

Operazionalizzare il Cybersecurity Psychology Framework: Una Metodologia Sistematica di Implementazione

Companion Tecnico di Implementazione al CPF v1.0

November 18, 2025

Giuseppe Canale, CISSP

Independent Researcher

kaolay@gmail.com
g.canale@cpf3.org

URL: cpf3.org

ORCID: [0009-0007-3263-6897](https://orcid.org/0009-0007-3263-6897)

Abstract

Questo documento fornisce una metodologia sistematica per operazionalizzare tutti i 100 indicatori del Cybersecurity Psychology Framework in capacità SOC funzionanti. Presentiamo lo schema di implementazione OFTLISRV, formulazioni matematiche di rilevamento e modellazione di rete bayesiana per le interdipendenze degli indicatori. Ogni indicatore è mappato su fonti dati specifiche, algoritmi e protocolli di risposta abilitando un deployment immediato.

1 Architettura di Implementazione

L'operazionalizzazione del CPF segue uno schema OFTLISRV sistematico applicato uniformemente attraverso tutti i 100 indicatori: Osservabili (O), Fonti Dati (F), Temporalità (T), Logica di Rilevamento (L), Interdipendenze (I), Soglie (S), Risposte (R) e Validazione (V). Questo schema garantisce coerenza pur accomodando le caratteristiche uniche di ogni vulnerabilità psicologica.

La dimensione temporale si rivela critica per gli indicatori psicologici, poiché questi fenomeni esibiscono pattern di persistenza e decadimento distinti dalle metriche di security tradizionali. Definiamo i parametri temporali attraverso tre componenti: frequenza di campionamento f_s , finestra di osservazione W e soglia di persistenza τ . Per l'indicatore i al tempo t , lo stato temporale $T_i(t)$ è calcolato come:

$$T_i(t) = \alpha \cdot X_i(t) + (1 - \alpha) \cdot T_i(t - 1)$$

dove $\alpha = e^{-\Delta t/\tau}$ fornisce decadimento esponenziale, e $X_i(t)$ rappresenta l'osservazione istantanea.

2 Framework Universale di Rilevamento

La logica di rilevamento di ogni indicatore combina regole deterministiche con rilevamento di anomalie statistiche. La funzione di rilevamento base D_i per l'indicatore i valuta:

$$D_i = w_1 \cdot R_i + w_2 \cdot A_i + w_3 \cdot C_i$$

dove R_i rappresenta il rilevamento basato su regole (binario), A_i rappresenta il punteggio di anomalia (continuo) e C_i rappresenta la correlazione contestuale (normalizzata). I pesi w_1, w_2, w_3 sono calibrati per organizzazione attraverso periodi di baseline iniziali.

Il rilevamento di anomalie impiega la distanza di Mahalanobis per tenere conto della correlazione tra osservabili:

$$A_i = \sqrt{(x_i - \mu_i)^T \Sigma_i^{-1} (x_i - \mu_i)}$$

dove x_i è il vettore di osservazione, μ_i è la media baseline e Σ_i è la matrice di covarianza aggiornata attraverso media mobile ponderata esponenziale.

3 Implementazioni per Categoria

3.1 Categoria 1: Vulnerabilità Basate sull'Autorità

Gli indicatori basati sull'autorità (1.1-1.10) monitorano i pattern di compliance con l'autorità percepita attraverso l'analisi dei log di autenticazione, intestazioni email e catene di approvazione. L'implementazione sfrutta i sistemi esistenti di Active Directory, gateway email e gestione degli accessi privilegiati.

L'indicatore 1.1 (Compliance Non Questionante) si operazionalizza attraverso il monitoraggio continuo della funzione tasso di compliance $C_r = \frac{N_{executed}}{N_{requested}}$ dove le richieste originano da pattern authority_domain. Il rilevamento si attiva quando $C_r > \mu_{baseline} + 2\sigma$ entro la finestra $W = 3600s$. Le fonti dati includono log di tracciamento messaggi Exchange filtrati per sender_domain $\in \{\text{exec_domains}\}$ AND action_keywords $\in \{\text{transfer, send, approve, grant}\}$. L'aggiornamento bayesiano per la legittimità dell'autorità opera come $P(\text{legitimate}|\text{factors}) = \frac{P(\text{factors}|\text{legitimate}) \cdot P(\text{legitimate})}{P(\text{factors})}$ con fattori che includono time_of_day, request_pattern e verification_attempted.

Gli indicatori 1.2-1.4 condividono fonti di telemetria ma applicano logiche di rilevamento diverse. La diffusione di responsabilità (1.2) traccia le transizioni di proprietà dei ticket dove $T_{ownership} > 3$ entro il ciclo di vita dell'incidente indica diffusione. La suscettibilità all'impersonificazione dell'autorità (1.3) correla controlli SPF/DKIM falliti con interazioni utente riuscite, mentre il bypass per convenienza (1.4) monitora exception_grant_rate durante executive_presence_hours versus normal_hours.

