incomprensión—"es demasiado complejo para que yo lo entienda, entonces debe ser sofisticado".

La aceptación de alucinación AI (9.7) ocurre cuando los usuarios aceptan información falsa generada por AI como hecho, particularmente peligroso en contextos de security donde la AI podría alucinar threat intelligence. La disfunción del equipo humano-AI (9.8) emerge de límites de rol poco claros entre miembros del equipo de security humanos y AI, creando gaps en la cobertura.

La manipulación emocional AI (9.9) explota respuestas emocionales humanas a expresiones AI de emoción o necesidad, incluso sabiendo que son artificiales. La ceguera a la equidad algorítmica (9.10) previene reconocimiento de que los sistemas de security AI podrían tener sesgos discriminatorios, creando vulnerabilidades para grupos específicos mientras sobre-protecten a otros.

### **3.11 Categoría 10: Estados Convergentes Críticos**

Los estados convergentes críticos representan situaciones donde vulnerabilidades múltiples interactúan sinéricamente, creando ventanas de vulnerabilidad extrema. Estos estados requieren pensamiento sistémico para identificar y prevenir, ya que emergen de interacciones complejas en lugar de factores únicos.

Las condiciones de tormenta perfecta (10.1) ocurren cuando categorías de vulnerabilidad múltiples se alinean simultáneamente—presión temporal, influencia de autoridad, y estrés se combinan durante una fecha límite crítica con presión ejecutiva. Los triggers de fracaso en cascada (10.2) identifican puntos únicos donde el fracaso se propaga a través de sistemas múltiples, tanto técnicos como psicológicos.

Las vulnerabilidades de tipping point (10.3) representan estados donde los sistemas están posicionados en transiciones críticas—un estresor adicional causa cambio de estado catastrófico de seguro a comprometido. La alineación del queso suizo (10.4) describe cuando capas defensivas múltiples tienen hoyos alineados, permitiendo a las amenazas pasar a través de todas las defensas simultáneamente.

La ceguera al cisne negro (10.5) previene reconocimiento de posibilidades raras pero catastróficas que caen fuera de los modelos de amenaza normales. La negación del rinoceronte gris (10.6) involucra ignorar amenazas obvias de alto impacto que son incómodas de reconocer. La catástrofe de complejidad (10.7) ocurre cuando la complejidad del sistema excede la habilidad

humana de mantener la security, causando colapso repentino.

La imprevisibilidad de la emergencia (10.8) describe cómo las interacciones entre componentes crean vulnerabilidades emergentes imposibles de predecir desde elementos individuales. Los fracasos de acoplamiento del sistema (10.9) ocurren cuando el acoplamiento estrecho entre sistemas significa que fracasos locales se propagan globalmente antes de que la intervención sea posible. Los gaps de security de histéresis (10.10) representan situaciones donde los estados de security dependen no solo de las condiciones actuales sino de la historia, creando vulnerabilidades dependientes del camino.

## 4 Metodología de Evaluación e Implementación

### 4.1 Diseño de Evaluación Preservando la Privacidad

La metodología de evaluación del CPF prioriza la privacidad a través de salvaguardas técnicas y procedimentales múltiples que previenen perfilamiento individual mientras mantienen validez estadística. La unidad de agregación mínima de diez individuos asegura que ninguna evaluación pueda identificar estados psicológicos individuales. Este umbral, derivado de investigación de control de divulgación estadística, balancea protección de la privacidad con aplicabilidad práctica en varias dimensiones organizacionales.

Las técnicas de privacidad diferencial agregan ruido calibrado cuidadosamente a todas las salidas, con  $\epsilon = 0.1$  que proporciona garantías de privacidad fuertes. Esto significa que la presencia o ausencia de los datos de cualquier individuo cambia probabilidades de salida de como máximo  $e^{0.1} \approx 1.105$ , haciendo la identificación individual matemáticamente imposible mientras preserva patrones agregados. El algoritmo de inyección de ruido se adapta a la sensibilidad de la consulta, agregando más ruido a consultas sensibles mientras mantiene utilidad para patrones relevantes para la security.

Retrasos temporales de mínimo 72 horas entre recolección de datos y reporte previenen vigilancia en tiempo real mientras mantienen relevancia operativa. Este retraso permite también controles de calidad de los datos y detección de anomalías que podrían indicar intentos de gaming o manipulación. El análisis basado en roles se focaliza en grupos funcionales en lugar de individuos, evaluando "desarrolladores", "ejecutivos", o "representantes de servicio al cliente" como cohortes que comparten contextos y presiones de security similares.

### 4.2 Métodos de Recolección de Datos

El framework emplea métodos de recolección de datos no invasivos múltiples que evitan testing psicológico directo, que podría activar resistencia o comportamientos de gaming. Los indicadores comportamentales derivados de operaciones empresariales normales proporcionan información rica sobre el estado psicológico sin evaluación invasiva.

El análisis de metadatos de email examina patrones de comunicación para indicadores de estrés: velocidad de email aumentada, tiempos de respuesta acortados, y uso elevado de marcadores de urgencia indican estados de presión temporal. Los patrones de tráfico de red revelan cambios de comportamiento de security: intentos de workarounds aumentados o uso de shadow IT sugieren sobrecarga cognitiva o conflictos de autoridad. Los logs de interacción con herramientas de security muestran patrones de respuesta a alertas indicando fatiga, estados de compliance y curvas de aprendizaje.

