

# Appunti Tesi

## 1 Collegamenti tra SSG e CRA per IACS

CRA for IACS = cybersecurity risk assessment for industrial automation control systems

SUC = system under consideration

Considerando un Stackelberg Security Game (SSG) a somma zero, è possibile derivare input e output del gioco:

- **Input:**

- insieme di asset condivisi (visibili sia dal difensore che dall'attaccante);
- insieme delle risorse limitate per il difensore, che corrispondono ad una restrizione sulle proprie possibilità di difesa dei target, ovvero al numero massimo di assets che il difensore riesce a proteggere contemporaneamente;
- insieme dei valori di payoff, che definiscono i guadagni/perdite sia per il difensore che per l'attaccante in caso di un attacco fallimentare o meno ad un determinato target, e tenendo in considerazione anche se un determinato asset è protetto o meno (il payoff di un determinato target definisce quanto esso è rilevante per il relativo agente razionale);
- strategia mista iniziale per il difensore, che definisce la probabilità di ogni target di essere protetto dal giocatore.

- **Output:** l'output del SSG è rappresentato da una condizione di equilibrio, detto Stackelberg Equilibrium (SE), rappresentata da una coppia di strategie (una per il difensore e una per l'attaccante) tale per cui: il difensore e l'attaccante applicano la propria strategia ottimale. Generalmente si considera lo Strong Stackelberg Equilibrium (SSE), nel quale si aggiunge la condizione per cui in caso di più strategie ottimali per l'attaccante, quest'ultimo applica la strategia ottimale che massimizzi il payoff del difensore.

Un SSG in sé si potrebbe considerare come costituito dalle funzioni di payoff/utilità, ovvero dalle funzioni che “combinano” le varie componenti del gioco per arrivare ad ottenere un SE.

Negli input del SSG possiamo considerare, o meno, la strategia iniziale del difensore; se non la si considera, allora la possiamo determinare randomicamente. Inoltre, non si considera negli input di un SSG la strategia dell'attaccante, in quanto agendo per secondo può essere contraddistinto da un tipo di strategia arbitraria a seconda dei casi che si vogliono considerare (es. strategia greedy). Nel modello classico di un SSG, l'attaccante può selezionare un unico bersaglio, quindi le risorse dell'attaccante sono implicitamente limitate. In caso di modelli più complessi, è possibile ipotizzare un attaccante con risorse multiple (contemporaneamente può attaccare più assets), ma comunque limitate. Un'assunzione da fare in un SSG è che il payoff di una generica azione dipende solo dall'asset attaccato, e dal fatto che questo sia difeso o meno dal difensore. Inoltre, il difensore

avrà come obiettivo quello di minimizzare il danno massimo che può subire, sapendo che l'attaccante risponderà in modo ottimale. È attraverso i valori di payoff che vengono rappresentati costi e benefici delle strategie di ognuno delle 2 tipologie di agenti. L'ipotesi del gioco a somma zero rende necessari solo i payoff di uno dei 2 agenti, in quanto quelli dell'altro sarebbero ricavabili di conseguenza, andando così a ridurre gli input del gioco.

Da qui si possono ricercare delle **similitudini** tra la struttura di un SSG e il processo di CRA per IACS.

Nella fase **ZCR 5.1** vengono identificate le minacce per il sistema e queste ultime in un SSG possono essere comparate alle varie tipologie di attaccanti che possono essere interessati nel bersagliamento degli asset condivisi. Il tipo di un attaccante in un modello di gioco, infatti, può essere considerato un input implicito a seconda del quale possono essere definiti i vari comportamenti di attacco.

Inizialmente era stato pensato di identificare l'impatto, ovvero l'elemento identificato nella sezione ZCR 5.3 del processo di CRA, con il valore di payoff del difensore. Questo ragionamento era stato fatto in quanto il concetto di "impatto" in un SSG potrebbe essere tradotto in "quanto valore ha quel determinato asset per il difensore", e questa immagine viene proprio rappresentata nel valore di payoff, legato al target, del difensore. Però, in seguito al dialogo con un'esperta di Resiltech, si è arrivati alla conclusione che questa associazione sarebbe andata in contrasto con la letteratura relativa all'impatto nel processo di CRA.

La fase **ZCR 5.4** del processo di CRA per IACS ha come obiettivo quello di identificare la "likelihood" di un attacco ad ogni asset del SUC. In un SSG, tale fattore, rappresentato dalla metrica di Attack Feasibility Rating (**AFR**), corrisponde alla probabilità che l'attaccante decida di intraprendere un atto offensivo nei confronti di uno specifico target. Quindi la likelihood associata ad ogni asset andrebbe a corrispondere ad un elemento della strategia mista dell'attaccante.

L'obiettivo principale della fase **ZCR 5.5** è quello di determinare il "**unmitigated cybersecurity risk**" per ogni minaccia, e quindi assegnare un "risk score" ad ogni asset tramite la combinazione dell'AFR e l'impatto associati ad un target. Poiché l'impatto in sé non è associabile con il payoff del difensore in un SSG, allora è possibile associare quest'ultimo elemento all'unmitigated risk. Inoltre, una volta arrivati a uno SSE, la strategia ottimale del difensore corrisponde ad un vettore in cui ogni elemento costituisce la probabilità che il difensore decida di proteggere un determinato target. Ognuno degli elementi di questo vettore può essere considerato come il SL-T per ogni asset considerato (fase ZCR 5.6)

Un elemento che non è presente nel modello classico degli SSGs, ma che bensì viene considerato nel CRA per IACS, è la presenza di **path tra asset** che l'attaccante può sfruttare per il raggiungimento del proprio obiettivo finale. Tra i target interdipendenti si suppone che, per una minaccia, gli asset più facilmente raggiungibili siano quelli che si trovano all'inizio di un path; questi però rappresentano anche quelli di minor valore per l'attaccante. Contrariamente a ciò, gli asset che si trovano verso la fine del percorso di attacco sono più difficili da raggiungere ma costituiscono anche il vero obiettivo dell'attaccante. Viste queste considerazioni, è possibile una rappresentazione dell'interdipendenza tra gli asset, che in un SSG sono effettivamente indipendenti tra loro, tramite l'assegnazione di specifici valori di payoff all'attaccante. In questo senso, si avreb-

bero valori di payoff alti per l'attaccante in corrispondenza di asset significativi per quest'ultimo, e che quindi si troverebbero in coda al path che l'attaccante percorrerebbe in un corrispettivo sistema; mentre invece si assegnerebbero valori di payoff più contenuti a target che invece non sono considerati dall'attaccante come possibile obiettivo finale, e che quindi si troverebbero in cima a degli ipotetici path di attacco. A tale proposito, sarebbe possibile pensare a una numerazione degli asset, tramite indice, che rappresenterebbe la disposizione dei target in un path del sistema e strutturare i valori di payoff del difensore in relazione a tali indici.