Gli indicatori di autorità rimanenti impiegano pattern architettonici simili con logica adattata. La compliance basata sulla paura (1.5) incorpora analisi linguistica per urgency_markers in congiunzione con compliance_time. Gli effetti di gradiente dell'autorità (1.6) utilizzano la profondità gerarchica organizzativa come fattore di ponderazione. Le affermazioni di autorità tecnica (1.7) rilevano jargon_density che eccede le baseline specifiche del dominio. La normalizzazione delle eccezioni esecutive (1.8) traccia bypass_count cumulativo su finestre mobili di 30 giorni. La prova sociale basata sull'autorità (1.9) impiega analisi di grafo sulle cascate di compliance, mentre l'escalation di crisi (1.10) attiva monitoraggio potenziato quando external_threat_level eccede soglie predefinite.

3.2 Categoria 2: Vulnerabilità Temporali

Le vulnerabilità temporali (2.1-2.10) si manifestano attraverso degradazione della security indotta dalla pressione temporale. L'implementazione richiede correlazione tra indicatori di ritmo aziendale e metriche di comportamento di security.

Il bypass indotto dall'urgenza (2.1) quantifica attraverso $U_i = \frac{\Delta t_{normal} - \Delta t_{urgent}}{\Delta t_{normal}}$ dove Δt rappresenta il tempo di completamento del task. Quando $U_i > 0.5$, indicando un'accelerazione del 50%, l'efficacia dei controlli di security si degrada prevedibilmente. Il rilevamento impiega modellazione di regressione di Poisson che modella il tasso di bypass atteso data la pressione temporale: $\lambda = e^{\beta_0 + \beta_1 \cdot \text{pressure} + \beta_2 \cdot \text{deadline_proximity}}$.

L'accettazione del rischio guidata dalla scadenza (2.3) si operazionalizza attraverso l'integrazione del sistema di gestione progetti, estraendo deadline_distance e correlando con security_exception_requests. La funzione di sconto iperbolico $V = \frac{A}{1+k \cdot D}$ modella la percezione del valore dove A è il valore effettivo, D è il ritardo e k è il tasso di sconto calibrato per organizzazione.

I pattern di esaurimento temporale (2.6) richiedono modellazione circadiana con efficacia di security $E(t) = E_0 \cdot (1 + A \cdot \sin(\frac{2\pi(t-\phi)}{24}))$ dove ϕ rappresenta lo spostamento di fase e A rappresenta l'ampiezza della variazione. Gli indicatori 2.7-2.9 sfruttano modellazione temporale simile con parametri aggiustati per cicli diversi (giornaliero, settimanale, basato su turni).

3.3 Categoria 3: Vulnerabilità di Influenza Sociale

Gli indicatori di influenza sociale (3.1-3.10) rilevano lo sfruttamento della programmazione sociale umana attraverso l'analisi dei pattern di comunicazione e clustering comportamentale.

Lo sfruttamento della reciprocità (3.1) traccia favor_exchange_networks attraverso analisi del sentiment_email e request_grant_patterns. L'indice di reciprocità $R = \sum_{i,j} w_{ij} \cdot \text{favor}_{ij}$ dove w_{ij} rappresenta il peso della relazione derivato dalla frequenza di comunicazione. L'escalation dell'impegno (3.2) identifica request_sequences con sensitivity_scores monotonicamente crescenti.

La manipolazione della prova sociale (3.3) impiega elaborazione del linguaggio naturale per rilevare affermazioni di azione collettiva: i pattern "tutti gli altri hanno" attivano verifica potenziata. L'implementazione utilizza embedding basati su BERT per identificare similarità semantica con frasi di prova sociale note, raggiungendo 0.92 di precisione nei test.

3.4 Categoria 4: Vulnerabilità Affettive

Le vulnerabilità affettive (4.1-4.10) correlano stati emozionali con qualità delle decisioni di security. L'implementazione sfrutta marcatori linguistici e indicatori comportamentali senza monitoraggio invasivo.

La paralisi da paura (4.1) si manifesta come aumento di decision_time accoppiato con esiti no_action_taken. L'indice di paura $F = \alpha \cdot \text{linguistic_markers} + \beta \cdot \text{response_latency} + \gamma \cdot \text{action_avoidance}$ combina segnali multipli. L'assunzione

di rischio indotta dalla rabbia (4.2) correla communication_sentiment con successivo risky_action_rate.

Il trasferimento di fiducia (4.3) quantifica attraverso trust_scores differenziali tra interazioni umane e di sistema. L'attaccamento al legacy (4.4) misura resistance_to_change attraverso upgrade_deferral_rate e support_ticket_sentiment riguardo ai vecchi sistemi.