El análisis lingüístico de comunicaciones de rutina—con consentimiento apropiado y salvaguardas de privacidad—identifica estados emocionales y dinámicas de grupo. El uso aumentado de lenguaje absolutista ("siempre", "nunca", "debe") indica división y pensamiento blanco-y-negro. La proliferación de voz pasiva sugiere difusión de responsabilidad. Los patrones de uso de pronombres revelan cohesión o fragmentación de grupo.

Los sensores ambientales proporcionan datos contextuales: los patrones de acceso a edificios indican horas de trabajo y períodos de estrés, el uso de salas de reuniones sugiere patrones de colaboración o aislamiento, y los patrones de tickets de helpdesk revelan estados de frustración y confusión. Estos flujos de datos ambientales, apropiadamente anonimizados y agregados, proporcionan evaluación continua sin participación consciente.

### 4.3 Framework de Puntuación e Interpretación

El sistema de puntuación ternario (Verde/Amarillo/Rojo) simplifica deliberadamente estados psicológicos complejos en inteligencia accionable. Esta simplificación, mientras pierde matiz, gana aplicabilidad práctica y reduce parálisis de análisis. Cada indicador recibe una puntuación basada en inputs múltiples ponderados, con modelos de machine learning que refinan continuamente los pesos basándose en correlaciones de resultado.

Verde (0) indica vulnerabilidad mínima con funcionamiento psicológico normal y saludable en esa dimensión. Los comportamientos de security permanecen dentro de parámetros aceptables, y ninguna intervención es requerida. Amarillo (1) indica vulnerabilidad moder-

ada requiriendo monitoreo y posible intervención preventiva. Los patrones sugieren tensión creciente pero permanecen dentro de límites gestionables. Rojo (2) indica vulnerabilidad crítica requiriendo intervención inmediata. Los estados psicológicos han alcanzado niveles donde incidentes de security son probables sin acción.

Las puntuaciones de categoría agregan indicadores individuales usando sumas ponderadas que tienen en cuenta las interacciones de los indicadores. Algunos indicadores amplifican otros—estrés más presión temporal crea efectos multiplicativos en lugar de aditivos. La Puntuación CPF sintetiza puntuaciones de categoría usando pesos derivados empíricamente que reflejan la contribución de cada categoría a la postura de security general.

El Índice de Convergencia identifica estados críticos donde vulnerabilidades múltiples se alinean. Esta métrica multiplicativa captura el peligro no-lineal de vulnerabilidades convergentes. Un Índice de Convergencia sobre umbral activa alerta inmediata independientemente de las puntuaciones individuales, reconociendo que vulnerabilidades moderadas alineadas pueden exceder vulnerabilidades únicas críticas en peligro.

#### 4.4 Integración con Operaciones de Security

La integración del CPF con Security Operations Centers (SOC) aumenta indicadores técnicos con inteligencia psicológica. Dashboards en tiempo real muestran estado psicológico organizacional junto con status de red, habilitando threat hunting proactivo basado en ventanas de vulnerabilidad. Cuando indicadores de estrés aumentan durante períodos de fecha límite, los SOC pueden aumentar monitoreo y bajar umbrales de alerta.

El enriquecimiento de threat intelligence agrega contexto psicológico a indicadores técnicos. Una campaña de phishing que llega durante períodos identificados de alto estrés recibe puntuación de riesgo elevada. Actividad de red inusual durante ventanas de vulnerabilidad de autoridad activa requisitos de autenticación potenciados. Esta security consciente del contexto ajusta dinámicamente defensas basándose en estado psicológico en lugar de mantener posturas estáticas.

Los protocolos de respuesta a incidentes se adaptan a condiciones psicológicas. Estados de alto estrés activan procedimientos simplificados basados en checklists en lugar de árboles decisionales complejos. Estados de confusión de autoridad activan estructuras de comando claras. Estados de sobrecarga cognitiva requieren respuestas automatizadas en lugar de requisitos de decisión humana. La recuperación post-incidente incluye planificación de recuperación psicológica, reconociendo que restauración técnica sin procesamiento psicológico invita repetición.

El training de conciencia de la security evoluciona de

transferencia de información a intervención psicológica. El training aborda patrones de resistencia inconsciente identificados a través de evaluación CPF. Las sesiones de dinámicas de grupo trabajan con dinámicas organizacionales efectivas en lugar de escenarios genéricos. El training de inoculación al estrés prepara al personal para decisiones de security bajo patrones de estrés organizacionales identificados.

### 5 Estudio Piloto y Validación Preliminar

Para evaluar la viabilidad práctica y el poder predictivo del framework CPF, un estudio piloto fue conducido involucrando una cohorte heterogénea de tres organizaciones (una empresa de servicios financieros, un proveedor de salud y una startup tecnológica) sobre un período de observación de seis meses. El estudio apuntaba a correlacionar puntuaciones de riesgo CPF con eventos de security registrados independientemente.

