

Governare l'autonomia degli agenti nella Cyber Security

Federico Cerutti



RICHMOND
**CYBER RESILIENCE
FORUM**

RIMINI
10-12 MAGGIO 2026
SPRING EDITION



VOTA LA CONFERENZA



Image by ChatGPT



Contesto



Memoria



Strumenti



Permessi



Azione

Images from Pexel

Investigations

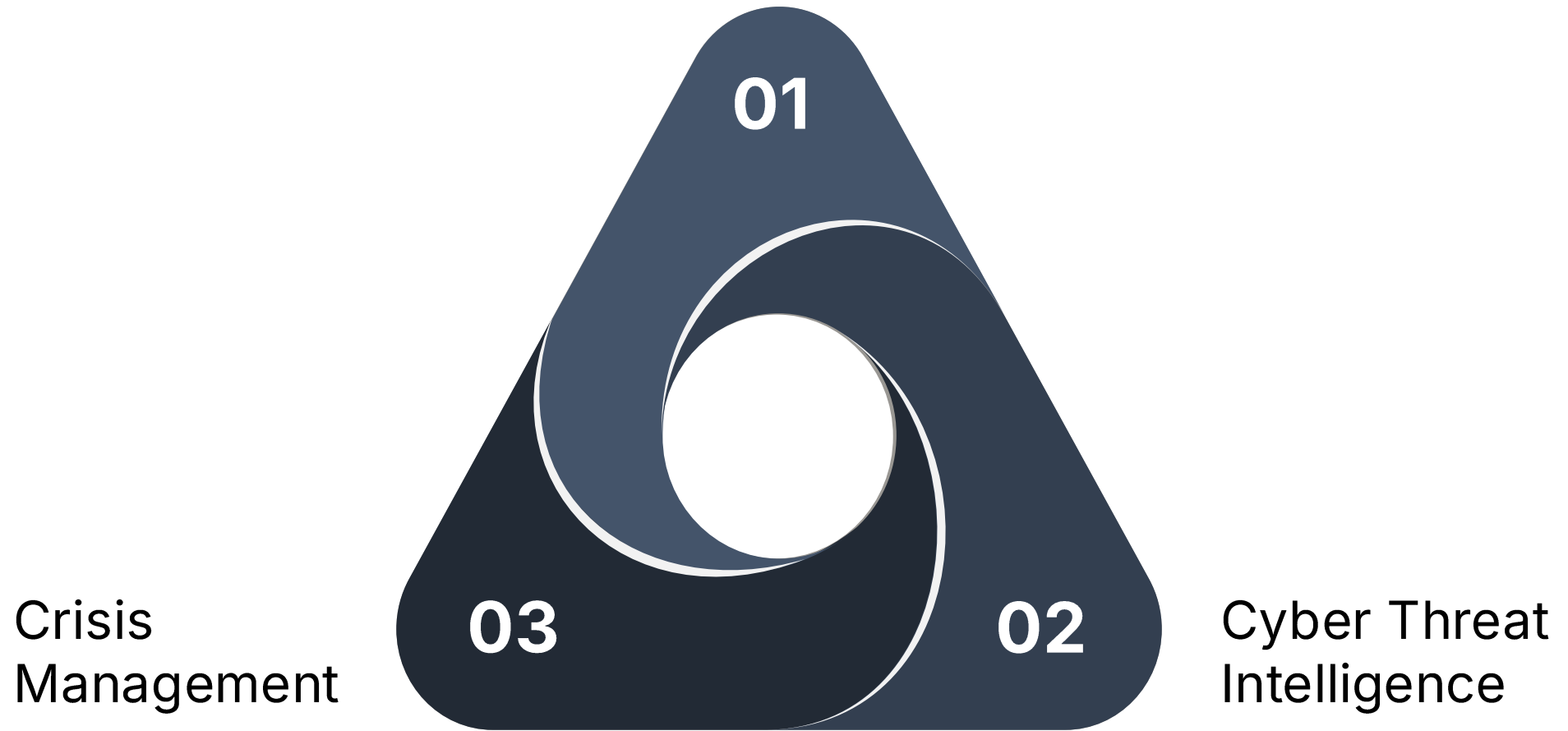




Image from Pexel



Image from Pexels



RICHMOND
CYBER RESILIENCE
FORUM

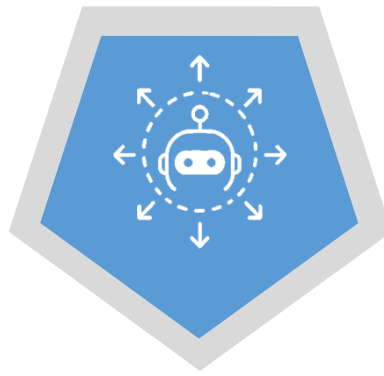


Image by Gemini

Il consiglio
sbagliato
