

# Agents IA et Raisonnement Causal pour la Décision 2026

Catégorie : Intelligence Artificielle | Lecture : 12 min | Publié le : 16/02/2026 | Auteur : Ayi NEDJIMI

Guide expert sur le raisonnement causal dans les agents IA : échelle de Pearl, graphes causaux, DAGs, modèles SCM, intégration LLM, applications.

---

## Table des Matières

1. 1. Introduction : Corrélation vs Causalité dans l'IA
2. 2. L'Échelle Causale de Pearl : Association, Intervention, Contrefactuels
3. 3. Graphes Causaux (DAGs) et Modèles SCM
4. 4. Intégration avec les Agents LLM : Prompting et Neuro-Symbolique
5. 5. Applications : Stratégie, Diagnostic Médical, Root Cause Analysis
6. 6. Méthodes de Découverte Causale : Constraint-Based et Score-Based
7. 7. Scénarios Contrefactuels "What-If" pour les Agents
8. 8. Limitations du Raisonnement Causal et Stratégies de Mitigation
9. 9. Benchmarks et Tendances Futures

Votre organisation est-elle prête à faire face aux attaques basées sur l'IA ?

## 1 Introduction : Corrélation vs Causalité dans l'IA

Les systèmes d'intelligence artificielle traditionnels, notamment les modèles de machine learning et les grands modèles de langage (LLM), excellent dans la détection de **corrélations** à partir de données massives. Cependant, la simple corrélation ne permet pas de répondre aux questions fondamentales de la décision stratégique : *Que se passerait-il si nous modifions cette variable ?* ou *Pourquoi cet événement s'est-il produit ?*

Le **raisonnement causal** représente un saut qualitatif majeur pour les agents IA. Contrairement aux approches purement statistiques qui observent des associations (X et Y varient ensemble), le raisonnement causal permet de modéliser des relations de **cause à effet** (X influence Y) et d'explorer des scénarios contrefactuels (si X avait été différent, Y aurait changé comment ?).

Cette distinction est cruciale dans des domaines comme la stratégie d'entreprise, le diagnostic médical, l'analyse financière ou la maintenance prédictive, où les décideurs doivent comprendre non seulement *ce qui s'est passé*, mais surtout *pourquoi* et *ce qui se passerait dans des conditions différentes*.

**Point clé :** Un agent IA équipé de raisonnement causal peut passer d'une simple prédiction statistique ("Il y a 80% de probabilité que Y augmente") à une explication actionnable ("Si nous réduisons X de 10%, Y diminuera de 15% parce que X cause directement Y via le mécanisme Z").

Critere	Description	Niveau de risque
Confidentialite	Protection des donnees d'entrainement et des prompts	Eleve
Integrite	Fiabilite des sorties et detection des hallucinations	Critique
Disponibilite	Resilience du service et gestion de la charge	Moyen
Conformite	Respect du RGPD, AI Act et politiques internes	Eleve

### Notre avis d'expert

Chez Ayi NEDJIMI Consultants, nous constatons que la majorité des organisations sous-estiment les risques liés aux modèles de langage déployés en production. La sécurité des LLM ne se limite pas au prompt engineering : elle exige une approche systémique couvrant les embeddings, les pipelines de données et les mécanismes de contrôle d'accès aux API.

## 2 L'Échelle Causale de Pearl : Les Trois Niveaux de Raisonnement

Judea Pearl, pionnier de l'inférence causale et lauréat du prix Turing, a formalisé une hiérarchie du raisonnement causal en trois niveaux distincts, chacun offrant des capacités croissantes pour les agents IA.

### Niveau 1 : Association (Seeing)

Le premier niveau concerne l'observation passive de données et la détection de **patterns statistiques**. Les questions typiques sont du type : "*Quelle est la probabilité de Y sachant X ?*" ( $P(Y|X)$ ).

C'est le domaine des modèles de machine learning classiques : régression, classification, clustering. Un agent à ce niveau peut identifier que les clients qui achètent A ont tendance à acheter B, mais ne peut pas affirmer que A *cause* l'achat de B.

## Niveau 2 : Intervention (Doing)

Le deuxième niveau introduit la notion d'**intervention active**. Les questions deviennent : "Que se passerait-il si nous fixions  $X$  à une certaine valeur ?" ( $P(Y | do(X=x))$ ).