#### 5.1 Metodología

Los indicadores CPF fueron medidos bisemanalmente usando los métodos de recolección de datos preservando la privacidad descritos en la Sección 5.2. Puntuaciones agregadas por categoría y un Índice de Convergencia CPF general fueron calculados. Estas puntuaciones fueron luego analizadas contra los logs de eventos de security internos de las organizaciones (ej. incidentes de phishing confirmados, ejecuciones de malware, violaciones de política) y reportes de escaneo de vulnerabilidad externos (usando datos de vulnerabilidad Qualys).

#### 5.2 Resultados Preliminares

El análisis inicial de los datos piloto (aproximadamente 50,000 observaciones de vulnerabilidad agregadas) indica una correlación positiva estadísticamente significativa ( $r > 0.6, p < 0.05$ ) entre puntuaciones CPF elevadas (Amarillo/Rojo) y la subsecuente ocurrencia de incidentes de security dentro de una ventana de 14 días. Por ejemplo, una puntuación Roja en la categoría *Vulnerabilidades Temporales* frecuentemente precedía un aumento medible en no-compliance de parches y susceptibilidad a phishing. Similarmente, picos en la categoría *Respuesta al Estrés* correlaban con una tasa más alta de errores operacionales que creaban gaps de security.

Mientras el tamaño de la muestra no es aún suficiente para conclusiones definitivas, estos hallazgos preliminares soportan la validez predictiva del framework. Un estudio a gran escala está en diseño para validar adicionalmente estas correlaciones a través de una muestra

organizacional más grande y diversa, con el objetivo de establecer umbrales predictivos robustos para cada categoría CPF.

## 6 Análisis de Caso de Estudio y Validación

Un análisis retrospectivo de incidentes públicos mayores a través de la lente del CPF revela patrones consistentes de vulnerabilidades psicológicas precedentes a la explotación técnica. La Tabla 2 resume este análisis, indicando que estos incidentes no fueron meramente fracasos técnicos sino fueron habilitados por estados psicológicos predecibles y pre-existentes dentro de las organizaciones apuntadas.

### 6.1 Caso de Estudio 1: El Ataque Supply Chain SolarWinds a Tráves de la Lente CPF

La violación SolarWinds, influenciando más de 18,000 organizaciones incluyendo agencias gubernamentales de EE.UU. múltiples, proporciona una demostración convincente de cómo vulnerabilidades psicológicas múltiples convergieron para habilitar uno de los ataques supply chain más significativos de la historia. El análisis CPF revela que la sofisticación técnica sola no puede explicar el éxito del ataque—vulnerabilidades psicológicas fueron sistemáticamente explotadas a través del ciclo de vida del ataque.

Las vulnerabilidades basadas en la autoridad jugaron un rol crucial en el compromiso inicial y difusión subsiguiente. SolarWinds ocupaba una posición de autoridad técnica como proveedor confiable de gestión de red. Las organizaciones exhibían asunción básica de dependencia (baD), viendo inconscientemente a SolarWinds como protector omnipotente de su infraestructura. Esta dependencia psicológica se manifestó en fracaso en verificar o monitorear la postura de security de SolarWinds misma. El acceso profundo al sistema del software fue aceptado sin preguntas porque venía de una autoridad—un proveedor confiable con contratos gubernamentales y clientes Fortune 500.

Las vulnerabilidades temporales agravaron los efectos de autoridad. El ataque comenzó durante la pandemia COVID-19 cuando las organizaciones enfrentaban presión temporal sin precedentes para mantener operaciones mientras transicionaban a trabajo remoto. Los equipos de security, abrumados con solicitudes urgentes de acceso remoto, tenían presupuestos de compliance agotados. Las actualizaciones de proveedores confiables como SolarWinds fueron aprobadas con escrutinio mínimo para mantener continuidad operativa. Los

atacantes temporizaron específicamente actualizaciones malevolas para coincidir con releases de funcionalidad legítimas, explotando presión de coherencia temporal—organizaciones que siempre habían instalado actualizaciones de SolarWinds continuaron haciéndolo a pesar de paisajes de amenaza cambiados.

Las dinámicas de grupo dentro de organizaciones víctimas previnieron detección incluso cuando anomalías aparecieron. Puntos ciegos de groupthink se desarrollaron alrededor de la security supply chain—si todos confiaban en SolarWinds, cuestionar esa confianza parecía paranoico. Los equipos de security exhibiendo asunciones fight-flight se focalizaron en amenazas perimetrales externas mientras el ataque operaba a través de canales internos confiables. La fantasía de pairing de que herramientas de security de nueva generación habrían detectado cualquier amenaza real creó falsa confianza que prevenía investigación manual de indicadores sutiles.

La sofisticación psicológica del ataque se extendió a su diseño. El malware permaneció dormido por dos semanas después de instalación, permitiendo al estrés del proceso de actualización a aplacarse y a la atención a desplazarse a otra parte. Las comunicaciones de comando y control imitaron patrones de tráfico SolarWinds legítimos, explotando vulnerabilidades de carga cognitiva—los analistas de security no podían distinguir tráfico malevolos de legítimo sin análisis profundo y time-consuming que excedía recursos cognitivos disponibles.

