

# Livre Blanc : Confluent Cloud

## Architecturer Confluent Cloud à l'Échelle pour les Systèmes Mission-Critical

Auteur : Directeur de l'Architecture d'Entreprise

Date : Octobre 2025

Sujet : Stratégies avancées pour le déploiement, la gouvernance et l'opérationnalisation de Confluent Cloud dans des environnements d'entreprise complexes.

### Résumé Exécutif

La transition des architectures statiques, fondées sur la donnée au repos, vers des systèmes dynamiques et événementiels représente le changement de paradigme le plus significatif de la dernière décennie pour l'infrastructure informatique d'entreprise. Alors que les organisations matures dans leur adoption d'Apache Kafka et de Confluent Cloud, l'objectif se déplace de la simple ingestion de données vers l'établissement d'un système nerveux central garantissant résilience, sécurité et gouvernance à une échelle mondiale.

Ce livre blanc sert de guide architectural définitif pour "L'Entreprise en Temps Réel". Il dissèque les piliers critiques nécessaires pour opérer Confluent Cloud à un niveau industriel. Notre analyse priorise les exigences non fonctionnelles qui distinguent les preuves de concept des plateformes critiques : la Haute Disponibilité (HA) et la Reprise après Sinistre (DR), une connectivité réseau rigoureuse, des pratiques de Platform Engineering robustes, une gouvernance des données applicable par contrat, et l'intégration émergente de l'IA Générationnelle (GenAI) dans les pipelines de streaming.

Ce document explore les compromis nuancés entre les Clusters Multi-Régions et le Cluster Linking, la mise en œuvre de modèles "Data Strangler" pour le délestage des mainframes, et l'opérationnalisation du raisonnement de l'IA via la journalisation "Chain of Thought". Il fournit la profondeur technique nécessaire aux architectes pour concevoir des systèmes non seulement temps réel, mais antifragiles.

### Partie I : Haute Disponibilité

Dans le domaine des systèmes distribués, la disponibilité n'est pas un état binaire mais une fonction de probabilité dépendant de la résilience de l'infrastructure, de la latence de réPLICATION et de la logique de basculement (failover). Pour l'Entreprise en Temps Réel, l'objectif est de minimiser le rayon d'impact de toute défaillance — qu'il s'agisse d'une panne

zonale ou d'une catastrophe régionale complète — tout en gérant les compromis inhérents entre la cohérence (RPO) et la disponibilité (RTO).

## 1.1 Modèles Architecturaux : Clusters Multi-Régions vs Cluster Linking

Une décision architecturale fondamentale lors de la conception pour la reprise après sinistre sur Confluent Cloud est le choix entre les "Stretch Clusters" synchrones (Clusters Multi-Régions ou MRC) et la réPLICATION asynchrone via Cluster Linking. Ce choix dicte la tolérance du système à la perte de données et ses caractéristiques de latence.

### 1.1.1 Clusters Multi-Régions (MRC) : L'Isomorphisme du Cluster Étendu

Les Clusters Multi-Régions, souvent appelés "stretch clusters", opèrent comme un cluster Kafka logique unique étendu sur plusieurs zones géographiques ou régions proches.<sup>1</sup> La caractéristique déterminante du MRC est sa capacité à offrir un **RPO = 0** (Recovery Point Objective) grâce à la réPLICATION synchrone. Dans cette architecture, une requête de production n'est acquittée (acks=all) que lorsque la donnée est écrite sur le leader et les réplicas synchronisés (ISR), qui sont distribués à travers les centres de données.<sup>3</sup>

Cependant, cette cohérence a un coût imposé par les lois de la physique ; le temps d'aller-retour (RTT) entre les centres de données s'ajoute directement à la latence du producteur.<sup>5</sup> Par conséquent, le MRC n'est généralement viable que lorsque la latence réseau entre les régions est stable et faible (typiquement inférieure à 100ms).<sup>6</sup>

Pour mitiger l'impact sur la latence d'écriture tout en maintenant une présence multi-régionale, Confluent Server introduit un nouveau type de réplica : les **Observateurs**. Les observateurs sont des réplicas asynchrones qui ne participent pas à l'ISR (In-Sync Replicas) et ne peuvent pas devenir leaders tant qu'ils n'ont pas rattrapé leur retard, mais ils permettent une scalabilité de lecture locale dans une région distante sans pénaliser la latence d'écriture du cluster primaire.<sup>2</sup>

Placement des Réplicas et Topologie :

Une configuration typique pour un MRC robuste implique une distribution 2.5 DC (Data Center) ou 3 DC pour éviter les situations de "split-brain" (cerveau divisé). La formule recommandée pour calculer le facteur de réPLICATION est :

\$\$\text{Replication Factor} = \text{Nombre de DC} \times \text{Réplicas par DC}\$\$

Par exemple, avec 3 DCs et 2 réplicas par DC, le facteur de réPLICATION global est de 6, garantissant une durabilité extrême.<sup>8</sup>

### 1.1.2 Cluster Linking : La RéPLICATION Cloud-Native

Le Cluster Linking représente l'évolution cloud-native de la géo-réPLICATION. Contrairement à MirrorMaker 2, qui nécessite l'opération d'un cluster Kafka Connect séparé ajoutant une

complexité opérationnelle et des coûts de latence, le Cluster Linking est intégré directement dans le binaire du broker.<sup>1</sup> Il permet une réPLICATION "byte-for-byte", préservant parfaitement les offsets entre le cluster source et le cluster de destination.<sup>1</sup>

Cette préservation des offsets est critique pour les scénarios de DR actif-passif car elle permet aux consommateurs de basculer vers le cluster de DR et de reprendre le traitement exactement au même offset, une capacité non nativement garantie par les outils de réPLICATION externes sans traduction complexe d'offsets.<sup>10</sup>

L'avantage architectural du Cluster Linking réside dans son couplage lâche. La source et la destination sont des clusters distincts, ce qui signifie qu'une corruption de métadonnées ou une défaillance catastrophique dans un cluster ne se propage pas nécessairement à l'autre.<sup>12</sup> Pour les déploiements mondiaux couvrant des continents (par exemple, AWS us-east-1 vers eu-central-1), le Cluster Linking est le modèle mandaté en raison de sa tolérance à une latence réseau élevée et de sa capacité à fonctionner de manière asynchrone, découplant la performance du producteur du lag de réPLICATION.<sup>1</sup>