Per quanto riguarda i valori di payoff dei 2 agenti di un SSG, dovranno essere fatte delle assunzioni, considerando che il SSG preso in considerazione è non a somma zero. In caso di attacco riuscito (obiettivo non difeso)

- Attaccante riceve un payoff positivo (successo dell'attacco):

$$U_A^u(t_i) > 0$$

Tipicamente proporzionale al valore dell'obiettivo (es: danno causato, guadagno ottenuto).

- Difensore subisce una perdita (fallimento difensivo):

$$U_D^u(t_i) < 0$$

Rappresenta il costo/danno subito per l'attacco riuscito.

In caso di attacco fallito (obiettivo difeso)

- Attaccante riceve un payoff negativo o nullo:

$$U_A^c(t_i) \leq 0$$

Può rappresentare perdita di risorse, rischio di cattura, o semplicemente "nessun guadagno".

- Difensore ottiene un payoff positivo o nullo:

$$U_D^c(t_i) \geq 0$$

Rappresenta il successo della protezione. Può essere zero (se difesa non produce valore diretto) o positivo (prevenzione di danni).

## 2 Path negli SSG

Questo lavoro, presentato in *Decision Support Systems 148 (2021) 113599* [1], propone un **sistema di supporto decisionale per la cyber-sicurezza** [2, 3]. L'obiettivo principale è assistere le organizzazioni nella selezione di un **portafoglio ottimale di controlli di sicurezza** per contrastare attacchi a fasi multiple (multi-stage attacks) [2, 3]. Gli autori sono Yunxiao Zhang e Pasquale Malacaria [1].

Il sistema è strutturato in diversi componenti interconnessi [2-4]:

- Un **ottimizzazione preventiva**: Serve per selezionare un portafoglio difensivo iniziale prima che gli attacchi si verifichino. Mira a minimizzare il rischio di sicurezza potenziale [2, 3, 5, 6].
- Un **meccanismo di apprendimento**: Utilizza approcci consolidati come gli Hidden Markov Models (HMMs) [4], basati su lavori precedenti [5, 7-9], per rilevare e stimare possibili attacchi in corso [2-4, 6]. L'inferenza basata su HMM restituisce un vettore di credenza (belief vector) sullo stato dell'attaccante [4, 10].
- Un **ottimizzazione online**: Questa componente seleziona un portafoglio ottimale per contrastare gli attacchi in corso rilevati dal meccanismo di apprendimento [2, 3]. Interviene quando viene rilevato un attacco e il difensore ha informazioni incomplete sullo stato esatto dell'attaccante [8, 10].

Per modellare attacchi a fasi multiple realistici, il sistema si avvale di un **grafo di attacco probabilistico** [7, 11].

- I nodi nel grafo rappresentano gli **"stati di privilegio"** dell'attaccante [7, 11, 12].
- Gli archi (edges) rappresentano le **vulnerabilità** che l'attaccante può sfruttare per passare tra gli stati [7, 11].
- Ogni arco è associato alla **probabilità di successo dello sfruttamento**, influenzata da una probabilità di base e dall'efficacia dei controlli applicati [13, 14].
- Il **rischio di sicurezza** è definito come la massima probabilità di successo che un attaccante raggiunga l'obiettivo percorrendo uno qualsiasi dei cammini possibili nel grafo [7, 12, 15].

L'ottimizzazione preventiva, relativa all'investimento iniziale in sicurezza, è formulata come un **gioco di Stackelberg standard** [5, 6]. In questo contesto, il difensore agisce come leader, scegliendo un portafoglio di sicurezza preventiva per minimizzare il rischio potenziale, mentre l'attaccante è il follower che mira a massimizzare il rischio selezionando il percorso più critico [5, 6]. Il problema viene risolto efficientemente convertendo il problema originale non lineare in un problema di programmazione lineare (LP) grazie alle proprietà delle matrici totalmente unimodulari e alla dualità forte [1, 16, 17].

L'ottimizzazione online, che gestisce la difesa contro gli attacchi in corso con informazioni incomplete sullo stato dell'attaccante (rappresentato da un belief vector) [8, 10], è formulata come un **gioco di Stackelberg Bayesiano** [8, 10, 18, 19]. L'obiettivo è minimizzare il rischio di sicurezza atteso, calcolato come l'aspettativa sui rischi per ciascun tipo di attaccante (la cui probabilità è data dal belief vector) [10, 20].

Un contributo tecnico chiave del lavoro è un **nuovo algoritmo efficiente per risolvere le ottimizzazioni online (giochi di Stackelberg Bayesiani) su grafi di attacco probabilistici** [10, 21]. I metodi classici come la trasformazione di Harsanyi [8, 18, 19, 22] sono inefficaci per grandi spazi di attacco a causa della potenziale matrice di payoff esponenzialmente grande [23, 24]. Anche approcci più recenti come DOBSS [4, 18, 25], HBGS [6, 18], e HUNTER [10, 18] non scalano sufficientemente bene per grafi di attacco di grandi dimensioni nel dominio della cyber-sicurezza [18, 26].

Il nuovo approccio sfrutta le proprietà delle **matrici totalmente unimodulari** e della **dualità forte** per convertire l'ottimizzazione online in un problema di **Programmazione Conica a**

**Numeri Interi Misti (MICP)** con coni esponenziali [2, 18, 21, 27]. Questo MICP può essere risolto efficientemente utilizzando risolutori esistenti (come MOSEK versione 9.2) [21, 28].

Gli autori dimostrano che il loro approccio per i giochi di Stackelberg Bayesiani su grafi di attacco è **molto efficiente**, richiedendo un numero significativamente inferiore di variabili di ottimizzazione rispetto a DOBSS e alla trasformazione di Harsanyi [25, 29, 30]. Le valutazioni numeriche mostrano che l'ottimizzazione online è efficiente anche per scenari realistici con un gran numero di nodi e tipi di attaccante [26, 28, 31], e fornisce **miglioramenti significativi nella mitigazione degli attacchi in corso** rispetto agli approcci precedenti [2, 32, 33].