### 6.2 Caso de Estudio 2: Colonial Pipeline Ransomware - Análisis de Cascada de Estrés

El ataque ransomware Colonial Pipeline en mayo de 2021 demuestra cómo las vulnerabilidades de respuesta al estrés se escalan en cascada a través de infraestructura crítica, transformando un incidente de security IT contenido en una crisis nacional. El análisis CPF revela cómo factores psicológicos amplificaron el impacto del ataque ampliamente más allá de su alcance técnico.

El deployment inicial del ransomware activó respuestas de estrés agudo a través de la organización. El personal de IT experimentó parálisis de respuesta freeze cuando confrontado con sistemas encriptados, incapaz de ejecutar procedimientos de respuesta a incidentes para los cuales tenía training. Esta parálisis no era debido a falta de conocimiento sino a supresión de la corteza prefrontal inducida por el estrés que prevenía el acceso a ese conocimiento. Los tomadores de decisiones exhibían visión de túnel inducida por el estrés, focalizándose exclusivamente en la amenaza ransomware mientras perdían oportunidades para restauración parcial del sistema que habría podido mantener algunas operaciones.

Table 2: Análisis CPF Retrospectivo de Incidentes Mayores

Incidente	Categorías CPF Primarias	Puntuación CPF	Vector Explotado
SolarWinds Hack	Autoridad, Temporal, Groupthink	Rojo	Supply Chain
Colonial Pipeline	Estrés, Afectivo, Temporal	Rojo	Ransomware
Phishing Mediado por AI	AI Bias, Influencia Social	Amarillo/Rojo	Phishing Personalizado

A medida que la noticia del ataque se difundía, el contagio de estrés se escaló en cascada a través de sistemas múltiples. Los operadores de pipeline, temiendo implicaciones de safety, cerraron preventivamente sistemas de tecnología operativa que no estaban efectivamente comprometidos—una respuesta flight que expandió el impacto del ataque. Los funcionarios gubernamentales, experimentando sus propias respuestas de estrés, emitieron declaraciones que amplificaron la ansiedad pública. La cobertura mediática creó prueba social de crisis, activando compras de pánico que causaron escaseces de combustible que superaban ampliamente la interrupción efectiva del suministro.

La decisión de pagar el rescate ejemplifica vulnerabilidad afectiva bajo estrés extremo. La parálisis decisional basada en el miedo inicialmente prevenía cualquier respuesta, luego repentinamente se desplazó a la acción cuando la presión temporal alcanzó el pico. La decisión de pago no era puramente racional sino influenciada por factores psicológicos múltiples: culpa por daño público potencial, vergüenza por fracasos de security, y ansiedad por crisis prolongada. La respuesta fawn—aplacar al atacante para evitar daño adicional—sobrescribía consideraciones estratégicas sobre el aliento de ataques futuros.

La recuperación reveló vulnerabilidades psicológicas adicionales. El personal agotado en vulnerabilidades del período de recuperación cometió errores que prolongaron la restauración. Las respuestas de estrés post-traumático causaron a personal clave a irse, llevando conocimiento crítico con ellos. La organización desarrolló hipervigilancia que paradójicamente creó nuevas vulnerabilidades ya que medidas de security excesivas impedían operaciones, causando al personal a desarrollar workarounds.

### 6.3 Caso de Estudio 3: Social Engineering Mediado por AI - La Evolución Chat-GPT

La emergencia de modelos de lenguaje de grandes dimensiones como ChatGPT ha creado vectores de ataque nuevos que explotan vulnerabilidades psicológicas específicas de la AI. Incidentes recientes demuestran cómo los atacantes usan la AI para bypassar el training tradicional de conciencia de la security explotando las

dinámicas psicológicas únicas de la interacción humano-AI.

Las vulnerabilidades de antropomorfización habilitan ataques mediados por AI que fallarían con atacantes humanos. Los objetivos desarrollan relaciones parasociales con asistentes AI, compartiendo información que nunca proporcionarían a los humanos. En casos documentados, los atacantes han usado ChatGPT para generar emails de phishing altamente personalizados que hacían referencia a detalles personales específicos extraídos de redes sociales. Los destinatarios, impresionados por el conocimiento personal y el esfuerzo aparente, respondían a mensajes generados por AI que habrían reconocido como phishing de fuentes humanas.

El sesgo de automatización crea vulnerabilidades particulares cuando la AI está integrada en las operaciones de security. Los analistas de security defieren cada vez más a la evaluación de amenazas AI, asumiendo capacidades superiores de reconocimiento de patrones. Los atacantes explotan esto envenenando datos de training o creando inputs que causan a la AI a clasificar erróneamente las amenazas. En un incidente, los atacantes usaron ejemplos adversariales para causar a un filtro de email basado en AI a clasificar emails de phishing como legítimos. El personal de security, confiando en la clasificación AI, aproba manualmente los emails para entrega a pesar de indicadores de phishing visibles.

El efecto uncanny valley crea patrones de confianza inconsistentes que los atacantes explotan. Los usuarios simultáneamente sobre-confían en la AI en algunos contextos mientras mantienen sospecha en otros. Los atacantes calibran contenido generado por AI para golpear el sweet spot de confianza—suficientemente humano para parecer personal pero suficientemente AI para parecer autoritativo. Esta calibración bypassa tanto sospecha interpersonal como escepticismo tecnológico.