#### **Tableau Comparatif Technique : MRC vs Cluster Linking**

<b>Caractéristique Technique</b>	<b>Multi-Region Cluster (MRC)</b>	<b>Cluster Linking</b>
<b>Mode de RéPLICATION</b>	Synchrone (par défaut pour ISR)	Asynchrone (natif)
<b>Objectif RPO</b>	Zéro (0) - Cohérence Forte	> 0 (Dépendant du Lag réseau)
<b>Modèle de Couplage</b>	Fort (Cluster Logique Unique)	Faible (Clusters Indépendants)
<b>Tolérance Latence Réseau</b>	Faible (< 50-100ms recommandé)	Élevée (Global/Inter-Cloud)
<b>Basculement Client (Failover)</b>	Transparent (Mise à jour métadonnées)	Explicite (Changement Bootstrap)
<b>Préservation des Offsets</b>	Native (Même cluster)	Native (Copie byte-for-byte) <sup>1</sup>

<b>Cas d'Usage Principal</b>	Transactions Financières, RPO=0	DR Cross-Région, Migration, Agrégation
------------------------------	------------------------------------	---

## 1.2 Définir et Calculer RPO et RTO dans les Contextes de Streaming

Quantifier la résilience nécessite des définitions précises du Recovery Point Objective (RPO) et du Recovery Time Objective (RTO) dans le contexte spécifique de Kafka.

Recovery Point Objective (RPO) :

Le RPO mesure la tolérance maximale à la perte de données. Dans le contexte du Cluster Linking, le RPO est déterminé par le mirroring lag (retard de miroir). Si la région primaire échoue, tout message produit sur le leader mais pas encore récupéré (fetch) par le lien de cluster est perdu.<sup>12</sup>

La formule pour estimer le RPO en secondes est approximativement :

$$\$\$RPO \approx \frac{\text{Mirroring Lag (bytes)}}{\text{Taux de RéPLICATION (bytes/sec)}} + \text{Latence Réseau}$$$$

Pour atteindre un RPO proche de zéro avec Cluster Linking, la bande passante entre les régions doit excéder le débit d'ingestion (ingress throughput) du cluster source afin d'empêcher l'accumulation de lag.<sup>14</sup> Il est impératif de surveiller la métrique `mirroring_lag` via l'API Metrics de Confluent Cloud.<sup>12</sup>

Recovery Time Objective (RTO) :

Le RTO mesure la durée de l'interruption de service. Le RTO n'est pas uniquement une propriété de l'infrastructure, mais dépend largement de la logique de basculement des clients.<sup>12</sup> Bien que le cluster de destination puisse être promu en mode écriture en quelques secondes à l'aide de la commande failover ou promote<sup>13</sup>, le RTO de bout en bout inclut le temps nécessaire aux producteurs et consommateurs pour détecter la panne, redémarrer et se reconnecter aux nouveaux serveurs d'amorçage (bootstrap servers).<sup>12</sup>

## 1.3 Patterns de Basculement Client et Implémentation

Atteindre un RTO faible exige une logique côté client sophistiquée. Coder en dur les serveurs d'amorçage est un anti-pattern majeur.<sup>12</sup> Une découverte dynamique et une configuration robuste sont requises.

### 1.3.1 Intégration de la Découverte de Service

Dans une architecture DR robuste, les clients doivent résoudre leurs serveurs d'amorçage via un mécanisme de découverte de service (ex: HashiCorp Consul, AWS Route53 avec des health checks, ou un sidecar de configuration).<sup>12</sup> Lors d'un événement de basculement, l'équipe des opérations met à jour l'entrée de découverte de service pour pointer vers le endpoint du cluster de DR.

### 1.3.2 Configuration des Clients Java pour la Résilience

Pour les clients Java, des configurations spécifiques dictent la rapidité avec laquelle ils détectent une perte de connexion et tentent de se reconnecter.

- **session.timeout.ms** : Définit le temps maximum que le broker attend pour un heartbeat du consommateur avant de le considérer comme mort. Augmenter cette valeur prévient les faux positifs lors de micro-coupures réseau mais retarde le rééquilibrage (rebalance) lors de pannes réelles. La valeur par défaut a été augmentée à 45000ms (45s) dans les versions récentes.<sup>17</sup>
- **client.dns.lookup** : Régler ce paramètre sur `use_all_dns_ips` permet au client de tenter des connexions vers toutes les adresses IP résolues, ce qui est crucial dans les environnements cloud où l'IP d'un seul broker peut changer ou devenir inatteignable.<sup>19</sup>
- **bootstrap.servers** : Fournir une liste d'au moins 3 brokers réduit le risque d'échec de connexion initiale.<sup>20</sup>

#### Exemple de Code : Logique de Basculement Dynamique (Java Conceptuel)

Java

```
// Exemple conceptuel de configuration consciente du basculement (Failover-Aware)
Properties props = new Properties();

// Dans une implémentation réelle, ceci serait récupéré depuis une source de config dynamique
// (ex: Consul, Vault, ou un appel API interne)
String activeBootstrap = configService.getKafkaEndpoint();

props.put(ProducerConfig.BOOTSTRAP_SERVERS_CONFIG, activeBootstrap);
props.put(ProducerConfig.ACKS_CONFIG, "all"); // Garantie de durabilité maximale
props.put(ProducerConfig.RETRIES_CONFIG, Integer.MAX_VALUE); // Réessais infinis pour
pannes transitoires
props.put(ProducerConfig.DELIVERY_TIMEOUT_MS_CONFIG, 120000); // 2 minutes pour livrer
avant erreur

// Configuration cruciale pour la résolution DNS robuste
props.put("client.dns.lookup", "use_all_dns_ips");

KafkaProducer<String, String> producer = new KafkaProducer<>(props);

try {
    // Envoi synchrone pour détection immédiate en cas de scénario critique
    producer.send(record).get();
} catch (TimeoutException e) {
    // Logique pour déclencher un Circuit Breaker ou vérifier le basculement DR
```

```

    // Si l'exception persiste, l'application peut se mettre en pause et poller le configService
    circuitBreaker.open();
    handleOutage(e);
}