L'opérateur  $do(\cdot)$  est fondamental : il représente une manipulation causale du système, pas une simple observation conditionnelle. Par exemple, "Si nous augmentons le prix de 5% ( $do(prix=1.05)$ ), comment les ventes vont-elles réagir ?" Cette question ne peut être résolue par des corrélations passives si le prix n'a jamais été testé à ce niveau.

## Niveau 3 : Contrefactuels (Imagining)

Le troisième niveau, le plus élaboré, permet de raisonner sur des **scénarios alternatifs** : "Si  $X$  avait été différent dans le passé,  $Y$  aurait-il changé ?" Ce type de raisonnement rétrospectif est essentiel pour comprendre les **causes racines** d'événements passés.

Exemple : "Si notre campagne marketing avait été lancée une semaine plus tôt, aurions-nous évité la baisse des ventes de Q3 ?" Cette question contrefactuelle nécessite un modèle causal complet du système business, incluant les mécanismes temporels et les confondants potentiels. Pour approfondir, consultez [Développement Intelligence Artificielle](#) |.

## 3 Graphes Causaux (DAGs) et Modèles Causaux Structurels (SCM)

Les **graphes causaux**, formellement appelés **Directed Acyclic Graphs (DAGs)**, constituent le langage mathématique fondamental pour représenter les relations causales dans un système. Dans un DAG, les nœuds représentent des variables et les arêtes orientées représentent des relations de causalité directe.

### Propriétés des DAGs Causaux

Un DAG causal respecte plusieurs propriétés fondamentales :

- **Directionnalité** : Les flèches indiquent la direction de la causalité ( $X \rightarrow Y$  signifie "X cause Y")
- **Acyclicité** : Pas de boucles causales (X ne peut pas causer Y qui cause Z qui cause X)
- **d-séparation** : Critère graphique pour déterminer l'indépendance conditionnelle entre variables
- **Colliders et confondants** : Structures spécifiques (fork, chain, collider) qui influencent l'inférence

### Modèles Causaux Structurels (SCM)

Un **Structural Causal Model (SCM)** enrichit le DAG en associant à chaque variable une **équation structurelle** qui décrit comment elle est générée à partir de ses parents causaux et d'un terme de bruit exogène.

Formellement, un SCM est défini par :

- Un ensemble de variables endogènes  $V = \{V_1, V_2, \dots, V_n\}$
- Un ensemble de variables exogènes  $U = \{U_1, U_2, \dots, U_m\}$  (non observées)
- Pour chaque  $V_i$ , une fonction structurelle :  $V_i = f_i(PA_i, U_i)$  où  $PA_i$  sont les parents de  $V_i$

Exemple concret en stratégie e-commerce :

```
Budget_Marketing = U_budget (exogène)
Trafic_Site = f_1(Budget_Marketing, U_trafic)
Taux_Conversion = f_2(UX_Design, U_conversion)
Revenus = f_3(Trafic_Site, Taux_Conversion, U_revenus)
```

Ce modèle permet de répondre à des questions comme : "Si nous augmentons le budget marketing de 20%, comment les revenus vont-ils évoluer, sachant que le taux de conversion dépend indépendamment de l'UX design ?"

### Cas concret

En février 2024, une entreprise de Hong Kong a perdu 25 millions de dollars après qu'un employé a été trompé par un deepfake vidéo lors d'une visioconférence. Les attaquants avaient recréé l'apparence et la voix du directeur financier à l'aide de modèles d'IA générative, démontrant les risques concrets de cette technologie en contexte corporate.

Comment garantir que vos modèles de machine learning ne deviennent pas des vecteurs d'attaque ?

## 4 Intégration avec les Agents LLM : Prompting Causal et Approches Neuro-Symboliques

L'intégration du raisonnement causal dans les agents basés sur des LLM représente un défi fascinant. Les LLM, par nature, sont des systèmes d'association statistique (niveau 1 de Pearl) qui capturent des corrélations dans leurs données d'entraînement. Comment les doter de capacités causales de niveau 2 et 3 ?