Il documento include un **caso di studio** basato su uno scenario di rete universitaria per illustrare l'applicazione pratica del sistema [34, 35].

In sintesi, il lavoro offre un sistema completo per la cyber-difesa, con un particolare focus su un metodo efficiente basato su MICP per risolvere il complesso problema dell'ottimizzazione delle contromisure contro attacchi in corso in presenza di informazioni incomplete sullo stato dell'attaccante [36].

### 3 Network Games e SSG

Questo documento di ricerca esamina i **giochi di sicurezza interdipendenti** (Interdependent Security - IDS) nel contesto dei **giochi di Stackelberg**, combinando elementi della teoria dei **giochi di rete** e della teoria dei giochi sequenziali (Stackelberg) [1, 2]. L'obiettivo è studiare le interazioni strategiche tra agenti interconnessi, in particolare nel dominio della **sicurezza informatica** [3].

Vengono distinti due tipi principali di giochi strategici utilizzati nel contesto della sicurezza [1]:

- I **giochi di rete** (o giochi IDS): Catturano le interazioni strategiche tra agenti interconnessi che agiscono **simultaneamente** [1, 3, 4]. Sono spesso usati per studiare gli **investimenti in sicurezza** e le esternalità [1, 3, 5-7]. Il concetto di soluzione tipico è l'**Equilibrio di Nash** (NE), ottenuto risolvendo un sistema di equazioni a punto fisso [1, 3]. In questi modelli, l'impatto dello sforzo di sicurezza è catturato nella funzione di utilità degli agenti, spesso come funzione crescente di uno "sforzo totale ponderato" che dipende anche dagli sforzi degli altri agenti [3, 8, 9]. Non sempre un avversario (attacker) è modellato esplicitamente [3].
- I **giochi di Stackelberg**: Modellano le **mosse sequenziali** tra agenti, tipicamente categorizzati come leader (primi a muovere) e follower (secondi a muovere) [1]. Il concetto di soluzione corrispondente è l'**Equilibrio Perfetto di Sottogioco** (SPE), ottenuto tipicamente tramite induzione all'indietro [1, 10, 11]. Sono spesso usati per modellare scenari attaccante-difensore [2, 5, 12, 13].

Lo studio propone un modello che combina l'interdipendenza dei giochi di rete con la sequenzialità dei giochi di Stackelberg [2]. Nello scenario considerato, vi è una rete di **difensori interconnessi che sono i primi a muovere simultaneamente**, e uno o più **attaccanti che sono i secondi a muovere simultaneamente**, rispondendo alle azioni dei difensori [2, 10, 14]. I difensori prendono le loro decisioni sia anticipando la migliore risposta degli attaccanti sia l'impatto delle azioni degli altri difensori attraverso il grafo di interazione [10]. La soluzione è un tipo speciale di SPE, derivato tramite induzione all'indietro, che ora coinvolge la risoluzione di sistemi di equazioni a punto fisso in entrambi gli stadi [10, 11, 15].

Il modello utilizza una **matrice G** per rappresentare le dipendenze tra gli agenti [14]. Questa matrice può essere suddivisa in blocchi [13]: **D** (dipendenza tra difensori), **A** (dipendenza tra

attaccanti), **B** (dipendenza dei difensori dagli attaccanti), e **C** (dipendenza degli attaccanti dai difensori). Ogni agente sceglie un livello di "sforzo" ( $x_i \in R_+$ ) [14]. Le funzioni di utilità ( $u_i$ ) sono basate sul concetto di "**sforzo totale ponderato**", dove il beneficio di un giocatore dipende dalla somma ponderata del proprio sforzo e di quello degli altri ( $x_i + \sum_{j \neq i} g_{ij}x_j$ ) [8, 9].

Un concetto chiave è il **valore autarchico** ( $q_i$ ), che rappresenta lo sforzo ottimale di un giocatore quando è isolato [9, 16]. Se il suo sforzo effettivo è inferiore al valore autarchico perché beneficia dello sforzo dei vicini, si parla di **free riding** [16]. Il grado di free riding individuale di un difensore è misurato dallo **score di free riding individuale** ( $s_i$ ), definito come il beneficio (spillover) ricevuto dagli sforzi dei vicini ( $\sum_{j \in D \setminus \{i\}} d_{ij}x_j$ ) [17, 18].

Una delle principali scoperte riguarda l'**equivalenza**, per una classe speciale di SPE chiamati **Equilibri Perfetti di Sottogioco Interni (I-SPE)**, tra il gioco dei difensori nello stadio iniziale del gioco sequenziale e un gioco di rete a mosse simultanee giocato sui difensori sulla base dello **Schur complement** ( $G/A = D - BA^{-1}C$ ) della matrice  $A$  nella matrice  $G$  [5, 19, 20]. Questo viene chiamato **gioco di rete ridotto** [5, 19-21]. In determinate **condizioni sulla struttura della rete** ( $G$  P-matrix,  $B$  e  $C$  non-positive,  $A$  Z-matrix) [22], esiste un **unico I-SPE** che può essere calcolato analiticamente [19]. Queste condizioni limitano la magnitudine delle dipendenze reciproche dello stesso tipo [22] e catturano l'effetto dannoso degli attaccanti sui difensori ( $B$  non-positive) e l'effetto deterrente dei difensori sugli attaccanti ( $C$  non-positive) [23].