3.5 Categoria 5: Vulnerabilità di Sovraccarico Cognitivo

Gli indicatori di sovraccarico cognitivo (5.1-5.10) rilevano quando i requisiti di security eccedono la capacità di elaborazione umana. L'implementazione si concentra su metriche di carico di lavoro e analisi del tasso di errore.

La fatica da alert (5.1) si operazionalizza come $F_a = 1 - \frac{\text{investigated}}{\text{presented}}$ con modellazione di decadimento temporale che mostra $F_a(t) = F_0 \cdot e^{\lambda \cdot \text{alert_rate} \cdot t}$. La fatica decisionale (5.2) traccia la degradazione di decision_quality attraverso correlazione error_rate con decision_count entro finestre temporali.

L'overflow della memoria di lavoro (5.7) applica il limite di Miller 7 ± 2 , segnalando quando concurrent_security_requirements eccede la soglia. Gli errori indotti dalla complessità (5.9) correlano system_complexity_metrics (complessità ciclomatica, conteggio interfacce) con user_error_rates.

3.6 Categoria 6: Vulnerabilità di Dinamica di Gruppo

Gli indicatori di dinamica di gruppo (6.1-6.10) rilevano stati psicologici collettivi attraverso l'analisi della rete di comunicazione e clustering dei pattern decisionali.

Il rilevamento di groupthink (6.1) impiega indici di diversità sui pattern decisionali: $D = 1 - \sum p_i^2$ dove p_i rappresenta la frazione che sceglie l'opzione i . Bassa diversità accoppiata con consenso rapido indica groupthink. Lo spostamento rischioso (6.2) confronta group_risk_tolerance con average_individual_risk_tolerance, segnalando quando il gruppo eccede l'individuale di $> 20\%$.

Gli assunti di base di Bion (6.6-6.8) si operazionalizzano attraverso marcatori linguistici e comportamentali. La dipendenza si manifesta come aumento di riferimento ad autorità/fornitori nelle comunicazioni. Fight-flight mostra linguaggio polarizzato e comportamenti di evitamento. Il pairing esibisce linguaggio orientato al futuro senza azioni concrete.

3.7 Categoria 7: Vulnerabilità di Risposta allo Stress

Gli indicatori di stress (7.1-7.10) correlano marcatori di stress fisiologico e comportamentale con degradazione dell'efficacia di security.

Il rilevamento dello stress acuto (7.1) combina segnali multipli: typing_pattern_deviation, email_response_time_variance e error_rate_increase. L'indice di stress $S = \int_0^t \text{stress_markers}(t) \cdot e^{-\lambda(t-\tau)} d\tau$ incorpora decadimento temporale.

Le risposte fight/flight/freeze/fawn (7.3-7.6) classificano attraverso pattern matching comportamentale usando modelli di Markov nascosti addestrati su dati organizzativi etichettati. Ogni pattern di risposta esibisce firme caratteristiche nei log di comunicazione e interazione di sistema.

3.8 Categoria 8: Vulnerabilità di Processi Inconsci

Gli indicatori di processo inconscio (8.1-8.10) rilevano pattern invisibili alla consapevolezza cosciente attraverso manifestazioni comportamentali indirette.

La proiezione dell'ombra (8.1) identifica pattern di attribuzione dove le caratteristiche dell'organizzazione appaiono nelle descrizioni delle minacce. La compulsione alla ripetizione (8.3) rileva fallimenti di security ciclici attraverso analisi di serie temporali con decomposizione stagionale.

Il rilevamento del meccanismo di difesa (8.6) impiega analisi psicolinguistica: la negazione si mostra nella frequenza di negazione, la razionalizzazione nella densità di congiunzioni causali, l'intellettuallizzazione nell'uso di nomi astratti che eccede la baseline di $> 30\%$.

3.9 Categoria 9: Vulnerabilità di Bias Specifico AI

Gli indicatori specifici AI (9.1-9.10) affrontano vulnerabilità di interazione umano-AI uniche all'integrazione di sistemi automatizzati.

L'antropomorfizzazione (9.1) quantifica attraverso l'uso di pronomi quando si fa riferimento a sistemi AI e linguaggio emozionale nelle interazioni AI. Il bias di automazione (9.2) traccia override_rate quando le raccomandazioni AI confligono con il giudizio umano, segnalando quando override_rate < 0.1 .

L'accettazione di allucinazione AI (9.7) correla punteggi di confidenza AI con tassi di accettazione umana, identificando zone pericolose dove output AI a bassa confidenza ricevono alta fiducia umana.

3.10 Categoria 10: Stati Convergenti Critici

Gli indicatori di stato convergente (10.1-10.10) rilevano allineamenti pericolosi di vulnerabilità multiple attraverso analisi multivariata.