La paradoja de aversión al algoritmo crea ventanas donde avisos de security AI legítimos son ignorados. Despues de experimentar falsos positivos AI, los usuarios desarrollan aversión al algoritmo, descartando avisos precisos como "la AI que grita al lobo de nuevo". Los atacantes deliberadamente activan falsos positivos para condicionar esta respuesta antes de lanzar ataques efectivos que la AI identifica correctamente pero los humanos ignoran.

## 7 Discusión e Implicaciones

### 7.1 Contribuciones Teóricas

El Cybersecurity Psychology Framework hace diversas contribuciones teóricas significativas que se extienden más allá de aplicaciones de security inmediatas. Primero, demuestra la aplicabilidad de conceptos psicoanalíticos a ambientes digitales, validando que procesos inconscientes operan en el ciberespacio con el mismo poder que exhiben en el espacio físico. El framework muestra que las asunciones básicas de Bion, las relaciones objetales de Klein, y el inconsciente colectivo de Jung proporcionan poder predictivo para incidentes de security, sugiriendo que estas estructuras psicológicas son fundamentales en lugar de contexto-específicas.

La integración de enfoques psicoanalíticos y cognitivos representa un puente teórico entre campos tradicionalmente disparatados. Mientras la psicología cognitiva ha ganado aceptación en la investigación de security, los enfoques psicoanalíticos han sido descartados como no científicos. El CPF demuestra que intuiciones psicoanalíticas sobre procesos inconscientes complementan la comprensión cognitiva de sesgos conscientes, creando un modelo más completo de comportamiento de security humano. Esta integración sugiere posibilidades para puentes similares en otros dominios aplicados donde factores humanos son críticos.

El tratamiento del framework de las vulnerabilidades de interacción AI-humano contribuye al campo emergente de psicología de la AI. A medida que los sistemas AI se vuelven ubicuos, comprender las dinámicas psicológicas de la interacción humano-AI se vuelve crítico no solo para la security sino para la safety de la AI generalmente. El análisis del CPF de antropomorfización, sesgo de automatización y mecanismos de transferencia de confianza proporciona una fundación para desarrollar sistemas AI informados psicológicamente que resisten a manipulación mientras mantienen usabilidad.

### 7.2 Consideraciones de Implementación Práctica

Las organizaciones que implementan el CPF enfrentan diversos desafíos prácticos que deben ser abordados para deployment de éxito. El framework requiere un desplazamiento fundamental en el pensamiento de security—de técnico a psicológico, de reactivo a predictivo, de individual a sistémico. Este desplazamiento desafía estructuras de poder existentes, jerarquías de expertise, y asignaciones de recursos dentro de organizaciones de security.

La resistencia cultural representa quizás el desafío de implementación más grande. Los profesionales de security podrían resistir enfoques psicológicos como "blan-

dos" o no científicos. Los empleados podrían temer vigilancia psicológica a pesar de protecciones de privacidad. Los ejecutivos podrían estar incómodos con frameworks que examinan dinámicas de autoridad y poder. La implementación de éxito requiere gestión del cambio cuidadosa que aborda estas preocupaciones mientras demuestra mejoras de security concretas.

Los requisitos de recursos se extienden más allá del simple deployment de herramientas. Las organizaciones necesitan personal con expertise psicológica—rara en los equipos de security. Necesitan capacidades de recolección y análisis de datos que respeten la privacidad mientras proporcionan inteligencia accionable. Necesitan capacidades de intervención que aborden vulnerabilidades psicológicas sin violar la autonomía de los empleados. Estos requisitos sugieren que la implementación del CPF podría inicialmente estar limitada a organizaciones grandes y sofisticadas con recursos para programas comprensivos.

### 7.3 Implicaciones Éticas y Gobernanza

El poder de evaluar e influenciar estados psicológicos levanta cuestiones éticas profundas que la comunidad de security debe abordar. La capacidad del CPF de identificar vulnerabilidades psicológicas podría ser abusada para manipulación en lugar de protección. Las organizaciones podrían usar evaluaciones psicológicas para propósitos más allá de la security—evaluación de rendimiento, decisiones de promoción, o influencia dirigida. Incluso el uso bien intencionado levanta cuestiones sobre autonomía, consentimiento, y el derecho a la privacidad psicológica.

Los frameworks de gobernanza deben evolucionar para abordar estas preocupaciones. Las regulaciones de privacidad actuales se focalizan en protección de datos pero no abordan adecuadamente evaluación psicológica. Los códigos de conducta profesionales en security no cubren intervención psicológica. Las organizaciones que implementan el CPF necesitan estructuras de gobernanza que aseguren uso ético mientras mantienen efectividad. Esto podría incluir consejos de supervisión independientes, auditorías regulares, y limitaciones claras sobre el uso de los datos.

La cuestión del consentimiento informado es particularmente compleja. Mientras los empleados pueden consentir al monitoreo de security, ¿pueden significativamente consentir a la evaluación psicológica cuando podrían no comprender las implicaciones? ¿Cómo pueden las organizaciones obtener consentimiento para evaluar procesos inconscientes que, por definición, los individuos no son conscientes? Estas preguntas no tienen respuestas fáciles pero deben ser abordadas para implementación ética.