Il documento analizza come la **sequenzialità** influenzi le azioni di equilibrio dei giocatori rispetto a un gioco a mosse simultanee equivalente [6, 24, 25]:

- Si osserva una **riduzione del valore autarchico** dei difensori ( $\tilde{q}_d \leq q_d$ ) nel gioco sequenziale [26, 27]. Questo deriva dall'internalizzazione da parte dei difensori dell'effetto indiretto del proprio sforzo attraverso gli attaccanti, percependo un costo effettivo maggiore per unità di sforzo [27, 28].
- La presenza degli attaccanti induce un **effetto di complementarità strategica** tra i difensori [29, 30]. Questo significa che un aumento dello sforzo di un difensore può portare un altro difensore dipendente ad aumentare il proprio sforzo [29]. In alcuni casi, le dipendenze di sostituzione strategica (alleati) nel gioco simultaneo possono diventare dipendenze di complementarità strategica (avversari) nel gioco sequenziale a causa di questo effetto indiretto [7, 29, 31].
- In determinate condizioni (matrice  $D$  non negativa, effetto di rete diretto tra difensori inferiore alla somma degli effetti di complementarità strategica indiretta attraverso gli attaccanti) [32], i difensori investono **meno** ( $x_d^{SPE} \leq x_d^{NE}$ ) e hanno **minori score di free riding individuale** ( $s^{SPE} \leq s^{NE}$ ) nel gioco sequenziale rispetto a quello simultaneo [31-33]. Ciononostante, possono godere di **utilità maggiori** nel gioco sequenziale [32-34]. Questo suggerisce che i difensori possono sfruttare la presenza degli attaccanti a proprio vantaggio [34].
- Al contrario, i difensori che sono "**free rider**" **più "gravi"** nel gioco simultaneo (cioè, fortemente supportati dagli sforzi altrui) possono essere "**esposti**" nel contesto sequenziale, riducendo la loro capacità di free ride e portandoli a **investire di più** [35-38]. Questo aumento dello sforzo da parte dei free rider esposti può, a sua volta, indurre l'attaccante ad **aumentare il proprio sforzo di attacco** [35, 37, 39, 40].
- Il cambiamento nello sforzo di equilibrio dei difensori è legato alla loro **alpha-centrality** nella rete ridotta [6, 41-43]. Difensori con maggiore alpha-centrality (con  $\alpha = -1$ ) tendono a vedere diminuire maggiormente il loro sforzo nel gioco sequenziale [43, 44].

- Confrontando con l'ottimo sociale, in certe strutture di rete, gli sforzi dei difensori nell'I-SPE del gioco sequenziale sono **superiori** a quelli socialmente ottimali ( $x_d^{SPE} \geq x_d^*$ ) [7, 41, 45, 46], indicando un **eccesso di investimento** (over-investment) [7, 47]. Questo è attribuito all'effetto di complementarità strategica indotto dalla sequenzialità [7, 47].

Il documento riconosce alcune **limitazioni** del modello [48, 49], come il fatto che l'attaccante eserciti un singolo sforzo non mirato che colpisce tutti i difensori in base alla forza della connessione, non catturando attacchi altamente mirati [49, 50]. Inoltre, l'analisi si concentra sugli equilibri interni, non considerando casi in cui alcuni giocatori potrebbero esercitare sforzo zero (azioni di confine o "boundary actions"), il che può rendere le funzioni di utilità non concave e complicare l'analisi [51-55].

In sintesi, il documento analizza un modello combinato di gioco di rete e Stackelberg nel contesto della sicurezza informatica [56], identificando condizioni per l'esistenza di equilibri interni e studiando come la sequenzialità e la struttura della rete influenzino l'investimento in sicurezza, il free riding e le utilità dei giocatori rispetto a scenari a mosse simultanee [56]. I risultati evidenziano effetti complessi, tra cui la riduzione del free riding generale in alcuni casi, ma anche l'esposizione di free rider "gravi" che sono costretti a investire di più, potenzialmente portando a un aumento dello sforzo dell'attaccante [56].

## 4 MITRE ATT&CK

Il framework MITRE ATT&CK è una base di conoscenza universalmente accessibile e costantemente aggiornata per modellare, rilevare, impedire e contrastare le minacce alla cybersecurity sulla base dei comportamenti antagonisti noti dei criminali informatici. MITRE ATT&CK (Adversarial Tactics, Techniques, and Common Knowledge) è un framework di conoscenza basato su osservazioni del mondo reale delle tecniche utilizzate da attori di minacce informatiche. È strutturato in:

- Tattiche: gli obiettivi strategici dell'attaccante (es. iniziare un accesso, mantenere la persistenza, eludere i controlli).
- Tecniche: i metodi specifici utilizzati per raggiungere ciascuna tattica.
- Sottotecniche: variazioni dettagliate delle tecniche.

Ha un approccio modulare (14 tattiche: Reconnaissance, Resource Development, Initial Access, Execution, Persistence, Privilege Escalation, Defense Evasion, Credential Access, Discovery, Lateral Movement, Collection, Command and Control, Exfiltration, Impact), un focus su detection e risposta ed è basato su tecniche reali e dettagliate.

La "kill chain di MITRE" è un modo informale di riferirsi all'intero ciclo di attacco così come descritto nel MITRE ATT&CK framework. È uno strumento essenziale per: analizzare il comportamento degli attaccanti, sviluppare capacità di detection, mappare tecniche note a gruppi APT, costruire difese mirate.

MITRE ATT&CK organizza tattiche e tecniche (e sottotecniche) dell'avversario in matrici. Ciascuna matrice include strategie e tecniche corrispondenti agli attacchi su domini specifici. Nel caso degli IACS è possibile prendere in considerazione la matrice ICS, che include tecniche utilizzate negli attacchi ai sistemi di controllo industriale, in particolare macchinari, dispositivi, sensori e reti

utilizzati per controllare o automatizzare le operazioni per fabbriche, utenze, sistemi di trasporto e altri fornitori di servizi critici.

L'uso di tale framework può essere integrato nell'implementazione del processo di CRA tramite SSG in una fase antecedente al gioco stesso. Infatti, MITRE ATT&CK potrebbe essere applicato per la realizzazione della fase ZCR 5.1 del processo di CRA per IACS, ovvero quella relativa all'individuazione delle minacce. A tale fine, la matrice ICS del MITRE ATT&CK potrebbe essere usata per la delineazione dei profili degli attaccanti e quindi anche delle possibili strategie di difesa iniziali che possano essere applicate. **COLLEGAMENTO DA RIVALUTARE IN SEGUITO ALL'USO DEL GRAFO DI ATTACCO PROBABILISTICO**

MITRE ATT&CK può essere integrato in un SSG per la definizione dello spazio di azione degli attaccanti (es. ogni tecnica MITRE ATT&CK può essere considerata una possibile azione dell'attaccante). Quindi le azioni del difensore costituiscono meccanismi di mitigazione o rilevamento delle tecniche ATT&CK. Integrando tale framework in un SSG si potrebbero ottenere alcuni vantaggi, quali:

- Data-driven: MITRE ATT&CK fornisce una base realistica e strutturata di tecniche usate dagli avversari.
- Decision-making ottimale: SSG permette di ragionare in termini di risorse limitate e attaccanti strategici.
- Adattabile: Può essere personalizzato in base alle minacce rilevanti per la tua organizzazione.