```

Comme illustré ci-dessus et dans les sources<sup>21</sup>, la logique de basculement peut être synchrone (bloquante sur `.get()`) ou asynchrone. Le mode synchrone permet une détection immédiate mais impacte le débit. Le pattern "Actif-Passif" exige généralement d'arrêter les producteurs dans la région primaire et de les redémarrer pointés vers la région secondaire pour préserver l'ordre et éviter les scénarios de "split-brain" où deux clusters accepteraient des écritures concurrentes.<sup>12</sup>

## 1.4 Dimensionnement pour la Performance et le Rattrapage (Catch-up)

Pour garantir que le cluster de DR peut gérer la charge lors d'un basculement, la planification de la capacité est essentielle. La limite de débit soutenu théorique ( $\$T_{\text{cluster}}$ ) peut être modélisée comme suit :

$$\$T_{\text{cluster}} \leq \min \left( \frac{T_{\text{storage}}}{N_{\text{brokers}} \cdot R}, \frac{T_{\text{net}}}{N_{\text{brokers}}} \right)$$

Où  $R$  est le facteur de réplication.<sup>22</sup> Dans un scénario de basculement, les consommateurs qui ont accumulé du retard ou qui redémarrent généreront un phénomène de "Thundering Herd" (troupeau en furie) de requêtes de lecture. Cela nécessite que le cluster de DR soit dimensionné non seulement pour l'ingestion en régime permanent, mais aussi pour des lectures de rattrapage à haut débit, exigeant souvent des capacités d'E/S disque supérieures (IOPS) pour gérer les lectures aléatoires hors du pagecache.<sup>15</sup>

## Partie II : Connectivité

Le "système nerveux" de l'Entreprise en Temps Réel repose sur une connectivité sécurisée et à faible latence. Alors que les entreprises évoluent vers des architectures hybrides et multi-cloud, la connectivité via l'Internet public est rarement suffisante en raison des exigences de conformité, de sécurité et de stabilité de latence.

### 2.1 Réseaux Privés : AWS PrivateLink et Transit Gateway

Confluent Cloud supporte des options de réseau privé robustes, notamment AWS PrivateLink et Transit Gateway (TGW).<sup>23</sup> Le choix entre ces technologies a des implications profondes sur la topologie réseau et la gestion DNS.

### 2.1.1 AWS PrivateLink : Sécurité Unidirectionnelle

PrivateLink est la méthode privilégiée pour une connectivité sécurisée et unidirectionnelle. Elle expose les brokers Confluent Cloud sous forme d'Interfaces Réseau Élastiques (ENI) au sein du VPC client, garantissant que le trafic ne traverse jamais l'Internet public.<sup>24</sup> Cette architecture simplifie la gestion des groupes de sécurité (Security Groups) car elle ne nécessite pas de peering de VPC complexe ni de gestion de chevauchement d'adresses IP (CIDR overlaps).

Le Défi du DNS "Split Horizon" :

Les clients Kafka dépendent de la configuration bootstrap.servers pour découvrir les brokers. Dans Confluent Cloud, ces URL de brokers (ex: lkc-xxxx.us-east-1.aws.confluent.cloud:9092) se résolvent par défaut vers des adresses IP publiques. Lors de l'utilisation de PrivateLink, ces mêmes URL doivent impérativement se résoudre vers les adresses IP privées des endpoints VPC (VPCE) situés dans le VPC du client.<sup>24</sup>

Cela nécessite une configuration **DNS Split Horizon**. Une zone hébergée privée (Private Hosted Zone) Route53 doit être créée dans le VPC client pour surcharger la résolution DNS publique du domaine Confluent Cloud.<sup>24</sup>

- **Point de Vigilance** : Un mode d'échec fréquent lors de la configuration est la création de la zone hébergée privée sans l'associer à tous les VPC clients, ou une mauvaise configuration des jeux d'enregistrements (Record Sets), conduisant les clients à se connecter par inadvertance via les endpoints publics ou à échouer lors du handshake SSL/TLS.<sup>26</sup>

Pour les environnements hybrides (On-Premise vers AWS), la résolution DNS se complexifie. Les serveurs DNS sur site doivent être capables de transférer les requêtes pour \*.confluent.cloud vers le résolveur Route53 du VPC via un **Route 53 Inbound Resolver Endpoint**. Une règle de transfert conditionnel (Forwarding Rule) sur le DNS on-premise est nécessaire pour diriger le trafic DNS vers les IPs du résolveur entrant AWS, qui a accès à la Private Hosted Zone.<sup>28</sup>

### 2.1.2 AWS Transit Gateway (TGW)

Pour des topologies réseau complexes impliquant de multiples VPC ou des centres de données sur site connectés via Direct Connect, AWS Transit Gateway agit comme un hub central.<sup>23</sup> Confluent Cloud permet des attachements TGW, autorisant un modèle "hub-and-spoke".

- **Routage** : Des routes statiques doivent être ajoutées à la table de routage TGW pour diriger le trafic destiné au bloc CIDR de Confluent Cloud vers l'attachement Confluent.<sup>23</sup>
- **Compromis** : TGW simplifie le routage transitif (VPC A \$\leftarrow\$ TGW \$\rightarrow\$ VPC B) mais introduit une latence de saut (hop) supplémentaire et des frais de traitement de données par Go comparé au peering VPC direct ou à PrivateLink.<sup>23</sup> Contrairement à PrivateLink, TGW nécessite une gestion stricte des plages d'adresses IP

pour éviter les conflits de CIDR.

## 2.2 Cloud Hybride et le Pattern "Strangler Fig"

Connecter les systèmes hérités (Legacy) sur site à Confluent Cloud est un cas d'usage primaire pour les architectures hybrides. Le **Strangler Fig Pattern** (ou Pattern de l'Étrangleur) est la stratégie dominante pour moderniser les applications monolithiques (ex: Mainframes) sans une réécriture risquée de type "Big Bang".<sup>30</sup>

### 2.2.1 L'Implémentation "Data Strangler"

Contrairement au Strangler Fig traditionnel basé sur les API, le *Data Strangler* exploite le **Change Data Capture (CDC)**.