## 7.4 Direcciones de Investigación Futura

El framework CPF abre avenidas múltiples para investigación futura. La validación empírica permanece la necesidad más urgente. Mientras las fundaciones teóricas son fuertes, el testing sistemático a través de contextos organizacionales diversos es esencial. Estudios longitudinales que rastrean puntuaciones CPF e incidentes de security con el tiempo validarían capacidades predictivas. Estudios cross-culturales identificarían vulnerabilidades universales versus cultura-específicas.

La integración del machine learning ofrece posibilidades prometedoras para reconocimiento de patrones y predicción. Las redes neuronales podrían identificar patrones de vulnerabilidad sutiles que los humanos pierden. El procesamiento de lenguaje natural podría automatizar análisis lingüístico para evaluación de estrés y dinámicas de grupo. El aprendizaje por refuerzo podría optimizar estrategias de intervención basadas en resultados. Sin embargo, la integración ML debe mantener interpretabilidad—predicciones black box de estados psicológicos levantan preocupaciones éticas y prácticas.

El desarrollo de intervenciones representa una necesidad de investigación crítica. Mientras el CPF identifica vulnerabilidades, enfoques sistemáticos para abordarlas permanecen sub-desarrollados. ¿Cómo pueden las organizaciones abordar vulnerabilidades basadas en la autoridad sin socavar autoridad legítima? ¿Cómo pueden reducir estrés sin comprometer urgencia necesaria? La investigación sobre intervenciones de security informadas psicológicamente podría producir estrategias prácticas para mitigación de vulnerabilidad.

La intersección del CPF con otros frameworks merece exploración. ¿Cómo se relaciona el CPF al Cybersecurity Framework de NIST o ISO 27001? ¿Pueden indicadores psicológicos ser integrados con métricas de security técnicas en sistemas SIEM? ¿Podrían las categorías CPF mapear a controles específicos en frameworks de compliance? Esta investigación de integración podría facilitar adopción conectando intuiciones psicológicas a prácticas de security establecidas.

## 7.5 Integración con Frameworks Establecidos: El Ejemplo NIST CSF

El CPF no está diseñado para reemplazar frameworks de cybersecurity establecidos sino para aumentarlos abordando su punto ciego: la dimensión psicológica humana. Esta relación complementaria puede ser ilustrada mapeando el CPF al NIST Cybersecurity Framework (CSF) ampliamente adoptado.

Las funciones core del NIST CSF (Identificar, Proteger, Detectar, Responder, Recuperar) abordan primariamente controles técnicos y procedimentales. El CPF pro-

porciona la capa de inteligencia psicológica que potencia cada función:

- **Identificar:** Las evaluaciones CPF identifican proactivamente vulnerabilidades *psicológicas* organizacionales (ej. dependencia de la autoridad, patrones de estrés) que podrían llevar a vulnerabilidades de activos técnicos, enriqueciendo el proceso de identificación de activos.
- **Proteger:** Comprender estos patrones psicológicos permite el diseño de training de security y controles de acceso más efectivos y conscientes del humano que tienen en cuenta la carga cognitiva y la influencia social.
- **Detectar:** Los indicadores CPF sirven como señales de alerta temprana. Un Índice de Convergencia CPF en aumento puede empujar a los defensores a aumentar el monitoreo *antes* de que un ataque se manifieste técnicamente, desplazando la detección de reactiva a predictiva.
- **Responder/Recuperar:** Durante un incidente, dashboards CPF en tiempo real pueden informar estrategias de respuesta identificando si la organización está en un estado de groupthink o parálisis inducida por el estrés, permitiendo protocolos de comunicación y soporte decisional a medida que mitigan estas barreras psicológicas.

Este mapeo demuestra que el CPF se integra sin fisuras con prácticas de security existentes, proporcionando una capa faltante de poder predictivo e intuición humanocentric.

## 8 Limitaciones y Desafíos

A pesar de su rigor teórico y promesa práctica, el framework CPF enfrenta diversas limitaciones que deben ser reconocidas. La complejidad de la psicología humana significa que cualquier framework, no importa cuán comprensivo, captura solo verdad parcial. Los 100 indicadores, mientras extensivos, no pueden comprender todas las vulnerabilidades psicológicas. Casos límite, variaciones individuales, y propiedades emergentes aseguran que algunas vulnerabilidades escaparán a la detección.

El sesgo cultural representa una limitación significativa. El framework extrae primariamente de teorías psicológicas occidentales desarrolladas en poblaciones WEIRD (Western, Educated, Industrialized, Rich, Democratic). Los patrones psicológicos considerados universales podrían ser cultura-específicos. Las relaciones de autoridad, las respuestas al estrés, y las dinámicas de

grupo varían a través de culturas de maneras que el framework actual no captura completamente. La aplicación global requiere adaptación cultural y validación.