## 5 Terminazione di un SSG

Uno *Stackelberg Security Game (SSG)* è un gioco sequenziale tra due agenti: un **difensore (leader)** e un **attaccante (follower)**. È utilizzato per modellare scenari di sicurezza, come la protezione di infrastrutture critiche, aeroporti o reti informatiche.

Il gioco termina nel momento in cui l'attaccante osserva la strategia del difensore e decide il suo attacco. Tuttavia, il significato di "terminare" può variare a seconda del contesto. Di seguito sono elencati i principali modi in cui può concludersi uno SSG.

### 1. Risoluzione computazionale del gioco

Il gioco termina *matematicamente* quando si trova una strategia ottimale per il difensore, assumendo che l'attaccante risponda razionalmente. Ciò avviene tramite:

- Programmazione lineare o algoritmi di ottimizzazione per calcolare la strategia mista del difensore.
- Utilizzo di algoritmi specifici come DOBSS (Decomposed Optimal Bayesian Stackelberg Solver) in presenza di attaccanti con tipi multipli.

Una volta trovata questa strategia, il gioco è considerato "risolto".

### 2. Esecuzione pratica del gioco

In un'applicazione reale (ad esempio in un aeroporto):

- Il difensore fissa una strategia di pattugliamento (ad esempio, pattugliare certe aree con determinate probabilità).



- L'attaccante osserva e agisce una sola volta (gioco a singolo turno).
- Il gioco termina quando l'attacco avviene e si osserva l'esito.

### 3. Fine simulazione o iterazione

In simulazioni o ambienti dinamici:

- Il gioco può terminare dopo un numero predefinito di turni.
- Oppure quando si raggiunge una *convergenza* delle strategie (nessun miglioramento significativo nei payoff).

Convergenza algoritmica (nei casi dinamici o simulati)

In ambienti iterativi o RL (come se usi Gymnasium) l'equilibrio si raggiunge quando le strategie si stabilizzano. Quindi si possono usare criteri di convergenza tipo:

- Cambiamento medio nei payoff minore di una certa soglia
- Strategie stabili per N iterazioni
- Gradiente della funzione obiettivo vicino a zero

## CONSIDERAZIONI DA VALUTARE IN SEGUITO ALL'USO DI UN GRAFO DI ATTACCO PROBABILISTICO

### 4. Condizioni di arresto in ambienti dinamici

In versioni estese, come *repeated SSGs* o *learning-based SSGs*, il gioco può terminare:

- Quando si raggiunge un certo livello di sicurezza desiderata.
- Quando l'attaccante cambia comportamento in modo non modellabile.
- Per limiti di tempo computazionale o risorse di simulazione.

Uno Stackelberg Security Game si considera terminato quando:

- Viene calcolata una strategia ottima del difensore (in modelli teorici).
- Si compie l'attacco (in scenari reali).
- **Si raggiunge una condizione di arresto definita (in simulazioni).**

L'output di un SSG è una condizione di equilibrio detta SSE, quindi potremmo dichiarare il gioco come finito nel momento in cui questa condizione viene individuata. Tra gli strumenti pratici per verificare il raggiungimento del SSE si può usare Gymnasium? Altrimenti si può usare framework specializzati, come ad esempio, Gambit che calcola e verifica equilibri?

## 6 Esempio SSG

Per la costruzione di un SSG è necessario definirne gli input.

Tema: simulazione del processo di CRA per IACS

Difensore: sistema informatico di un'azienda

Attaccante/i: individuati con MITRE ATT&CK, usando la matrice per ICS

Input:

- Asset condivisi: numero da definire (  $n = 4$  )
- Valori di payoff per attaccante e difensore da definire
- Risorse difensore:  $\lfloor n/2 \rfloor$
- Strategia iniziale difensore: da definire

Terminazione: Si raggiunge una condizione di arresto definita (strategie stabili, ovvero invarianti, per  $N = 10$  iterazioni)

Assets contestualizzati: Engineering Workstation, Operator Work Station, Main Historian Server, Synchronisation Server.

Attaccanti contestualizzati (Matrice ICS):

- Exploitation of Remote Services (T0866): Initial Access e Lateral Movement
- Network Sniffing (T0842): Discovery
- Damage to Property (T0879): Impact
- Change Credential (T0892): Inhibit Response Function

TABELLA DA MODIFICARE A SECONDA DELLE NECESSITÀ

<b>Targets</b>	<b>Difensore</b>		<b>Attaccante 1</b>		<b>Attaccante 2</b>		<b>Attaccante 3</b>		<b>Attaccante 4</b>	
$i$	$u_d^c(\cdot)$	$u_d^u(\cdot)$	$u_a^c(\cdot)$	$u_a^u(\cdot)$	$u_a^c(\cdot)$	$u_a^u(\cdot)$	$u_a^c(\cdot)$	$u_a^u(\cdot)$	$u_a^c(\cdot)$	$u_a^u(\cdot)$
1	0	-1	0	1	0	1	2	5	2	5
2	0	-1	0	1	0	0	9	1	2	5
3	0	-1	0	1	0	0	9	1	2	5
4	0	-1	0	1	0	0	9	1	2	5

