

Sécurisation des outils d'IA Générative

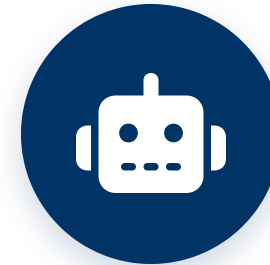
Guide des bonnes pratiques en entreprise

Stratégies et solutions pour protéger les données sensibles et assurer une utilisation conforme des technologies d'intelligence artificielle générative dans un environnement professionnel.

 Protection des données

 Conformité RGPD

 Gouvernance



Protection



Sécurité



Contrôle
d'accès



Configuration



Confidentialité

Table des matières

- 1  Introduction aux enjeux de sécurisation
- 2  Identification et analyse des risques
- 3  Classification & protection des données
- 4  Techniques d'anonymisation
- 5  Politiques de sécurité et contrôle d'accès
- 6  Paramétrage sécurisé des outils
- 7  Schéma – Vecteurs d'attaque sur l'IA générative
- 8  Bonnes pratiques opérationnelles
- 9  Processus de validation humaine
- 10  Gestion des incidents
- 11  Formation et sensibilisation
- 12  Veille technologique et réglementaire
- 13  Étude de cas : RATP
- 14+  Focus pratiques, recommandations, annexes



Introduction : Enjeux de la sécurisation autour de l'IA générative

L'essor de l'IA générative en entreprise

L'adoption massive des outils d'IA générative (ChatGPT, Claude, Gemini, etc.) transforme les méthodes de travail en offrant des capacités inédites d'analyse, de génération et d'assistance.

Avantages

- ✓ Productivité accrue
- ✓ Automatisation de tâches complexes
- ✓ Aide à la décision
- ✓ Innovation dans les processus métier

Risques et spécificités

Risques associés

- ⚠ Fuite de données confidentielles
- ⚠ Violation de la propriété intellectuelle
- ⚠ Non-conformité réglementaire
- ⚠ Perte de contrôle sur les informations sensibles

Spécificités des risques IA générative

- 🧠 Apprentissage continu : Mémorisation involontaire d'informations
- 🔗 Diffusion incontrôlée : Données utilisables pour améliorer les modèles
- 📦 Complexité technique : Systèmes "boîte noire" difficiles à évaluer

💡 Ces enjeux nécessitent une approche structurée combinant mesures techniques, organisationnelles et formation

Contexte réglementaire : RGPD, IA Act, secteur

Double exigence réglementaire



RGPD

- Protection des données personnelles
- Consentement et transparence
- Minimisation des données
- Droit à l'explication des décisions automatisées



IA Act

- Classification par niveau de risque
- Exigences techniques et de traçabilité
- Obligations de transparence
- Surveillance humaine des systèmes

Contraintes sectorielles spécifiques



Santé

HDS, HIPAA, confidentialité patient, recherche clinique



Finance

DSP2, MIFID II, LCB-FT, conformité bancaire



Transport

Données de mobilité, sécurité, billettique



Télécom

Code européen, protection des communications



L'utilisation de l'IA générative requiert une approche réglementaire intégrée combinant la protection des données personnelles, le cadre IA et les exigences sectorielles spécifiques

Approche méthodologique de sécurisation



i Une démarche cyclique et itérative pour maîtriser les risques liés à l'IA générative

Typologie des risques : divulgation, inférence, mémorisation



Risques de divulgation

- > Fuite directe via prompts/réponses
- > Inférence d'informations sensibles
- > Cross-contamination entre contextes



Risques de mémorisation

- > Intégration des données d'entraînement
- > Reproduction de contenus confidentiels
- > Apprentissage de patterns organisationnels







Risques réglementaires

- > Non-conformité RGPD
- > Manquement IA Act
- > Sanctions sectorielles spécifiques



Impact sur les individus

-  Atteinte à la vie privée
-  Usurpation d'identité
-  Discrimination potentielle
-  Préjudice financier



Impact organisationnel





-  Réputation dégradée
-  Pertes financières
-  Fuite de propriété intellectuelle
-  Perturbation opérationnelle

Schéma des vecteurs d'attaque sur l'IA générative



LÉGENDE DES VECTEURS D'ATTAQUE

■ Prompt Injection

■ Extraction de Données

🛡️ Défense recommandée: filtrage des prompts, tokenisation sécurisée et monitoring