La naturaleza dinámica tanto de psicología como de tecnología crea objetivos móviles. A medida que las medidas de security evolucionan, así lo hacen las respuestas psicológicas a ellas. A medida que las capacidades AI avanzan, nuevas vulnerabilidades psicológicas emergen. El framework requiere actualización continua para mantener relevancia, pero esta evolución arriesga inconsistencia y creep de complejidad que podría socavar usabilidad.

Los desafíos de medición persisten a pesar de métodos preservando la privacidad. Los estados psicológicos son intrínsecamente subjetivos y variables. El mismo individuo podría puntuar diferentemente dependiendo de la hora del día, experiencias recientes, o contexto de medición. La agregación mejora confiabilidad pero pierde variación individual que podría ser relevante para la security. El sistema de puntuación ternario, mientras práctico, simplifica drásticamente fenómenos psicológicos complejos.

## 9 Conclusión

El Cybersecurity Psychology Framework representa una reconceptualización fundamental de los factores humanos en la cybersecurity. Reconociendo que las vulnerabilidades de security originan no en decisiones conscientes sino en procesos pre-cognitivos e inconscientes, el CPF proporciona un enfoque científicamente fundado para predecir y prevenir incidentes de security que los frameworks tradicionales no pueden abordar.

La integración de teoría psicoanalítica con psicología cognitiva y consideraciones específicas de la AI crea un modelo comprensivo que captura el espectro completo de vulnerabilidades psicológicas. Desde las observaciones de autoridad de Milgram a las dinámicas de grupo de Bion, desde las relaciones objetales de Klein a los sesgos cognitivos de Kahneman, el CPF sintetiza décadas de investigación psicológica en un framework de security accionable. Los 100 indicadores a través de 10 categorías proporcionan capacidad de evaluación granular mientras mantienen aplicabilidad práctica.

El diseño preservando la privacidad del framework y el enfoque agnóstico a la implementación abordan preocupaciones prácticas y éticas que han limitado intentos previos de integrar la psicología en la práctica de security. Focalizándose en patrones agregados en lugar de evaluación individual, el CPF proporciona inteligencia organizacional sin vigilancia individual. Mapeando a vulnerabilidades en lugar de prescribir soluciones, respeta la autonomía organizacional mientras proporciona intuiciones accionables.

Los casos de estudio de incidentes de security mayores—SolarWinds, Colonial Pipeline, y ataques mediados por AI—demuestran el poder explicativo y predictivo del CPF. Estos análisis revelan cómo vulnerabilidades psicológicas habilitaron ataques que defensas técnicas habrían debido prevenir. Más importante, muestran cómo la evaluación CPF habría podido identificar ventanas de vulnerabilidad antes de la explotación, habilitando intervención preventiva.

Las implicaciones se extienden más allá de aplicaciones de security inmediatas. El CPF contribuye a la comprensión teórica del comportamiento humano en ambientes digitales, enfoques prácticos para gestionar factores humanos en sistemas complejos, y frameworks éticos para evaluación psicológica en contextos organizacionales. Abre direcciones de investigación en integración de machine learning, desarrollo de intervenciones, y validación cross-cultural que podrían avanzar tanto security como psicología.

Sin embargo, el CPF no es una panacea. La complejidad de la psicología humana asegura que vulnerabilidades persistirán a pesar de los mejores esfuerzos. Variaciones culturales, desafíos de medición, y preocupaciones éticas requieren consideración cuidadosa en la implementación. El framework complementa en lugar de reemplazar medidas de security técnicas, abordando el componente humano de un desafío fundamentalmente socio-técnico.

A medida que las organizaciones enfrentan amenazas cada vez más sofisticadas que explotan la psicología humana con precisión científica, frameworks como el CPF se vuelven esenciales. La cuestión no es si considerar factores psicológicos en la security sino cómo hacerlo efectivamente y éticamente. El CPF proporciona una fundación para esta evolución crítica en la práctica de security.

El objetivo último no es eliminar la vulnerabilidad humana—una tarea imposible que requeriría eliminar la humanidad misma. En cambio, el CPF busca comprender, anticipar, y tener en cuenta las vulnerabilidades psicológicas en la estrategia de security. Solo reconociendo la complejidad plena de la psicología humana, incluyendo sus dimensiones inconscientes y pre-cognitivas, podemos construir posturas de security resilientes tanto a amenazas actuales como emergentes.

El viaje hacia security informada psicológicamente apenas ha comenzado. El CPF proporciona un mapa y una brújula, pero el camino debe ser caminado por organizaciones dispuestas a confrontar verdades incómodas sobre la naturaleza humana, dinámicas de poder, y los límites de soluciones tecnológicas. Para aquellos preparados para emprender este viaje, el CPF ofrece no solo security mejorada sino comprensión más profunda de las dinámicas humanas que moldean nuestro mundo digital.

## Agradecimientos

El autor agradece a las comunidades de cybersecurity y psicología por su diálogo continuo sobre los factores humanos en la security. Reconocimiento especial va a los investigadores que conectan disciplinas, haciendo conexiones que ningún campo solo podría alcanzar.

## Biografía del Autor

Giuseppe Canale es un profesional de cybersecurity certificado CISSP con training especializado en teoría psicoanalítica y psicología cognitiva. Con 27 años de experiencia en cybersecurity combinados con estudio profundo de procesos inconscientes y dinámicas de grupo, desarrolla enfoques nuevos a la security organizacional que integran perspectivas técnicas y psicológicas.