1. **Interception (Intercept)** : Des outils comme Kafka Connect (avec Debezium ou des connecteurs mainframe propriétaires comme IBM InfoSphere ou Qlik Replicate) capturent les changements de base de données (DB2, VSAM, Oracle) en temps réel directement depuis les logs de transaction.<sup>32</sup>
2. **Routage (Route)** : Ces changements sont streamés vers des topics Confluent Cloud via une connexion sécurisée (VPN ou Direct Connect).
3. **Nouveaux Services** : Les microservices modernes consomment ces flux pour construire leurs propres magasins de données optimisés pour la lecture (ex: MongoDB, Elastic, S3).<sup>32</sup>
4. **Décommissionnement (Decommission)** : Une fois que le chemin de lecture est entièrement migré vers les microservices, le chemin d'écriture peut être basculé progressivement, permettant à terme de retirer le système hérité.<sup>30</sup>

Pattern Outbox vs. CDC :

Bien que le CDC soit transparent pour l'application héritée, il couple la structure du flux au schéma interne de la base de données (couplage fort). Le Pattern Outbox est souvent supérieur pour les nouveaux développements ou les refontes partielles. L'application écrit un événement dans une table "Outbox" dédiée au sein de la même transaction locale que ses données métier. Debezium capture ensuite uniquement cette table Outbox.<sup>34</sup> Cela garantit que le schéma de l'événement est explicite, stable, et découpé du modèle de données interne, offrant une meilleure gouvernance et une conception "event-first".<sup>34</sup>

---

## Partie III : Libre Service

Opérer Confluent Cloud à l'échelle exige une transition des opérations manuelles ("ClickOps") vers des pratiques rigoureuses de Platform Engineering. Cela implique de traiter la plateforme de streaming comme un produit interne offert aux équipes de développement.

### 3.1 Infrastructure as Code (IaC) avec Terraform

Le fournisseur Terraform Confluent (terraform-provider-confluent) permet la gestion déclarative de l'ensemble de l'écosystème Kafka.<sup>36</sup> Cela inclut non seulement l'infrastructure physique (Clusters, Réseaux) mais aussi les ressources logiques (Topics, ACLs, Schémas, Connecteurs).

#### Ressources Clés et Bonnes Pratiques :

- **confluent\_kafka\_cluster** : Définit la capacité du cluster (CKUs) et les détails du fournisseur cloud.<sup>36</sup> Il est crucial d'utiliser le bloc lifecycle { prevent\_destroy = true } pour les clusters de production afin d'éviter les suppressions accidentnelles.<sup>36</sup>
- **confluent\_private\_link\_access & confluent\_private\_link\_attachment** : Automatise la plomberie réseau complexe discutée dans la Partie II, réduisant les erreurs humaines lors de la configuration DNS et réseau.<sup>24</sup>
- **confluent\_kafka\_topic & confluent\_role\_binding** : Gère les ressources de données et le RBAC (Role-Based Access Control).<sup>37</sup>

#### Gestion de l'État Terraform :

Pour prévenir la "dérive de l'état" (State Drift) et les conflits dans les grandes équipes, des backends d'état distants (ex: S3 avec verrouillage DynamoDB) doivent être utilisés.<sup>37</sup> Une approche en couches est recommandée : isoler l'état de l'infrastructure (réseaux, clusters) de l'état de l'application (topics, schémas). Cela empêche qu'un changement de topic par une équipe applicative ne modifie accidentellement la configuration du cluster de production.<sup>37</sup>

### 3.2 GitOps : Monorepo vs. Polyrepo

La structure des dépôts Git définissant ces ressources Terraform est un sujet de débat architectural significatif.

Caractéristique	Monorepo	Polyrepo
<b>Centralisation</b>	Source de vérité unique ; outillage standardisé plus simple. <sup>41</sup>	Décentralisé ; s'aligne avec les frontières des microservices. <sup>42</sup>
<b>Coordination</b>	Changements atomiques à travers plusieurs services (ex: maj schéma producteur/consommateur) <sup>42</sup> .	Coordination plus difficile ; les dépendances doivent être gérées via versioning. <sup>43</sup>
<b>Contrôle d'Accès</b>	Nécessite des permissions complexes basées sur les répertoires	Isolation naturelle ; l'accès au dépôt équivaut à l'accès au service. <sup>43</sup>

	(CODEOWNERS). <sup>43</sup>	
<b>Scalabilité CI/CD</b>	Peut souffrir de goulets d'étranglement CI/CD si non optimisé (ex: ne builder que les chemins modifiés). <sup>43</sup>	Scale naturellement ; pipelines CI/CD indépendants. <sup>41</sup>

**Recommandation :** Pour l'équipe plateforme Kafka gérant l'infrastructure partagée (Clusters, Réseaux), une approche **Monorepo** est avantageuse pour maintenir la cohérence.<sup>44</sup> Cependant, pour les équipes applicatives définissant leurs propres topics et ACLs, une approche **Polyrepo** (ou fédérée) permet une plus grande autonomie, à condition que des garde-fous CI/CD stricts soient en place via des outils comme kafka-gitops ou des modules Terraform validés.<sup>43</sup>

### 3.3 Self-Service et Intégration Backstage

Pour passer à l'échelle au-delà de quelques dizaines d'équipes, une plateforme de self-service est essentielle. **Backstage**, le portail développeur open-source, s'impose comme l'interface standard.<sup>47</sup>

- **Plugin Integration** : Un plugin Backstage personnalisé peut interagir avec l'API Confluent Cloud (ou déclencher des exécutions Terraform) pour permettre aux développeurs de demander de nouveaux topics ou schémas via une interface conviviale.<sup>48</sup>
- **Catalogue Logiciel** : Le catalogue de Backstage peut ingérer des données depuis le Stream Catalog de Confluent, visualisant la propriété (Ownership) — montrant quelle équipe possède le topic-A et quels services le consomment.<sup>47</sup>
- **Templates Scaffolder** : Les modèles Scaffolder peuvent provisionner un "Chemin Doré" (Golden Path) pour les nouvelles applications de streaming, configurant automatiquement le dépôt Git, le Terraform pour les topics, et les ACLs nécessaires selon les standards de l'entreprise.<sup>48</sup>