Classification et protection des données sensibles

Niveaux de classification des données



Public

Données librement communicables sans restriction



Interne

Données pour usage interne à l'organisation



Confidentiel

Données dont la divulgation pourrait nuire



Secret

Données causant un préjudice grave si divulguées



Top Secret

Données critiques pour la survie de l'organisation

Critères de classification



Origine

Source de la donnée (interne, externe, partenaire)



Sensibilité

Niveau de sensibilité intrinsèque de l'information



Impact potentiel

Conséquences possibles en cas de divulgation



Durée de vie

Période pendant laquelle la donnée reste sensible









Réglementation

Contraintes légales et réglementaires applicables

i La classification adéquate des données est la première étape cruciale pour définir les mesures de protection appropriées avec les outils d'IA générative

Exemples de données RGPD/IA Act et niveau de protection





■ Protection Base
 ■ Protection Renforcée
 ■ Protection Maximale

Type de données	Exemples	Cadre réglementaire	Niveau de protection	Mesures principales
 Données personnelles courantes	Nom, adresse, email, téléphone	RGPD	BASE	<ul style="list-style-type: none"> • Chiffrement en transit • Contrôle d'accès (RBAC) • Journalisation des accès
 Données biométriques	Empreintes, reconnaissance faciale, vocale, ADN	RGPD (Catégorie particulière)	MAXIMALE	<ul style="list-style-type: none"> • Isolation physique et logique • Chiffrement de bout en bout • Audit continu, accès restreint
 Données de santé	Dossiers médicaux, handicap, arrêts maladie	RGPD (Catégorie particulière)	MAXIMALE	<ul style="list-style-type: none"> • Habilitations spécifiques <ul style="list-style-type: none"> • Chiffrement avancé • Approbation préalable
 Ressources humaines	Évaluations, recrutement, données disciplinaires	IA Act RGPD	RENFORCÉE	<ul style="list-style-type: none"> • MFA obligatoire • Chiffrement avancé • Monitoring temps réel
 Propriété intellectuelle	Brevets, secrets commerciaux, R&D, algorithmes	IA Act Droit PI	MAXIMALE	<ul style="list-style-type: none"> • Environnement isolé • Contrôles d'accès stricts • Journalisation détaillée
 Infrastructure critique	Architecture système, configurations réseau	IA Act NIS 2	RENFORCÉE	<ul style="list-style-type: none"> • Ségrégation réseau • Monitoring continu • Tests d'intrusion réguliers

 Le niveau de protection doit être adapté aux risques spécifiques identifiés pour chaque type de données

Techniques d'anonymisation des données

Techniques de base

-  **Suppression directe**
Retrait complet des identifiants directs (noms, adresses, etc.)
-  **Pseudonymisation**
Remplacement des identifiants par des codes avec table de correspondance sécurisée
-  **Généralisation**
Réduction de précision (âge exact → tranches d'âge, adresse → ville)
-  **Ajout de bruit**
Introduction de modifications aléatoires dans les données numériques

Exemple concret




Données originales:

Jean Dupont, 34 ans
Paris (75012)
45000€/an





Données anonymisées:

ID_47XB, 30-40 ans
Paris
40-50K€/an

Méthodes avancées

-  **K-anonymat**
Chaque individu est indistinguishable d'au moins k-1 autres dans le dataset
-  **L-diversité**
Extension du k-anonymat garantissant une diversité dans les attributs sensibles
-  **Confidentialité différentielle**
Ajout de bruit calibré avec garanties mathématiques de protection

Critères de sélection

-  **Type de données**
structurées vs. non-structurées
-  **Utilité analytique**
précision requise pour analyses
-  **Niveau de risque**
sensibilité et menaces
-  **Exigences légales**
RGPD, secteur spécifique

 La combinaison de plusieurs techniques permet d'optimiser l'équilibre entre protection des données et préservation de leur utilité

Panorama visuel des techniques d'anonymisation avancées

Suppression directe

Retrait complet des identifiants

Protection Faible

✓ **Avantages**

- Simple et rapide
- Facile à mettre en œuvre

✗ **Limites**

- Vulnérable aux attaques par recoupement
- Protection minimale

Exemple
Suppression des noms, emails, téléphones dans les données clients

Pseudonymisation

Remplacement par codes avec table de correspondance

Protection Moyenne

✓ **Avantages**

- Maintient les liens entre données
- Réversible si nécessaire

✗ **Limites**

- Protection de la table de correspondance
- Ne protège pas contre toutes les inférences