Tabella 1: Tabella Payoff SSG

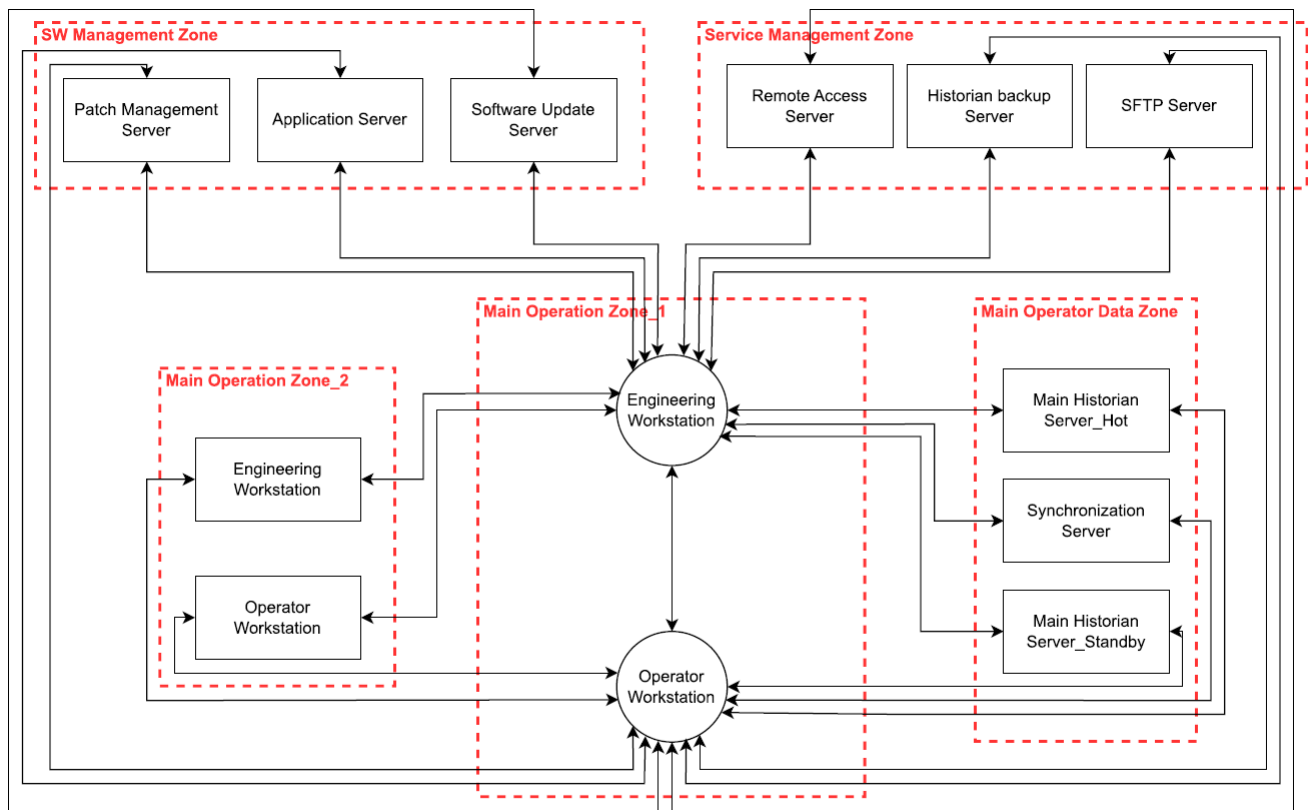
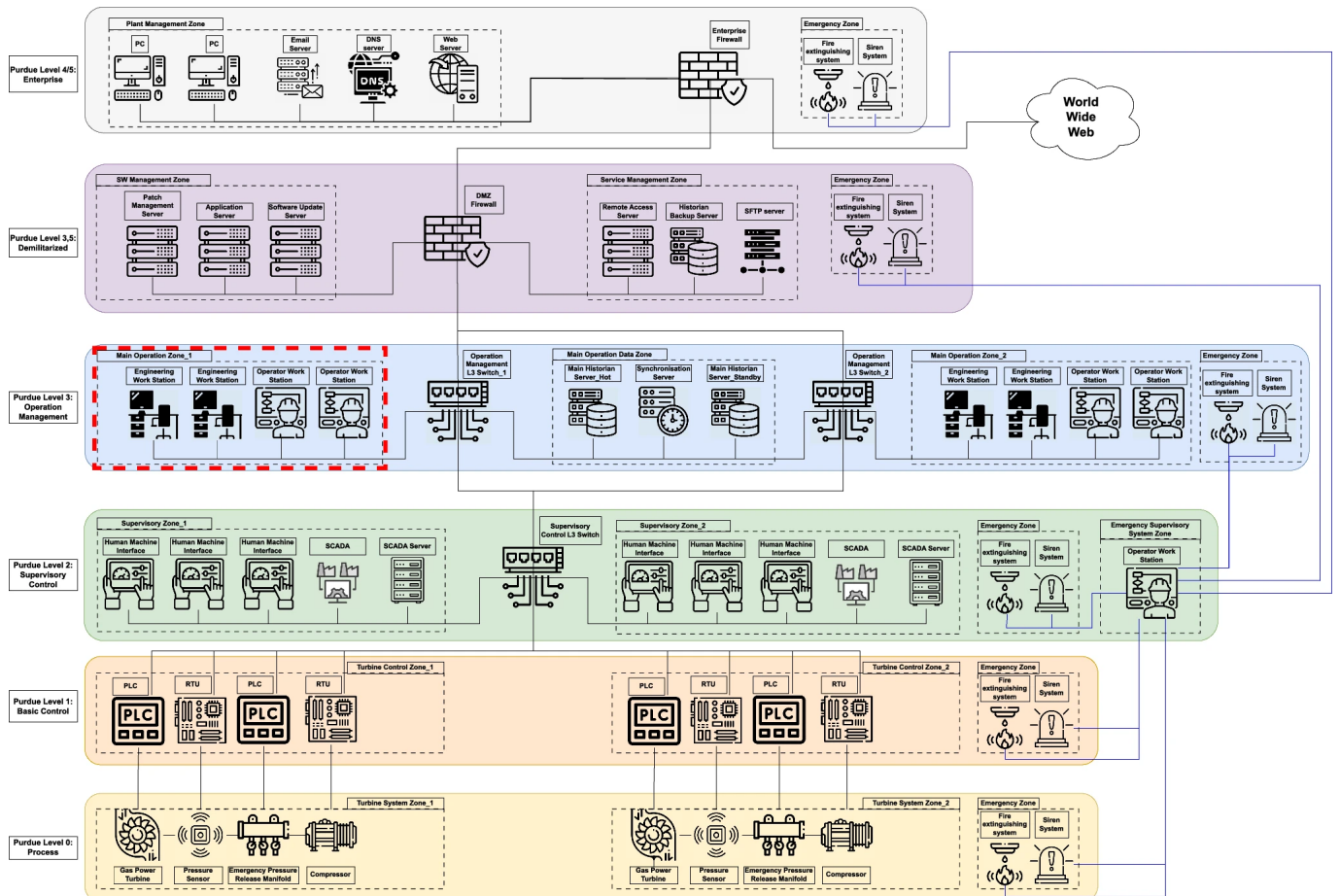