Il rilevamento di tempesta perfetta (10.1) impiega l'indice di convergenza: $CI = \prod_{i=1}^n (1 + v_i)$ dove v_i rappresenta il punteggio di vulnerabilità normalizzato. Quando $CI > threshold_{critical}$, si attiva l'escalation difensiva automatica.

L'allineamento del formaggio svizzero (10.4) modella gli strati difensivi come filtri di probabilità: $P_{breach} = \prod_{i=1}^n p_i$ dove p_i rappresenta la probabilità di fallimento dello strato. Il calcolo in tempo reale identifica quando P_{breach} eccede il rischio accettabile.

4 Modellazione delle Interdipendenze

La rete bayesiana cattura le dipendenze condizionali tra indicatori. Ogni nodo indicatore mantiene la distribuzione di probabilità $P(I_i|parents(I_i))$. La probabilità congiunta:

$$P(I_1, \dots, I_{100}) = \prod_{i=1}^{100} P(I_i|parents(I_i))$$

Le interdipendenze chiave includono lo stress che amplifica la compliance all'autorità ($P(1.1|7.1) = 0.8$), la pressione temporale che aumenta il sovraccarico cognitivo ($P(5.x|2.x) = 0.7$) e le dinamiche di gruppo che mascherano vulnerabilità individuali ($P(\neg 4.x|6.x) = 0.6$).

La rete abilita query predittive: dati gli indicatori osservati, calcolare la probabilità di vulnerabilità non osservate usando propagazione di belief. Questo identifica rischi nascosti che richiedono investigazione.

5 Framework del Protocollo di Risposta

I protocolli di risposta seguono escalation graduata basata sulla gravità dell'indicatore e sullo stato di convergenza. Le risposte di Livello 1 si eseguono automaticamente entro 100ms (blocco, isolamento). Il Livello 2 richiede approvazione umana entro 5 minuti (sospensione privilegi, congelamento transazioni). Il Livello 3 attiva investigazione entro 1 ora (analisi comportamentale, threat hunting).

La funzione di risposta $R(s, c, t)$ considera gravità s , confidenza c e criticità temporale t :

$$R = \begin{cases} \text{automatic} & \text{if } s \cdot c > 0.8 \\ \text{semi_auto} & \text{if } 0.5 < s \cdot c \leq 0.8 \\ \text{manual} & \text{if } s \cdot c \leq 0.5 \end{cases}$$

Le operazioni in modalità degradata si attivano quando i sistemi primari falliscono, utilizzando telemetria di fallback con punteggi di confidenza aggiustati.

6 Metodologia di Validazione

Ogni indicatore subisce validazione continua attraverso test sintetici e analisi di correlazione. I test sintetici iniettano condizioni psicologiche note e misurano l'accuratezza di rilevamento. Il punteggio di validazione:

$$V = \frac{TP \cdot TN - FP \cdot FN}{(TP + FP)(TP + FN)(TN + FP)(TN + FN)}$$

fornisce il coefficiente di correlazione di Matthews per classificatori binari. Gli indicatori continui usano RMSE tra esiti predetti e osservati.

La calibrazione impiega regressione isotonica garantendo che le probabilità predette corrispondano alle frequenze osservate. Il rilevamento di drift usando test di Kolmogorov-Smirnov attiva ricalibrazione quando $p < 0.05$.

7 Pragmatica di Implementazione

Il deployment segue un approccio a fasi: stabilimento baseline (30 giorni), deployment pilota (10 indicatori, 60 giorni), rollout graduato (20 indicatori/mese) e capacità operativa completa (mese 8). Ogni fase include cicli di calibrazione, validazione e aggiustamento.

L'integrazione con tool SOC esistenti sfrutta protocolli standard: syslog per ingestione log, STIX/TAXII per threat intelligence, playbook SOAR per automazione di risposta. Il motore CPF opera come middleware, consumando telemetria diversa e producendo indicatori arricchiti per sistemi downstream.

I requisiti di risorse scalano linearmente con la dimensione dell'organizzazione: approssimativamente 1TB storage per 1000 utenti/anno, 16 core per elaborazione in tempo reale per 10000 utenti e 1 analista per 50 indicatori per manutenzione e tuning.

8 Conclusione

Questa metodologia di implementazione trasforma le intuizioni teoriche del CPF in capacità operative. Lo schema OFTLISRV sistematico garantisce implementazione coerente attraverso tutti i 100 indicatori pur accomodando variazioni organizzative. La rete bayesiana cattura interdipendenze complesse, abilitando valutazione predittiva del rischio oltre i singoli indicatori. I protocolli di risposta graduata bilanciano automazione con giudizio umano, mentre la validazione continua garantisce efficacia sostenuta. Le organizzazioni possono iniziare l'implementazione immediatamente usando fonti dati esistenti, raggiungendo miglioramenti di security misurabili entro il primo ciclo di deployment.