Exemple
Jean Dupont → ID_45789 (santé), User123 (logs)

Généralisation

Réduction de la précision des données

Protection Moyenne

✓ **Avantages**

- Préserve l'utilité statistique
- Flexible selon le besoin

✗ **Limites**

- Perte de précision
- Difficile à calibrer

Exemple
Âge 34 ans → tranche 30-40 ans; Adresse → Code postal

Ajout de bruit

Modifications aléatoires des valeurs

Protection Moyenne-Élevée

✓ **Avantages**

- Conserve les distributions
- Paramétrage précis possible

✗ **Limites**

- Complexité technique
- Impact sur précision analytique

Exemple
Revenus 45K€ → valeur aléatoire dans plage 42-48K€

K-anonymat

Indistinguabilité dans des groupes de k individus

Protection Élevée

✓ **Avantages**

- Garantie formelle
- Mesurable et vérifiable

✗ **Limites**

- Sensible aux attaques par homogénéité
- Coûteux pour grands k

Exemple
Pour k=5, chaque combinaison d'attributs apparaît ≥ 5 fois

Confidentialité différentielle

Bruit calibré mathématiquement

Protection Très élevée

✓ **Avantages**

- Garantie mathématique forte
- Résistance aux attaques avancées

✗ **Limites**

- Très complexe à implémenter
- Impact significatif sur l'utilité avancées

Exemple
 ϵ -DP garantit l'indistinguabilité avec budget de bruit ϵ



Politiques de sécurité et gouvernance de l'accès

Architecture de sécurité



Modèle Zero Trust

"Ne jamais faire confiance, toujours vérifier" - Authentification continue et principe du moindre privilège



Authentification MFA

Authentification multi-facteurs obligatoire pour tous les accès aux outils d'IA générative



Contrôle RBAC

Accès basé sur les rôles avec permissions spécifiques selon les besoins métiers

Gouvernance et politique d'usage



Comité IA

Instance décisionnelle (CISO, DPO, Innovation) pour l'évaluation et l'approbation des outils



Cycle de vie des accès

Provisioning contrôlé, révision trimestrielle, déprovisioning immédiat en cas de départ



Politique documentaire

Directives claires, procédures standardisées, traçabilité des validations et des usages

Classification des usages IA

Usage autorisé sans restriction

Contenu créatif, documents génériques, traduction de textes non sensibles

Usage contrôlé

Analyses de données agrégées, rapports sur données publiques, avec validation

Usage interdit

Données personnelles sensibles, propriété intellectuelle critique, décisions automatisées

Architecture de restriction d'accès pour l'IA générative

Modèle Zero Trust : "Ne jamais faire confiance, toujours vérifier"



Couche d'authentification

MFA obligatoire | SSO sécurisé | Authentification contextuelle

Protection périmétrique



Gestion des rôles (RBAC)

Moindre privilège | Ségrégation des droits | Attribution dynamique



Utilisateur

Fonctions basiques



Analyste

Analyses avancées



Data Scientist

Fine-tuning

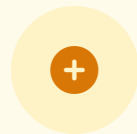


Admin

Configuration



Cycle de vie des accès



Provisioning

Validation manager + sécurité



Revue

Audit trimestriel



Déprovisioning

Automatique & immédiat

Flux de données sécurisé



Utilisateur authentifié



Proxy filtrant



Anonymisation



IA générative



Validation sortie

Paramétrage sécurisé des outils d'IA générative

Paramètres de confidentialité critiques



Désactivation de l'historique

Empêche la conservation de vos conversations et requêtes

ChatGPT: Settings → Data Controls → Chat History → Off
Claude: Conversation Memory → Disabled



Opt-out des données d'entraînement

Refus d'utilisation de vos données pour améliorer le modèle

OpenAI: Settings → Data Controls → Improve the model → Off
Anthropic: Data Retention → Minimal/None



Configuration API sécurisée

Utilisation de l'API avec contrôles renforcés

- Clés API avec rotation régulière
- IP whitelisting + TLS 1.3
- Rate limiting + Logs complets

Options de déploiement sécurisé



Déploiement on-premise

Contrôle total des données Llama 2 / Mistral Infrastructure GPU



Cloud privé virtuel (VPC)