Fig. 3 DFD of the main operation zone 1



SUC per esempio per creazione attack graph (Purdue level 3 e 3.5 del file Brancat et al)

- Patch Management Server (PMS): sistema che automatizza il processo di distribuzione degli aggiornamenti software (patch) ai dispositivi in una rete, inclusi server e workstation.
- Application Server (AS): server che fornisce l'infrastruttura e le funzionalità logiche di supporto, sviluppo ed esecuzione di applicazioni nonché altri componenti server in un contesto distribuito.
- Software Update Server (SUS): server centrale che distribuisce aggiornamenti software (patch, nuove versioni, fix di sicurezza) ai dispositivi di una rete, evitando che ogni dispositivo scarichi gli aggiornamenti direttamente da Internet.
- Remote Access Server (RAS): è un server che consente agli utenti di connettersi da remoto a una rete privata o aziendale tramite una connessione sicura, come una VPN, un desktop remoto o un protocollo di accesso remoto (es. RDP, SSH).
- SFTP Server (SFTPS): un server che utilizza il protocollo SFTP per il trasferimento sicuro di file su una rete.
- DMZ Firewall (F): separa la rete interna di un'organizzazione da una zona demilitarizzata (DMZ), che ospita servizi accessibili da Internet. Questo crea un ulteriore livello di sicurezza, impedendo che eventuali compromissioni nella DMZ si propaghino alla rete interna.
- Synchronization Server (SS): è un sistema o componente di rete che ha il compito di mantenere dati o stati coerenti tra più dispositivi, server o applicazioni. In altre parole, assicura che le informazioni siano uguali e aggiornate su tutti i nodi coinvolti.
- Main Historian Server (MHS): è un tipo speciale di server utilizzato per raccogliere, archiviare, gestire e analizzare dati di processo nel tempo (dati storici), tipicamente in contesti industriali, di automazione o impianti SCADA (Supervisory Control and Data Acquisition).
- Operation Management L3 Switch (S3):
- Operator Workstation (OWS): è una postazione di lavoro (computer) usata dagli operatori di impianto per monitorare e controllare in tempo reale i processi industriali o di automazione
- Engineering Workstation (EWS): è un computer utilizzato in ambienti industriali o di automazione che serve per configurare, gestire e mantenere i sistemi di controllo.

Parti del SUC e relative vulnerabilità (CVSS v4.0 Calculator):

- Patch Management Server: CVE-2024-55660 (CVSS 6.9) (EPSS 0.00092) (V1)
- Application Server: CVE-2025-5140 (CVSS 5.3) (EPSS 0.00038) (V2)
- Software Update Server: CVE-2024-6326 (CVSS 1.8) (EPSS 9e-05) (V3)
- Remote Access Server: CVE-2025-5309 (CVSS 8.6) (EPSS 0.0033) (V4)
- SFTP Server: CVE-2024-52801 (CVSS 5.3) (EPSS 0.00067) (V5)
- DMZ Firewall: CVE-2024-2550 (CVSS 8.7) (EPSS 0.00115) (V6)

- Synchronization Server (SS): CVE-2024-11023 (CVSS 5.2) (EPSS 0.00036) (V7)
- Main Historian Server: CVE-2024-6456 (CVSS 8.5) (EPSS 0.00302) (V8)
- Operation Management L3 Switch: CVE-2023-20048 (CVSS 7.1, v3.01) (EPSS 0.02029) (V9)
- Operator Workstation: CVE-2023-44487 (CVSS 8.7) (EPSS 0.94469) (V10) (Denial of Service:  
CVSS:4.0/AV:N/AC:L/AT:N/PR:N/UI:N/VC:N/VI:N/VA:H/SC:N/SI:N/SA:N),  
CVE-2025-2260 (CVSS 7.1) (EPSS 0.00036) (V11)
- Engineering Workstation: CVE-2025-0327 (CVSS 8.5) (EPSS 0.00025) (V12), CVE-2025-2223 (CVSS 8.4) (EPSS 0.00026) (V13)

## 7 Idea Finale

### Trascrizione PPT

La valutazione del rischio in cybersecurity è essenziale per proteggere proattivamente i sistemi industriali di automazione e controllo (IACS) dai rischi informatici. La teoria dei giochi offre un approccio strutturato per modellare le interazioni tra difensori e attaccanti. Gli **Stackelberg Security Games** (SSG) sono particolarmente adatti per ottimizzare l'allocazione delle risorse difensive.

**Obiettivo:** Introdurre gli SSG, il loro modello, i pro e contro, con un esempio semplice.

### Obiettivo degli Stackelberg Security Games

**Scopo:** Ottimizzare l'allocazione delle risorse del difensore per minimizzare i rischi informatici contro un avversario razionale.

Il modello SSG prevede una dinamica leader-follower:

- Il **difensore** (leader) si impegna per primo in una strategia di sicurezza.
- L'**attaccante** (follower) risponde in modo ottimale, assumendo di conoscere la strategia del difensore.

Ciò è coerente con la necessità della cybersecurity di proteggere risorse critiche con risorse limitate (budget limitato per le contromisure).

**Obiettivo:** Identificare quali contromisure implementare e dove, per minimizzare i rischi informatici, rispettando i vincoli di budget.

### Modello SSG: Input

- **Grafo di attacco:**
  - **Nodi:** stadi dell'attacco (es. shell remota come utente Apache).
  - **Archi:** vulnerabilità da sfruttare per passare al nodo successivo (es. CVE-2022-22720 – Apache HTTP Server 2.4.52).

- Ogni arco ha un numero associato che rappresenta la probabilità che la vulnerabilità venga sfruttata (versione precedente da usare per un confronto nelle considerazioni della tesi: numero associato = rischio non mitigato può essere ottenuto da EPSS \* CVSS).
- **Budget del difensore:** somma di denaro da investire nella difesa (es. 10k).
- **Contromisure:** requisiti di sicurezza con relativi costi ed efficacia sugli archi (es. segmentazione di rete, 1k, riduzione del rischio del 40% su CVE-2022-22720).

## Modello SSG: Iterazioni

1. **Mossa del difensore:** sceglie per primo una strategia selezionando le contromisure da applicare e dove, rispettando i vincoli di budget.
2. **Mossa dell'attaccante:** osserva la strategia del difensore e sceglie il percorso nel grafo che gli offre il miglior ritorno.