### 3.4 Gestion des Défaillances : Dead Letter Queues (DLQ)

Automatiser la gestion des échecs est critique pour la fiabilité. Dans Kafka Connect, l'activation d'une Dead Letter Queue (DLQ) empêche un message unique "poison pill" (ex: JSON malformé) de bloquer tout un pipeline.<sup>49</sup>

- **Configuration** : Les paramètres errors.tolerance=all et errors.deadletterqueue.topic.name doivent être configurés.<sup>49</sup>
- **Enrichissement des En-têtes** : L'utilisation de errors.deadletterqueue.context.headers.enable=true ajoute des métadonnées au message échoué (raison de l'échec, stack trace, offset d'origine) dans les en-têtes Kafka. Cela facilite le débogage sans avoir à inspecter manuellement le contenu du payload.<sup>49</sup>

- **Stratégie de Rejeu (Replay)** : Dumper les messages dans une DLQ est insuffisant. Une stratégie de *rejeu* est requise. Des outils (personnalisés ou open-source) doivent être mis en place pour consommer la DLQ, permettre l'inspection/correction manuelle ou automatique, et re-produire le message vers le topic original ou un topic de tentative.<sup>51</sup>
- 

## Partie IV : Gouvernance

Dans une Entreprise en Temps Réel, Kafka sert souvent de système nerveux central. Si les données qui y circulent sont corrompues ou mal définies, l'entreprise entière souffre. Le principe "Garbage In, Garbage Out" à la vitesse du temps réel est catastrophique.

### 4.1 Schema Registry et Contrats de Données (Data Contracts)

Le Confluent Schema Registry a évolué d'un référentiel passif vers un point d'application actif des politiques. Les **Contrats de Données** formalisent l'accord entre producteurs et consommateurs, allant au-delà de la simple validation structurelle (int vs string) pour inclure des règles sémantiques.<sup>19</sup>

CEL (Common Expression Language) :

Confluent Schema Registry supporte désormais le langage CEL pour définir des règles de validation complexes directement au sein du schéma.<sup>19</sup> Cela permet une gouvernance "Shift Left", où le message du producteur est rejeté avant même d'atteindre le broker s'il viole une règle métier.

**Exemple : Contrat de Données avec Règle CEL (Représentation JSON)**

JSON

```
{
  "schemaType": "AVRO",
  "schema": "{\"type\":\"record\",\"name\":\"Transaction\",\"fields\":[{\"name\":\"amount\",\"type\":\"double\"},{\"name\":\"currency\",\"type\":\"string\"}]}",
  "rules": []
}
```

Cette configuration<sup>19</sup> garantit qu'aucune transaction avec un montant négatif ou un code devise invalide n'entrera jamais dans le topic, protégeant ainsi les systèmes analytiques et opérationnels en aval contre les problèmes de qualité des données.

## 4.2 Lignage des Données (Data Lineage) et Catalogue de Flux

Comprendre le flux de données est une exigence réglementaire (GDPR, traçabilité) et une nécessité opérationnelle.

- **Stream Lineage :** La console Confluent Cloud et l'API fournissent un graphe visuel des nodes (topics, apps) et edges (flux de données).<sup>55</sup> L'API renvoie une structure JSON définissant ces relations, ce qui permet l'intégration avec des catalogues de données d'entreprise comme Collibra, DataHub ou Apache Atlas.<sup>56</sup>
  - Le format JSON de l'API Lineage expose les entités (source, destination) et les métriques associées, permettant de construire des vues de dépendance personnalisées.<sup>55</sup>
- **Stream Catalog :** Il permet le balisage (tagging) des données sensibles (ex: PII, PCI). Ces tags peuvent ensuite piloter l'application de politiques, telles que le chiffrement automatique ou le masquage de champs à l'aide de Single Message Transforms (SMT) comme MaskField basés sur la présence du tag.<sup>58</sup>

**Insight d'Intégration :** En combinant Terraform (pour appliquer les tags) et Schema Registry (pour appliquer les règles), la gouvernance devient codifiée ("Governance as Code"). Un standard architectural peut mandater que tous les topics tagués Environment:Production doivent impérativement avoir un Contrat de Données avec compatibility=BACKWARD pour assurer la stabilité des consommateurs.<sup>60</sup>

---

## Partie V : Agentique

La convergence de l'Architecture Événementielle (EDA) et de l'IA Générative crée un nouveau paradigme : l'**Entreprise en Temps Réel Augmentée par l'IA**. Kafka et Flink ne se contentent plus de déplacer des données ; ils deviennent le moteur de contexte pour les LLM (Large Language Models).

### 5.1 Retrieval Augmented Generation (RAG) en Temps Réel

Les architectures RAG traditionnelles s'appuient sur des ETL par lots pour mettre à jour les bases de données vectorielles, ce qui conduit à un contexte obsolète. Une architecture RAG en streaming garantit que le LLM a accès à l'état le plus actuel du monde.<sup>61</sup>

**Le Pipeline Flink SQL pour RAG :**

1. **Ingestion :** Les mises à jour de produits ou événements clients affluent dans Kafka.
2. **Embedding (Vectorisation) :** Flink SQL invoque un endpoint de modèle (ex: OpenAI, AWS Bedrock, Google VertexAI) pour générer des embeddings vectoriels pour les données textuelles en temps réel via des fonctions définies par l'utilisateur (UDF) ou des appels natifs ML\_PREDICT.<sup>16</sup>

3. **Stockage** : Les embeddings sont déversés (sunk) dans une Base de Données Vectorielle (Pinecone, MongoDB Atlas, Weaviate) via Kafka Connect.<sup>63</sup>
4. **Requête** : Lorsqu'un utilisateur interroge le LLM, l'application interroge la Vector DB pour obtenir le contexte le plus récent.

Syntaxe Flink SQL CREATE MODEL :

La capacité de définir des modèles d'IA directement en SQL démocratise l'ingénierie de l'IA.