Isolation réseau Azure OpenAI dédié AWS Bedrock privé

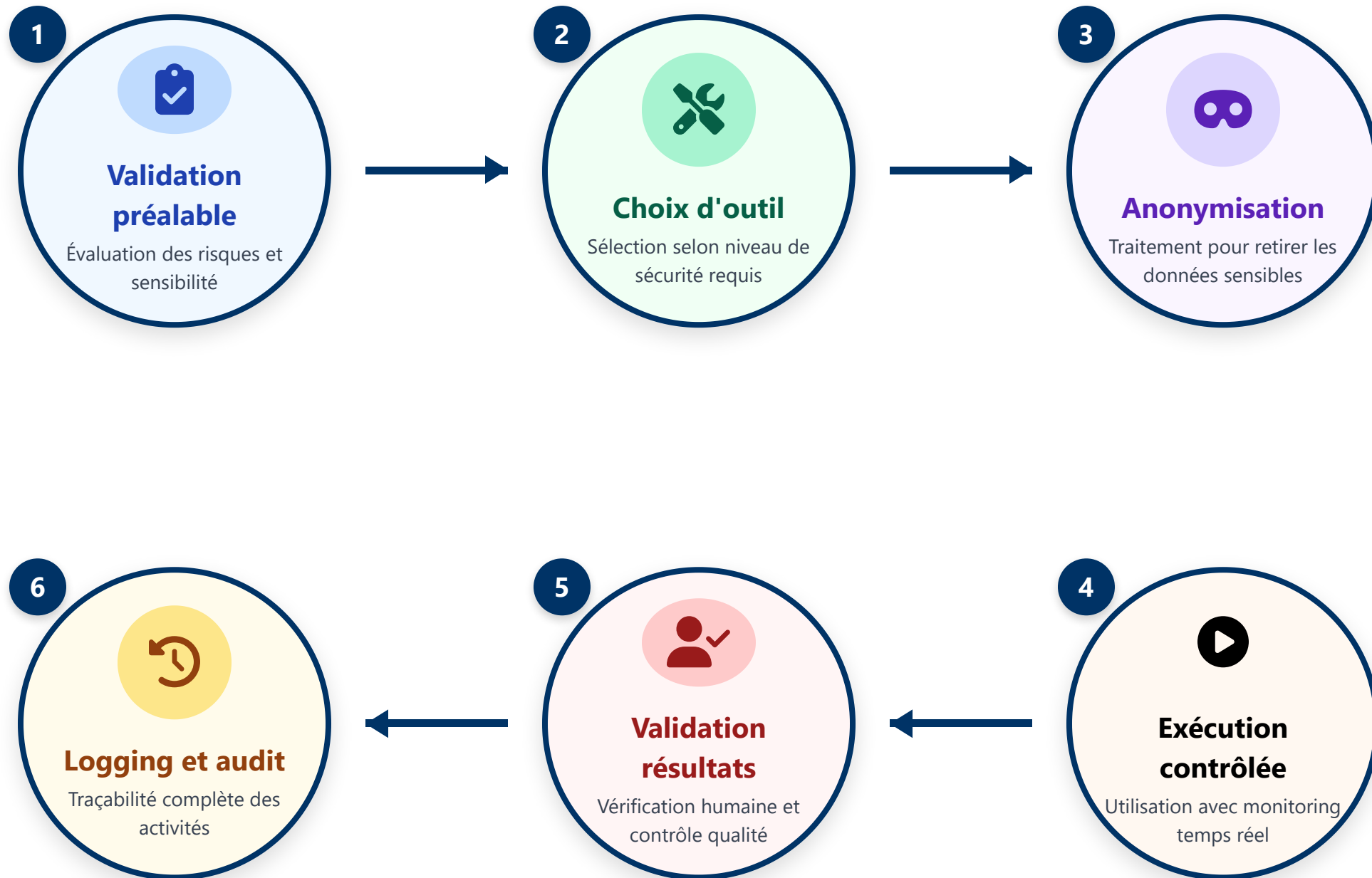


Architecture hybride

Données sensibles on-premise Passerelle sécurisée Orchestration centralisée

🛡️ Vérifier régulièrement les paramètres : les politiques de confidentialité des fournisseurs évoluent fréquemment

Workflow sécurisé : bonnes pratiques opérationnelles





Un processus standardisé réduit les risques d'exposition de données sensibles




Validation humaine et contrôles manuels

Valeur ajoutée de l'expertise humaine

Limites de l'automatisation

-  Détection contextuelle limitée des informations sensibles
-  Taux élevé de faux positifs/négatifs

Avantages de l'expertise humaine

-  Analyse contextuelle approfondie
-  Détection d'anomalies subtiles
-  Arbitrage utilité/risque en situations complexes