Il ritorno dell'attaccante su un percorso corrisponde al rischio del difensore su quel percorso (gioco a somma zero).

Formula del ritorno dell'attaccante per ogni iterazione:

$$\max_{\forall \sigma} \left( \sum_{\forall e \in \sigma} R(e) \right)$$

dove:

- $\sigma$ : percorsi nel grafo
- $e$ : archi
- $R(e)$ : rischio dell'arco

Il difensore aggiusta la propria strategia per ridurre il rischio dell'iterazione precedente. Le iterazioni continuano fino al raggiungimento dell'equilibrio.

## Modello SSG: Output

L'equilibrio è raggiunto quando qualsiasi strategia difensiva diversa da quella corrente darebbe all'attaccante un ritorno maggiore rispetto all'iterazione precedente.

In altre parole, il gioco continua finché il difensore non trova una strategia tale che il ritorno per l'attaccante sia minore o uguale a quello dell'iterazione precedente.

**Strategia del difensore:** la migliore allocazione delle risorse che minimizza i rischi lungo i percorsi, rispettando i vincoli di budget.

**Strategia dell'attaccante:** il percorso che alla fine gli offre il miglior ritorno.

Se il ritorno finale dell'attaccante non è considerato un rischio tollerabile (è necessario fissare una soglia), il difensore deve aumentare il budget da investire nella difesa.

## Considerazioni Post-Riunione (uso di CVSS ed EPSS obsoleto)

CVE → Il Common Vulnerabilities and Exposures, o CVE, è un dizionario di vulnerabilità e falle di sicurezza note pubblicamente mantenuto dalla MITRE Corporation.

CVSS → Il Common Vulnerability Scoring System (CVSS) è una norma tecnica aperta per valutare la gravità delle vulnerabilità di sicurezza di un sistema informatico. CVSS assegna un punteggio di gravità alle vulnerabilità, consentendo a chi si occupa di rispondere all'emergenza di stabilire la priorità di risposte e risorse in base al livello di minaccia. I punteggi vengono calcolati con una formula che dipende da diverse metriche che approssimano la facilità e l'impatto di un exploit. Il punteggio è espresso in una scala da 0 a 10, dove 10 indica il livello di vulnerabilità più grave.

EPSS → L'Exploit Prediction Scoring System (EPSS) è uno strumento basato sui dati per stimare la probabilità che una vulnerabilità software venga sfruttata in natura. Sebbene altri standard di settore siano stati utili per catturare le caratteristiche innate di una vulnerabilità e fornire misure di gravità, la loro capacità di valutare la minaccia è limitata. L'EPSS colma questa lacuna poiché utilizza le informazioni attuali sulle minacce provenienti da CVE e dati di exploit reali. Il modello EPSS produce un punteggio di probabilità compreso tra 0 e 1 (0 e 100%). Maggiore è il punteggio, maggiore è la probabilità che una vulnerabilità venga sfruttata.

Ciò su cui ci si concentra, in relazione al CRA per IACS, è la fase di sviluppo di un "Preventive Security Plan". A tale proposito si sviluppa il grafo di attacco probabilistico.

valore arco attack graph (rischio di attacco) = CVSS score (impatto) \* EPSS score (probabilità dell'attacco) → payoff attaccante (opposto a quello del difensore) → da mappare a livello qualitativo

Rischio raggiungimento obiettivo finale da parte dell'attaccante = valore arco attack graph maggiore nel path intrapreso dall'attaccante stesso

Nel grafo di attacco probabilistico c'è la possibilità di avere più nodi iniziali così come più nodi finali, ovvero possono essere presenti più punti di accesso al sistema, così come più obiettivi per vari attaccanti. I primi all'interno del grafo possono essere rappresentati come nodi con soli archi uscenti, mentre i secondi come nodi con soli archi entranti. Quest'ultimo fatto non vuole dire che da quell'asset non ci siano più collegamenti ad altre componenti del sistema, ma che tutti gli altri possibili collegamenti sono "effimeri" per l'attaccante. Il fatto che i nodi finali siano contraddistinti dai soli archi entranti è dovuto al fatto che, nel momento in cui un attaccante li raggiungesse, esso terminerebbe la propria penetrazione all'interno del sistema. Tale condizione, però, varia a seconda del profilo dell'attaccante stesso.

First → sito per CVSS e EPSS

Rimane da fare un'analisi a risorse infinite

## Nuove Considerazioni Post-Riunione (19/06/2025)

Considerare di fare 2 esempi, uno con il SUC del paper di Brancati et al e un altro basandosi sul sistema di esempio illustrato nel paper di ADVISE.

STRIDE GPT → Attack Tree, Per Attack graph usare NetworkX

Nuove considerazioni sui pesi degli archi del grafo: peso arco = probabilità che la vulnerabilità rappresentata dall'arco venga sfruttata da un attaccante.

Le vecchie considerazioni (peso arco = CVSS \* EPSS) possono essere usate per un confronto nelle considerazioni della tesi.

Analisi delle tipologie di attaccanti sfruttando il modello STRIDE e la tabella della TAL (La TAL aiuta rapidamente i risk manager a identificare con precisione e a comprendere l'importanza degli agenti di minaccia rilevanti. La libreria è composta da 22 archetipi standardizzati definiti utilizzando otto attributi comuni; gli archetipi rappresentano agenti di minaccia esterni e interni, che vanno dalle spie industriali ai dipendenti non formati).

Per ogni arco del grafo, considerando una tipologia di attaccante, si segue un determinato processo di livello qualitativo:

- Si parte dalla probabilità di attraversamento dell'arco ( = sfruttamento della vulnerabilità) più alta
- Si chiede all'owner del sistema se per quella vulnerabilità sono già state applicate delle contromisure di quelle disponibili
- A seconda delle contromisure già implementate si va a diminuire di conseguenza la probabilità di sfruttamento della vulnerabilità.

Le probabilità possono essere mappate a piacere dal livello qualitativo a quello quantitativo.

Le tipologie di attaccanti dipendono dal nodo dal quale l'arco parte, e quindi da quale struttura all'interno del sistema informatico è rappresentata dal nodo stesso.

Così facendo, si può evitare di usare per gli archi delle vulnerabilità troppo specifiche.