### Approche 1 : Causal Prompting

Une première approche consiste à utiliser des **techniques de prompting spécialisées** pour guider le LLM vers un raisonnement causal. Cela implique :

- **Questions causales explicites** : "Quelle est la CAUSE de X ?" plutôt que "X et Y sont-ils corrélés ?"
- **Chain-of-Thought causale** : Forcer le modèle à expliciter les mécanismes causaux étape par étape
- **Few-shot avec exemples causaux** : Fournir des exemples annotés de raisonnement causal correct

Exemple de prompt causal pour un agent LLM :

Tu es un expert en analyse causale. Analyse le système suivant :

- Budget marketing : 100k€
- Trafic site : 50k visiteurs
- Taux conversion : 2%
- Revenus : 200k€

Question : Si nous augmentons le budget marketing de 20%, quel sera l'impact sur les revenus ? Réponds en suivant ces étapes :

1. Identifie le graphe causal (quelles variables causent quelles autres)
2. Identifie les variables confondantes potentielles
3. Applique l'opérateur  $do(\text{budget\_marketing} = 120k)$
4. Calcule l'effet causal total sur les revenus
5. Indique les hypothèses et incertitudes

## Approche 2 : Architecture Neuro-Symbolique Hybride

Une approche plus robuste consiste à combiner la puissance des LLM avec des **moteurs d'inférence causale symboliques**. L'architecture typique :

- **1. LLM pour l'extraction de structure** : Le LLM analyse le contexte et propose un DAG causal initial basé sur sa connaissance du domaine
- **2. Moteur causal pour l'inférence** : Une bibliothèque comme DoWhy ou CausalNex effectue les calculs d'inférence causale rigoureux
- **3. LLM pour l'interprétation** : Le LLM traduit les résultats techniques en explications naturelles et recommandations actionnables

Voici un exemple d'implémentation avec **DoWhy** (bibliothèque Python développée par Microsoft Research) : Pour approfondir, consultez [Apprentissage Fédéré et Privacy-Preserving ML en Cybersécurité](#).

```

import dowhy
from dowhy import CausalModel
import pandas as pd
import numpy as np

# 1. Données observationnelles (historiques e-commerce)
data = pd.DataFrame({
    'budget_marketing': np.random.uniform(50, 150, 1000),
    'design_ux_score': np.random.uniform(1, 10, 1000),
    'saison': np.random.choice(['ete', 'hiver', 'printemps', 'automne'], 1000),
})

# 2. Définition du graphe causal (DAG)
model = CausalModel(
    data=data,
    treatment='budget_marketing',
    outcome='revenus',
    graph="""
    digraph {
        budget_marketing -> trafic;
        design_ux_score -> taux_conversion;
        saison -> budget_marketing;
        saison -> taux_conversion;
        trafic -> revenus;
        taux_conversion -> revenus;
    }
    """
)

# 3. Identification de l'effet causal (do-calculus)
identified_estimand = model.identify_effect(proceed_when_unidentifiable=True)

# 4. Estimation de l'effet causal
estimate = model.estimate_effect(
    identified_estimand,
    method_name="backdoor.linear_regression",
    control_value=100, # Baseline budget
    treatment_value=120 # Intervention: +20% budget
)
print(f"Effet causal estimé : {estimate.value} € de revenus supplémentaires")

# 5. Validation par réfutation (tests de robustesse)
refute = model.refute_estimate(
    identified_estimand,
    estimate,
    method_name="random_common_cause"
)

```

Ce code démontre les étapes clés : définition du DAG, identification de l'effet causal via le backdoor criterion, estimation quantitative de l'intervention, et validation par réfutation pour tester la robustesse des hypothèses.

## 5 Applications Pratiques : Stratégie, Diagnostic Médical, Root Cause Analysis

### Stratégie Business et Optimisation Marketing

Dans le domaine du marketing digital, les agents causaux permettent de dépasser les limites de l'attribution multi-touch traditionnelle. Au lieu de simples corrélations entre canaux et conversions, un agent causal peut :

- Identifier les canaux qui *causent réellement* des conversions vs ceux corrélés à des conversions
- Simuler des interventions budgétaires (do-calculus) avant de les déployer en production
- Détecter les effets de synergie causale entre canaux (ex: TV + Social cause un lift supérieur à leur somme)

### Diagnostic Médical et Aide à la Décision Clinique

Le raisonnement causal est fondamental en médecine. Un agent IA médical équipé de capacités causales peut :

- Distinguer les symptômes qui sont des **causes** d'une maladie vs de simples **comorbidités**
- Prédire l'effet d'un traitement (intervention) sur un patient spécifique, en tenant compte des confondants (âge, comorbidités, génétique)
- Raisonnement contrefactuel : "Si ce patient avait reçu le traitement A plutôt que B, son pronostic aurait-il été meilleur ?"