si discute

L'azione
sbagliata
si paga

Perdita del controllo



Compromissione



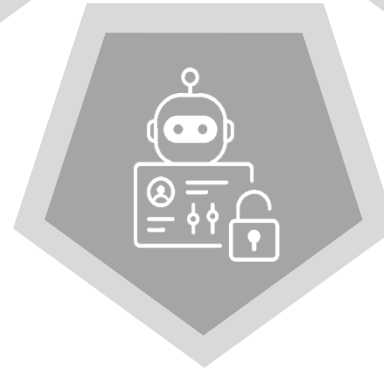
Goal drift



Manipolazione della memoria



Abuso delle autorizzazioni

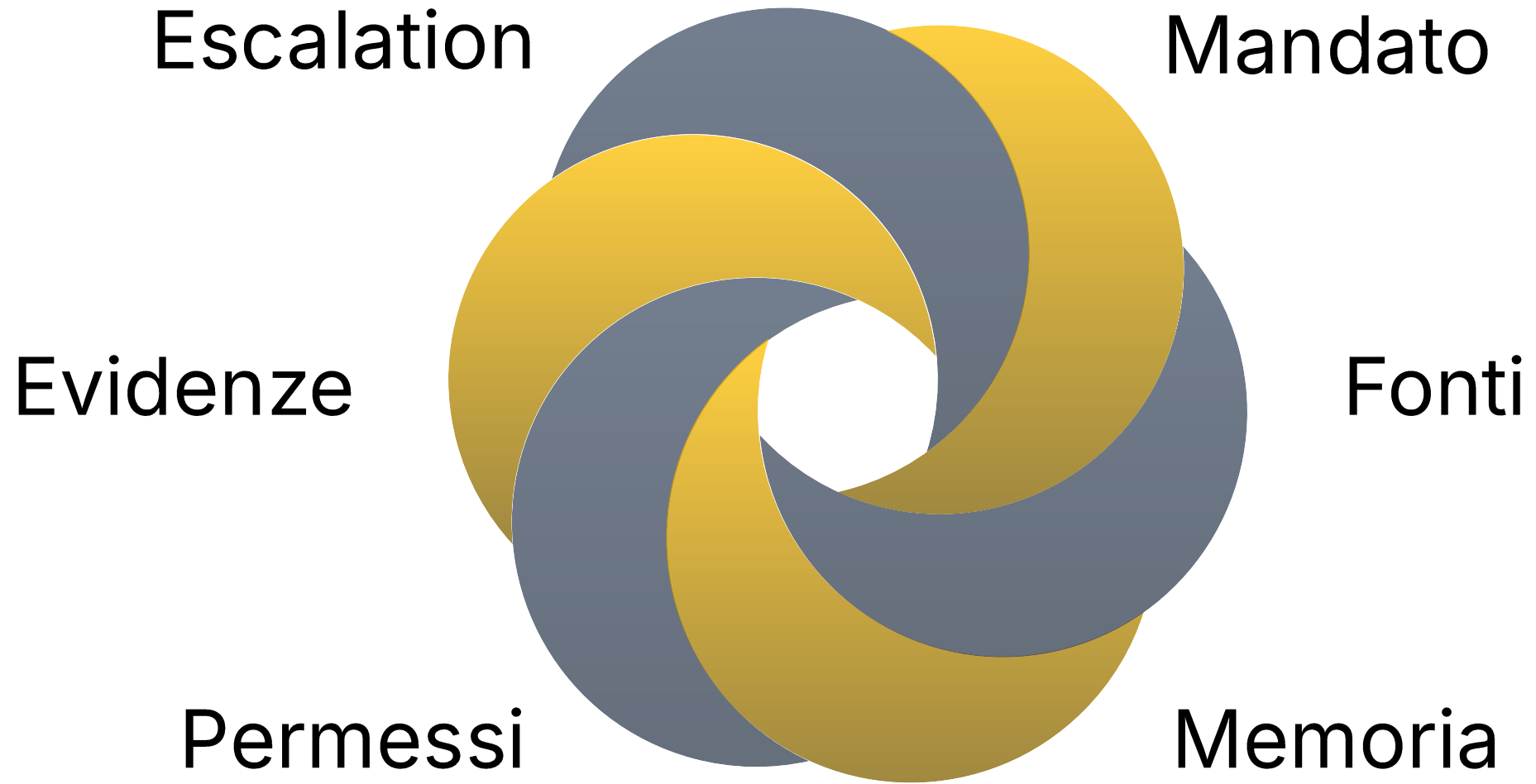


Icons by ChatGPT



La
conoscenza
decide
l'azione

Image by ChatGPT with a concept from Pexel



Non è
un *tool*
da usare

È
firmare
una delega

Il valore dell'agente è
nella conoscenza
che lo guida.

Grazie

RICHMONDITALIA
HUMAN2HUMANEVENTS



VOTA LA CONFERENZA