SQL

```
CREATE MODEL `embedding_model`
INPUT (input_text STRING)
OUTPUT (vector_embedding ARRAY<FLOAT>)
WITH (
  'TASK' = 'embedding',
  'PROVIDER' = 'openai',
  'OPENAI.CONNECTION' = 'openai-conn-id'
);

-- Utilisation du modèle pour enrichir le flux en temps réel
INSERT INTO vector_topic
SELECT product_id, embedding_model(description)
FROM product_updates;
```

Cette approche déclarative <sup>16</sup> permet aux développeurs SQL de construire des pipelines d'IA sophistiqués sans gérer de microservices Python complexes, réduisant la latence entre la création de la donnée et sa disponibilité pour l'inférence.

## 5.2 Journalisation "Chain of Thought" et Observabilité de l'IA

À mesure que les agents IA deviennent plus autonomes, le débogage de leur processus décisionnel devient critique. Kafka agit comme la "Boîte Noire" (Flight Recorder) immuable pour les agents IA.

Pattern "Chain of Thought" (CoT) :

Les agents chargés de raisonnements complexes (ex: "Analyser cette demande de prêt") décomposent le problème en étapes. Chaque étape — la pensée, l'outil utilisé, le résultat — doit être émise comme un événement structuré vers un topic Kafka.<sup>66</sup>

Schéma de Log JSON Proposé pour le Raisonnement IA :

Pour rendre ces données interrogables et conformes aux standards émergents comme les Conventions Sémantiques GenAI d'OpenTelemetry 68, un schéma structuré est requis :

JSON

```
{  
  "trace_id": "Oaf7651916cd43dd8448eb211c80319c",  
  "agent_id": "loan-processor-v2",  
  "timestamp": "2023-10-27T10:00:00Z",  
  "interaction_type": "chain_of_thought",  
  "step_index": 2,  
  "thought_process": {  
    "reasoning": "Revenu du demandeur > 50k, vérification du score de crédit nécessaire.",  
    "tool_invoked": "credit_bureau_api",  
    "tool_input": {"ssn_hash":  
      "e3b0c44298fc1c149afb4c8996fb92427ae41e4649b934ca495991b7852b855"}  
  },  
  "model_config": {  
    "provider": "openai",  
    "model": "gpt-4",  
    "temperature": 0.1  
  },  
  "gen_ai.system.token_count": 150  
}
```

- **Insight de Troisième Ordre :** En journalisant le "processus de pensée" de manière structurée dans Kafka, les organisations peuvent effectuer des analyses hors ligne pour détecter la "dérive des hallucinations" ou les biais dans les décisions des agents au fil du temps. Cela transforme Kafka d'un simple tuyau de données en un outil de conformité et de gouvernance pour l'IA.<sup>70</sup> L'utilisation de sorties structurées (Structured Outputs) garantit que les logs générés par le LLM respectent strictement le schéma JSON défini, facilitant leur ingestion et analyse automatique.<sup>73</sup>

---

## Conclusion

L'Entreprise en Temps Réel ne se construit pas uniquement sur des outils, mais sur des patterns architecturaux rigoureux. Elle exige un changement de mentalité : passer de "maintenir les lumières allumées" à "l'ingénierie du chaos" en utilisant le Cluster Linking et des configurations clients résilientes. Elle nécessite de traiter la plateforme comme un produit via le GitOps et le self-service (Platform Engineering). Elle impose une couche de gouvernance qui applique l'intention métier par le biais de Contrats de Données stricts. Enfin, elle positionne la plateforme de streaming comme le cortex cognitif de l'organisation, alimentant

le contexte temps réel pour la GenAI et enregistrant son raisonnement pour la postérité. En adhérant à ces principes architecturaux, Confluent Cloud devient non seulement un courtier de messages, mais le système nerveux central de l'entreprise numérique moderne.

## Ouvrages cités

1. Overview of Cluster Linking Confluent Platform, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/platform/current/multi-dc-deployments/cluster-linking/index.html>
2. Multi-Region Kafka using Synchronous Replication for Disaster Recovery with Zero Data Loss (RPO=0) - Kai Waehner, dernier accès : décembre 12, 2025,  
<https://www.kai-waehner.de/blog/2025/08/04/multi-region-kafka-using-synchronous-replication-for-disaster-recovery-with-zero-data-loss-rpo0/>
3. Kafka Replication and Committed Messages - Confluent Documentation, dernier accès : décembre 12, 2025, <https://docs.confluent.io/kafka/design/replication.html>
4. Tutorial: Multi-Region Clusters on Confluent Platform, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/platform/current/multi-dc-deployments/multi-region-tutorial.html>
5. Planning for Apache Kafka - Cloudera Documentation, dernier accès : décembre 12, 2025,  
<https://docs.cloudera.com/runtime/7.3.1/kafka-planning/kafka-planning.pdf>
6. AWS Kafka - Guide to Design & Kafka Deployment Considerations - Confluent, dernier accès : décembre 12, 2025,  
<https://www.confluent.io/blog/design-and-deployment-considerations-for-deploying-apache-kafka-on-aws/>
7. Configure Multi-Region Clusters in Confluent Platform, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/platform/current/multi-dc-deployments/multi-region.html>
8. Apache Kafka stretch cluster reference architecture | Cloudera on Cloud, dernier accès : décembre 12, 2025,  
<https://docs.cloudera.com/runtime/7.3.1/kafka-planning/topics/kafka-planning-stretch-reference-arch.html>
9. Geo-replication with Cluster Linking on Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/multi-cloud/cluster-linking/index.html>
10. The Hitchhiker's Guide to Disaster Recovery and Multi-Region Kafka : r/apachekafka, dernier accès : décembre 12, 2025,  
[https://www.reddit.com/r/apachekafka/comments/1l81az8/the\\_hitchhikers\\_guide\\_to\\_disaster\\_recovery\\_and/](https://www.reddit.com/r/apachekafka/comments/1l81az8/the_hitchhikers_guide_to_disaster_recovery_and/)
11. The Hitchhiker's Guide to Disaster Recovery and Multi-Region Kafka - WarpStream, dernier accès : décembre 12, 2025,  
<https://www.warpstream.com/blog/the-hitchhikers-guide-to-disaster-recovery-a>