Exemple : Un agent analyse un patient diabétique avec hypertension. Le graphe causal révèle que l'hypertension est partiellement *causée* par le diabète (via l'inflammation vasculaire), mais aussi influencée par l'âge et le mode de vie. L'agent peut alors recommander un traitement ciblant la cause racine (contrôle glycémique) plutôt que seulement les symptômes (antihypertenseurs).

### Root Cause Analysis en Maintenance et Production

Dans les systèmes industriels complexes (usines, datacenters, infrastructures IT), identifier la **cause racine** d'une défaillance est crucial pour éviter les récurrences. Un agent causal peut :

- Construire un DAG des dépendances système (composant A alimente composant B qui contrôle C)
- Lors d'une panne, remonter le graphe causal pour identifier le nœud source de la cascade de défaillances
- Raisonnement contrefactuel : "Si le capteur X avait été remplacé avant sa durée de vie maximale, la panne aurait-elle été évitée ?"

## 6 Méthodes de Découverte Causale : Constraint-Based et Score-Based

Dans les sections précédentes, nous avons supposé que le graphe causal (DAG) était connu ou spécifié par un expert. Mais que faire lorsque nous disposons uniquement de données observationnelles, sans connaissance a priori de la structure causale ? C'est le domaine de la **découverte causale automatique**.

### Approches Constraint-Based (PC, FCI)

Les algorithmes constraint-based, comme **PC (Peter-Clark)** et **FCI (Fast Causal Inference)**, exploitent les tests d'indépendance conditionnelle pour inférer la structure du DAG.

Principe : Si  $X$  et  $Y$  sont indépendants conditionnellement à  $Z$ , alors  $Z$  est un parent commun ou un collider. L'algorithme teste systématiquement toutes les combinaisons pour éliminer les arêtes incompatibles avec les données.

### Approches Score-Based (GES, NOTEARS)

Les méthodes score-based, comme **GES (Greedy Equivalence Search)** ou **NOTEARS** (plus récent, 2018), formulent la découverte causale comme un problème d'optimisation : trouver le DAG qui maximise un score (ex: BIC, likelihood) tout en respectant la contrainte d'acyclicité.

NOTEARS est particulièrement innovant : il reformule la contrainte d'acyclicité en une contrainte continue différentiable, permettant l'utilisation de gradient descent pour optimiser le graphe.

Exemple avec **CausalNex** (bibliothèque de découverte causale par QuantumBlack/McKinsey) : Pour approfondir, consultez [AI Worms et Propagation Autonome : Menaces Émergentes 2026](#).

```

from causalnex.structure.notears import from_pandas
from causalnex.network import BayesianNetwork
import pandas as pd

# Données observationnelles (sans connaissance du DAG)
data = pd.DataFrame({
    'trafic': [100, 150, 200, 120, 180],
    'budget': [50, 75, 100, 60, 90],
    'conversion': [0.02, 0.025, 0.03, 0.022, 0.028],
    'revenus': [200, 375, 600, 264, 504]
})

# 1. Découverte automatique du DAG via NOTEARS
sm = from_pandas(data, w_threshold=0.3)
print("DAG découvert automatiquement :", sm.edges())

# 2. Apprentissage des probabilités conditionnelles
bn = BayesianNetwork(sm)
bn.fit_node_states(data)
bn.fit_cpds(data, method="BayesianEstimator", bayes_prior="K2")

# 3. Inférence : effet d'une intervention sur le budget
from causalnex.inference import InferenceEngine
ie = InferenceEngine(bn)
marginals_baseline = ie.query()['revenus']
marginals_intervention = ie.query(do={'budget': 150})['revenus']

print(f"Effet causal moyen : {marginals_intervention.mean() -
marginals_baseline.mean()} €")

```

## 7 Scénarios Contrefactuels "What-If" pour les Agents

Le raisonnement contrefactuel (niveau 3 de Pearl) est sans doute la capacité la plus complexe et la plus utile pour les agents décisionnels. Il permet de répondre à des questions du type : *"Étant donné ce qui s'est passé, que se serait-il passé si nous avions agi différemment ?"*

### Formalisation des Contrefactuels

Mathématiquement, un contrefactuel s'écrit :  $P(Y_x = y \mid X' = x', Y' = y')$ , qui se lit : "Quelle serait la probabilité que Y prenne la valeur y si X avait été fixé à x, sachant que dans la réalité observée X a pris la valeur x' et Y a pris la valeur y' ?"