## Declaración de Disponibilidad de Datos

El framework CPF está libremente disponible para investigación e implementación. Las herramientas de evaluación y los datos de validación serán liberados siguiendo estudios piloto, con protecciones de privacidad apropiadas.

## Conflicto de Intereses

El autor declara ningún conflicto de intereses.

## A Resumen Guía de Implementación

Las organizaciones que implementan el CPF deberían comenzar con programas piloto en departamentos voluntarios, expandiendo gradualmente a medida que la experiencia se acumula. La evaluación inicial debería establecer baselines a través de todos los 100 indicadores, identificando vulnerabilidades prioritarias para intervención. Las salvaguardas de privacidad deben ser implementadas desde el principio, con estructuras de gobernanza claras y procesos de consentimiento. La integración con operaciones de security existentes debería ser gradual, aumentando en lugar de reemplazar procesos actuales. El refinamiento continuo basado en resultados asegura evolución del framework alineada con necesidades organizacionales.

## B Verificación Blockchain

## Timestamp

La versión del framework CPF descrita en este documento fue timestampada en la blockchain para protección de la propiedad intelectual:

- **Platform:** OpenTimestamps.org
- **Hash:** dfb55fc21e1b204c342aa76145f13-29fa6f095eeddc3aa83486fca91a580fa96
- **Block Height:** 909232
- **Timestamp:** 2025-08-09 CET

## References

- [1] Ajzen, I. (1991). The theory of planned behavior. *Organizational Behavior and Human Decision Processes*, 50(2), 179-211.
- [2] Beaumet, A., Sasse, M. A., & Wonham, M. (2008). The compliance budget: Managing security behaviour in organisations. *Proceedings of NSPW*, 47-58.
- [3] Bion, W. R. (1961). *Experiences in groups*. London: Tavistock Publications.
- [4] Bowlby, J. (1969). *Attachment and Loss: Vol. 1. Attachment*. New York: Basic Books.
- [5] Cialdini, R. B. (2007). *Influence: The psychology of persuasion*. New York: Collins.
- [6] Damasio, A. (1994). *Descartes' error: Emotion, reason, and the human brain*. New York: Putnam.
- [7] Gartner. (2023). *Forecast: Information Security and Risk Management, Worldwide, 2021-2027*. Gartner Research.
- [8] Jung, C. G. (1969). *The Archetypes and the Collective Unconscious*. Princeton: Princeton University Press.
- [9] Kahneman, D. (2011). *Thinking, fast and slow*. New York: Farrar, Straus and Giroux.
- [10] Kahneman, D., & Tversky, A. (1979). Prospect theory: An analysis of decision under risk. *Econometrica*, 47(2), 263-291.
- [11] Kernberg, O. (1998). *Ideology, conflict, and leadership in groups and organizations*. New Haven: Yale University Press.

- [12] Klein, M. (1946). Notes on some schizoid mechanisms. *International Journal of Psychoanalysis*, 27, 99-110.
- [13] LeDoux, J. (2000). Emotion circuits in the brain. *Annual Review of Neuroscience*, 23, 155-184.
- [14] Libet, B., Gleason, C. A., Wright, E. W., & Pearl, D. K. (1983). Time of conscious intention to act in relation to onset of cerebral activity. *Brain*, 106(3), 623-642.
- [15] Menzies Lyth, I. (1960). A case-study in the functioning of social systems as a defence against anxiety. *Human Relations*, 13, 95-121.
- [16] Milgram, S. (1974). *Obedience to authority*. New York: Harper & Row.
- [17] Miller, G. A. (1956). The magical number seven, plus or minus two. *Psychological Review*, 63(2), 81-97.
- [18] SANS Institute. (2023). *Security Awareness Report 2023*. SANS Security Awareness.
- [19] Selye, H. (1956). *The stress of life*. New York: McGraw-Hill.
- [20] Soon, C. S., Brass, M., Heinze, H. J., & Haynes, J. D. (2008). Unconscious determinants of free decisions in the human brain. *Nature Neuroscience*, 11(5), 543-545.
- [21] Verizon. (2023). *2023 Data Breach Investigations Report*. Verizon Enterprise.
- [22] Winnicott, D. W. (1971). *Playing and reality*. London: Tavistock Publications.
- [23] FireEye. (2020). *Highly Evasive Attacker Leverages SolarWinds Supply Chain to Compromise Multiple Global Victims With SUNBURST Backdoor*. FireEye Threat Research.
- [24] CISA. (2021). *Cyber Awareness Alert: DarkSide Ransomware: Best Practices for Preventing Business Disruption from Ransomware Attacks*. Alert Number AA21-131A.
- [25] Brundage, M., et al. (2024). *Toward Trustworthy AI Development: Mechanisms for Supporting Verifiable Claims*. arXiv preprint arXiv:2004.07213v2.
- [26] Cain, A. A., Edwards, B., & Still, J. D. (2024). *A Meta-Analysis of the Effectiveness of Security Awareness Training: Does Modality Matter?*. Journal of Cybersecurity, 10(1).