## nd-multi-region-kafka

12. Cluster Linking Disaster Recovery and Failover on Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/multi-cloud/cluster-linking/dr-failover.html>
13. Multi-Geo Replication in Apache Kafka - Confluent, dernier accès : décembre 12, 2025, <https://www.confluent.io/blog/multi-geo-replication-in-apache-kafka/>
14. Chapter 3. Kafka broker configuration tuning - Red Hat Documentation, dernier accès : décembre 12, 2025,  
[https://docs.redhat.com/en/documentation/red\\_hat\\_streams\\_for\\_apache\\_kafka/2.6/html/kafka\\_configuration\\_tuning/con-broker-config-properties-str](https://docs.redhat.com/en/documentation/red_hat_streams_for_apache_kafka/2.6/html/kafka_configuration_tuning/con-broker-config-properties-str)
15. Kafka Capacity planning : r/apachekafka - Reddit, dernier accès : décembre 12, 2025,  
[https://www.reddit.com/r/apachekafka/comments/1pbjbaz/kafka\\_capacity\\_planning/](https://www.reddit.com/r/apachekafka/comments/1pbjbaz/kafka_capacity_planning/)
16. SQL CREATE MODEL Statement in Confluent Cloud for Apache Flink, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/flink/reference/statements/create-model.html>
17. Kafka Consumer Configuration Reference for Confluent Platform, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/platform/current/installation/configuration/consumer-configs.html>
18. Kafka 3.8 Documentation, dernier accès : décembre 12, 2025,  
<https://kafka.apache.org/38/documentation/>
19. Data Contracts for Schema Registry on Confluent Platform, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/platform/current/schema-registry/fundamentals/data-contracts.html>
20. Broker offline and client failover - Amazon Managed Streaming for Apache Kafka, dernier accès : décembre 12, 2025,  
<https://docs.aws.amazon.com/msk/latest/developerguide/troubleshooting-offline-broker-clientfailover.html>
21. Multi-Cluster Fault-Tolerant Kafka Producer | by Sharad - Medium, dernier accès : décembre 12, 2025,  
<https://medium.com/@sharadb7f72044/multi-cluster-fault-tolerant-kafka-producer-c7e7b3f72044>
22. Best practices for right-sizing your Apache Kafka clusters to optimize performance and cost, dernier accès : décembre 12, 2025,  
<https://aws.amazon.com/blogs/big-data/best-practices-for-right-sizing-your-apache-kafka-clusters-to-optimize-performance-and-cost/>
23. Use AWS Transit Gateway with Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/networking/aws-transit-gateway.html>
24. Create an AWS PrivateLink connection to Confluent Cloud, dernier accès : décembre 12, 2025,

<https://docs.confluent.io/cloud/current/networking/private-links/aws-privatelink.html>

25. Use Azure Private Link for Dedicated Clusters on Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/networking/private-links/azure-privatelink.html>
26. Keeping your Kubernetes Backup Data Private with Azure Private Link - CloudCasa, dernier accès : décembre 12, 2025,  
<https://cloudcasa.io/blog/kubernetes-backup-data-kept-private-with-azure-private-link/>
27. Use the Confluent Cloud Console with Private Networking, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/networking/ccloud-console-access.html>
28. Route 53 Resolver endpoints and forwarding rules - Hybrid Cloud DNS Options for Amazon VPC - AWS Documentation, dernier accès : décembre 12, 2025,  
<https://docs.aws.amazon.com/whitepapers/latest/hybrid-cloud-dns-options-for-vpc/route-53-resolver-endpoints-and-forwarding-rules.html>
29. Forwarding inbound DNS queries to your VPCs - Amazon Route 53, dernier accès : décembre 12, 2025,  
<https://docs.aws.amazon.com/Route53/latest/DeveloperGuide/resolver-forwarding-inbound-queries.html>
30. Strangler Fig Pattern - Azure Architecture Center | Microsoft Learn, dernier accès : décembre 12, 2025,  
<https://learn.microsoft.com/en-us/azure/architecture/patterns/strangler-fig>
31. Strangler Fig - Confluent Developer, dernier accès : décembre 12, 2025,  
<https://developer.confluent.io/patterns/compositional-patterns/strangler-fig/>
32. CIO Guide: Strangler Fig Pattern Architecture | Enterprise Modernization | Salfati Group, dernier accès : décembre 12, 2025,  
<https://salfati.group/topics/strangler-fig-pattern>
33. Advanced Strategies for Strangler Fig Pattern Implementation - Capstone IT Solutions, dernier accès : décembre 12, 2025,  
<https://capstone-s.com/strangler-fig-pattern-implementation/>
34. Change Data Streaming Patterns With Debezium & Apache Flink - Speaker Deck, dernier accès : décembre 12, 2025,  
<https://speakerdeck.com/gunnarmorling/change-data-streaming-patterns-with-debezium-and-apache-flink>
35. Modernizing a Legacy Application Using the Strangler Fig Pattern - Curotec, dernier accès : décembre 12, 2025,  
<https://www.curotec.com/insights/modernizing-a-legacy-application-using-the-strangler-fig-pattern/>
36. confluent\_kafka\_cluster | Resources | confluentinc/confluent - Terraform Registry, dernier accès : décembre 12, 2025,  
[https://registry.terraform.io/providers/confluentinc/confluent/latest/docs/resources/confluent\\_kafka\\_cluster](https://registry.terraform.io/providers/confluentinc/confluent/latest/docs/resources/confluent_kafka_cluster)
37. Manage Kafka with Terraform: Why & How - AutoMQ, dernier accès : décembre