Le calcul contrefactuel nécessite trois étapes (algorithme de Pearl) :

- 1. **Abduction** : Inférer les valeurs des variables exogènes U à partir des observations (X', Y')
- 2. **Action** : Modifier le modèle en fixant  $X = x$  (intervention  $do(X=x)$ )
- 3. **Prédiction** : Calculer  $Y_x$  en utilisant les valeurs des U inférées à l'étape 1

## Agents Autonomes et Apprentissage Contrefactuel

Les agents IA peuvent utiliser le raisonnement contrefactuel pour **l'apprentissage par renforcement off-policy**. Au lieu d'explorer aléatoirement l'espace d'actions (coûteux et risqué), l'agent peut :

- Analyser les trajectoires passées et générer des contrefactuels : "Si j'avais choisi l'action  $A_2$  au lieu de  $A_1$ , quel aurait été le résultat ?"
- Apprendre des regrets causaux : améliorer la politique en identifiant les décisions sous-optimales *causalement*
- Sécurité : tester des actions potentiellement risquées en simulation contrefactuelle avant déploiement réel

## 8 Limitations du Raisonnement Causal et Stratégies de Mitigation

Malgré sa puissance, le raisonnement causal comporte des limitations importantes que tout praticien doit connaître.

### Limitation 1 : Hypothèses Non Testables

De nombreuses hypothèses causales sont **non testables empiriquement** avec des données observationnelles seules. Par exemple, l'hypothèse "il n'existe pas de confondant non observé" ne peut jamais être prouvée avec certitude.

**Mitigation** : Utiliser des analyses de sensibilité pour quantifier comment les conclusions changeraient si les hypothèses étaient violées. DoWhy offre des méthodes de réfutation (placebo treatment, random common cause) pour tester la robustesse.

### Limitation 2 : Équivalence de Markov

Plusieurs DAGs différents peuvent générer les mêmes distributions de probabilité observables (classe d'équivalence de Markov). Les données seules ne permettent pas toujours de distinguer  $X \rightarrow Y$  de  $Y \rightarrow X$ .

**Mitigation** : Intégrer de la connaissance du domaine (contraintes temporelles, impossibilités physiques) pour éliminer les DAGs incompatibles. Utiliser des expériences randomisées contrôlées (A/B tests) quand possible pour casser l'équivalence.

### Limitation 3 : Complexité Computationnelle

L'apprentissage de DAGs est NP-hard. Pour des systèmes avec des dizaines ou centaines de variables, la recherche exhaustive devient impraticable.

**Mitigation** : Utiliser des approches hiérarchiques (découper le système en modules causaux indépendants), des algorithmes d'approximation (NOTEARS, gradient-based), ou des contraintes de sparsité (imposer un nombre maximal de parents par nœud). Pour approfondir, consultez [Mixture of Experts : Architecture LLM de 2026](#).

## 9 Benchmarks, Évaluation et Tendances Futures

### Benchmarks de Raisonnement Causal

L'évaluation des capacités causales des agents IA reste un défi. Plusieurs benchmarks récents émergent :

- **Causalbench (2023)** : Benchmark de découverte causale sur des données biologiques (réseaux de gènes)
- **CLADDER (2024)** : Dataset de questions causales en langage naturel pour évaluer les LLM
- **CausalWorld (RL)** : Environnement de simulation pour agents RL avec structure causale explicite

### Tendances et Directions de Recherche 2026

Les développements récents et à venir incluent :

- **LLM causaux natifs** : Modèles pré-entraînés avec objectifs causaux (causal language modeling) plutôt que seulement prédictifs
- **Causal world models** : Agents qui apprennent des représentations causales de leur environnement pour généraliser à des contextes non observés
- **Causalité temporelle** : Extension aux séries temporelles et systèmes dynamiques (causal inference sur les DAGs temporels)
- **Fairness causale** : Utilisation de graphes causaux pour définir et garantir l'équité des décisions IA (éliminer les discriminations causales)

### Défis Organisationnels et Adoption en Entreprise

L'adoption du raisonnement causal en entreprise nécessite de surmonter plusieurs barrières :