Processus de validation multi-niveaux

Niveau 1 : Auto-contrôle utilisateur

Vérification préliminaire avant soumission

Niveau 2 : Validation par un pair

Revue croisée et partage d'expérience

Niveau 3 : Validation experte

Analyse approfondie sécurité/conformité

Outils d'aide à la validation

 NER avancée

 Scoring de risque

 La validation humaine n'est pas une option mais une nécessité pour garantir l'utilisation responsable de l'IA générative

Gestion des incidents et réponse aux violations



Violations de données personnelles

Exposition directe, fuite de données, inférences



Violations de confidentialité

Propriété intellectuelle, stratégies, informations sensibles



Incidents techniques

Compromission d'accès, prompt injection, exploitation



Monitoring en temps réel

Surveillance continue, détection d'anomalies, alertes

Équipe de réponse (IRT)



Responsable incident



Expert sécurité



DPO



Juriste



Communication



Expert métier



Formation, sensibilisation et culture sécurité

Formation adaptée par profil



Utilisateurs finaux

Sensibilisation aux risques, bonnes pratiques quotidiennes



Managers & Décideurs

Gouvernance, responsabilités, arbitrages sécurité/usage



Experts techniques & DPO

Sécurité avancée, anonymisation, conformité

Méthodes pédagogiques innovantes



Gamification

Simulations d'incident, compétitions, défis de sécurité



Immersion

Scénarios interactifs, exercices pratiques adaptés



Communautés

Partage d'expérience, mentoring entre pairs



Certification

Parcours certifiants adaptés au niveau d'expertise



Cycle d'amélioration continue



Formation



Pratique



Évaluation

Veille technologique et réglementaire

Organisation de la veille



Comité de veille stratégique

DPO, CISO, Responsable Innovation IA, Juriste spécialisé



Réseau opérationnel

Experts techniques, correspondants réglementaires, partenaires

Sources principales

Sources institutionnelles



Commission européenne (DG Connect, DG Justice)



CNIL, CEPD, ENISA (agence cybersécurité)



NIST, ISO/IEC, OWASP (normalisation)

Suivi des évolutions

Évolution du cadre européen

- Actes délégués et d'exécution de l'IA Act
- Nouveaux avis et jurisprudence RGPD
- Réglementations sectorielles connexes

Évolution technologique

- Nouvelles capacités des modèles IA
- Technologies de protection émergentes
- Outils de sécurisation et audit de modèles

Adaptation des pratiques



Évaluation d'impact

Analyse des nouvelles obligations et risques



Roadmap d'évolution

Priorisation actions court/moyen/long terme

 Une veille active permet d'anticiper les évolutions réglementaires et d'adapter proactivement les mesures de sécurité

Étude de cas : RATP – Anonymisation et gouvernance



Contexte

Partage de données de validation des titres de transport (15M validations/jour) avec des chercheurs universitaires pour optimiser les flux de mobilité urbaine tout en protégeant la vie privée des usagers.

Anonymisation multicouches

1

Suppression identifiants directs

Numéros Navigo, codes PIN, données bancaires, noms

2

Pseudonymisation contrôlée

HMAC sécurisé avec rotation des clés

3

Généralisation spatio-temporelle

Granularité heure (vs minute), districts (vs stations précises)

4

K-anonymat (k=5) + Confidentialité différentielle

Groupes indistinguables et ajout de bruit calibré

Pipeline sécurisé



Ingestion



Anonymisation



Distribution

Gouvernance



Comité d'éthique des données



Accord juridique avec clauses strictes



Dashboard de suivi temps réel



Tests de réidentification réguliers

Enseignements et recommandations



Commencer petit et tester intensivement avant déploiement complet



Investir dans la gouvernance avec implication des usagers



Équilibrer protection et valeur scientifique des données

Conclusion : recommandations et synthèse



La sécurisation de l'IA générative nécessite un **équilibre entre innovation et protection**, soutenu par une approche méthodique et une culture de sécurité partagée.



Approche par les risques

Évaluation adaptée à chaque contexte d'utilisation



Gouvernance et validation

Processus clairs et responsabilités définies



Anonymisation adaptative

Techniques proportionnelles à la sensibilité des données



Culture de sécurité

Formation continue et sensibilisation de tous les acteurs



**La sécurisation n'est pas un frein mais un accélérateur
de confiance et d'adoption responsable de l'IA générative**