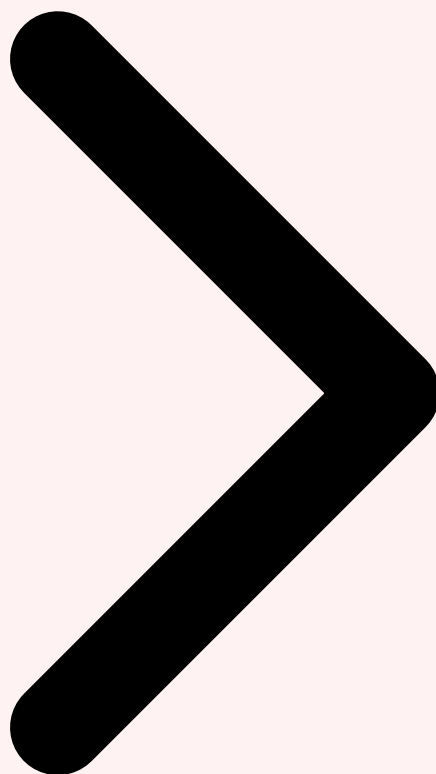
- **Formation** : Les équipes data doivent acquérir des compétences en inférence causale, au-delà du ML classique
- **Collaboration domaine-data** : La construction de DAGs nécessite l'expertise métier, pas seulement algorithmique
- **Infrastructure** : Intégration des outils causaux (DoWhy, CausalNex) dans les pipelines MLOps existants

**Conclusion** : Le raisonnement causal représente un saut qualitatif majeur pour les agents IA. En passant de la simple détection de patterns (niveau 1) à la capacité d'interventions (niveau 2) et de contrefactuels (niveau 3), les agents deviennent de véritables partenaires de décision stratégique.

Les entreprises qui maîtrisent cette transition — en combinant la puissance des LLM avec des moteurs d'inférence causale robustes — disposeront d'un avantage compétitif durable pour naviguer dans des environnements complexes et incertains.



[Retour au sommaire](#) [Raisonnement Causal IA](#) [Tous les articles IA](#)



### **Besoin d'un accompagnement expert en IA causale ?**

Nos consultants vous accompagnent dans l'intégration du raisonnement causal dans vos systèmes IA et agents autonomes. Devis personnalisé sous 24h.

### **Références et ressources externes**

- OWASP LLM Top 10 — Les 10 risques majeurs pour les applications LLM
- MITRE ATLAS — Framework de menaces pour les systèmes d'intelligence artificielle
- NIST AI RMF — AI Risk Management Framework du NIST
- arXiv — Archive ouverte de publications scientifiques en IA
- HuggingFace Docs — Documentation de référence pour les modèles de ML

Pour approfondir ce sujet, consultez notre outil open-source ai-threat-detection qui facilite la détection de menaces basée sur l'IA.

**Sources et références :** [ArXiv IA](#) · [Hugging Face Papers](#)


## FAQ

---

### Qu'est-ce que Agents IA et Raisonnement Causal pour la Décision 2026 ?

Le concept de Agents IA et Raisonnement Causal pour la Décision 2026 est détaillé dans les premières sections de cet article, qui couvrent les fondamentaux, les enjeux et le contexte opérationnel. Pour un accompagnement sur ce sujet, [contactez nos experts](#).

### Pourquoi Agents IA et Raisonnement Causal pour la Décision 2026 est-il important en cybersécurité ?

La compréhension de Agents IA et Raisonnement Causal pour la Décision 2026 permet aux équipes de sécurité d'améliorer leur posture défensive. Les sections «  Table des Matières » et « 1 Introduction : Corrélation vs Causalité dans l'IA » détaillent les raisons de cette importance. Pour un accompagnement sur ce sujet, [contactez nos experts](#).

### Comment mettre en œuvre les recommandations de cet article ?

Les recommandations pratiques sont détaillées tout au long de l'article, avec des commandes, des outils et des méthodologies éprouvées. La section « Conclusion » fournit une synthèse actionnable. Pour un accompagnement sur ce sujet, [contactez nos experts](#).

## Conclusion

---

Cet article a couvert les aspects essentiels de les concepts clés abordés. La mise en pratique de ces recommandations permet de renforcer significativement la posture de sécurité de votre organisation.

---

**Ayi NEDJIMI Consultants** — Expert cybersécurité offensive & intelligence artificielle

ayinedjimi-consultants.fr · ayi@ayinedjimi-consultants.fr

© 2026 — Reproduction interdite sans autorisation.