12, 2025,

<https://www.automq.com/blog/manage-kafka-with-terraform-why-and-how>

38. confluent\_private\_link\_attachment | Resources | confluentinc/confluent - Terraform Registry, dernier accès : décembre 12, 2025,  
[https://registry.terraform.io/providers/confluentinc/confluent/latest/docs/resource/confluent\\_private\\_link\\_attachment](https://registry.terraform.io/providers/confluentinc/confluent/latest/docs/resource/confluent_private_link_attachment)
39. Build an Event Streaming Platform with GitOps (Exercise) - Confluent Developer, dernier accès : décembre 12, 2025,  
<https://developer.confluent.io/courses/data-streaming-systems/build-an-event-streaming-platform-with-gitops-exercise/>
40. Best Practices for Confluent Terraform Provider, dernier accès : décembre 12, 2025,  
<https://www.confluent.io/blog/best-practices-confluent-terraform-provider/>
41. Is monorepo or polyrepo better for large-scale applications? - Graphite, dernier accès : décembre 12, 2025,  
<https://graphite.com/guides/monorepo-vs-polyrepo-large-scale-applications>
42. Monorepos vs. Polyrepos: Which one fits your use case? - LogRocket Blog, dernier accès : décembre 12, 2025,  
<https://blog.logrocket.com/monorepos-vs-polyrepos-which-one-fits-your-use-case/>
43. The GitOps Repository Structure: Monorepo vs. Polyrepo and Best Practices | by Rahul Kumar Singh | Google Cloud - Medium, dernier accès : décembre 12, 2025,  
<https://medium.com/google-cloud/the-gitops-repository-structure-monorepo-vs-polyrepo-and-best-practices-17399ae6f3f4>
44. How to set up your GitOps directory structure - Red Hat Developer, dernier accès : décembre 12, 2025,  
<https://developers.redhat.com/articles/2022/09/07/how-set-your-gitops-directory-structure>
45. Terraform Provider for Confluent - GitHub, dernier accès : décembre 12, 2025,  
<https://github.com/confluentinc/terraform-provider-confluent>
46. devshawn/kafka-gitops: Manage Apache Kafka topics and generate ACLs through a desired state file. - GitHub, dernier accès : décembre 12, 2025,  
<https://github.com/devshawn/kafka-gitops>
47. Building a Centralized Cloud Resource Management Hub with Backstage Plugins - Firefly, dernier accès : décembre 12, 2025,  
<https://www.firefly.ai/academy/building-a-centralized-cloud-resource-management-hub-with-backstage-plugins>
48. Building a Self-Service Confluent Developer Platform with Backstage and Terraform, dernier accès : décembre 12, 2025,  
<https://danielwessendorf.com/blog/2025-03-16-confluent-self-service-with-backstage-part-1/>
49. Kafka Dead Letter Queue Best Practices: Examples, Retries, and Monitoring - Superstream, dernier accès : décembre 12, 2025,  
<https://www.superstream.ai/blog/kafka-dead-letter-queue>
50. Kafka Connect Deep Dive – Error Handling and Dead Letter Queues | Confluent,

dernier accès : décembre 12, 2025,

<https://www.confluent.io/blog/kafka-connect-deep-dive-error-handling-dead-letter-queues/>

51. Building a Dead Letter Queue (DLQ) for Reliable Kafka-to-Database Pipelines - Medium, dernier accès : décembre 12, 2025,  
<https://medium.com/@sonal.sadafal/building-a-dead-letter-queue-dlq-for-reliable-kafka-to-database-pipelines-a4b0f11b125a>
52. 3 Best Practices to Effectively Manage Failed Messages - CrowdStrike, dernier accès : décembre 12, 2025,  
<https://www.crowdstrike.com/en-us/blog/three-best-practices-to-effectively-manage-failed-messages/>
53. Solving Event Losses with Kafka Dead Letter Queue | by Orçun Yılmaz - Medium, dernier accès : décembre 12, 2025,  
<https://medium.com/@orcunyilmazoy/solving-event-losses-with-kafka-dead-letter-queue-d47538b27697>
54. Using Data Contracts with Confluent Schema Registry, dernier accès : décembre 12, 2025,  
<https://www.confluent.io/blog/data-contracts-confluent-schema-registry/>
55. Track Data with Stream Lineage on Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/stream-governance/stream-lineage.html>
56. Deployment Environment Variables - DataHub, dernier accès : décembre 12, 2025, <https://docs.datahub.com/docs/deploy/environment-vars>
57. Confluent Cloud API Reference Documentation, dernier accès : décembre 12, 2025, <https://docs.confluent.io/cloud/current/api.html>
58. Kafka Connect MaskField SMT Usage Reference for Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/connectors/transforms/maskfield.html>
59. Configure Single Message Transforms for Kafka Connectors in Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/connectors/single-message-transforms.html>
60. Architectural considerations for streaming applications on Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://docs.confluent.io/cloud/current/client-apps/architecture.html>
61. Real-Time GenAI with RAG using Apache Kafka and Flink to Prevent Hallucinations, dernier accès : décembre 12, 2025,  
<https://www.kai-waehner.de/blog/2024/05/30/real-time-genai-with-rag-using-apache-kafka-and-flink-to-prevent-hallucinations/>
62. What is Retrieval Augmented Generation (RAG)? - Confluent, dernier accès : décembre 12, 2025,  
<https://www.confluent.io/learn/retrieval-augmented-generation-rag/>
63. How to generate vector embeddings for RAG with Flink SQL in Confluent Cloud, dernier accès : décembre 12, 2025,  
<https://developer.confluent.io/confluent-tutorials/gen-ai-vector-embedding/flinks>

ql/

64. Add RAG to Your Flink AI Flow Using Vector Search with Pinecone - Confluent, dernier accès : décembre 12, 2025,  
<https://www.confluent.io/blog/flink-ai-rag-with-federated-search/>
65. Introducing Real-Time Embeddings: Any Model, Any Vector Database - Confluent, dernier accès : décembre 12, 2025,  
<https://www.confluent.io/blog/flink-action-create-vector-embeddings/>
66. The Ultimate 2025 LLM Cheat-sheet | by Kafkafranz - Medium, dernier accès : décembre 12, 2025,  
<https://medium.com/@kafkafranz495/the-ultimate-2025-llm-handbook-4a545d2f5c59>
67. Security Monitoring for AI Agents and MCP, dernier accès : décembre 12, 2025,  
<https://realm.security/security-monitoring-for-ai-agents-and-mcp/>
68. Agent 365 observability schema reference | Microsoft Learn, dernier accès : décembre 12, 2025,  
<https://learn.microsoft.com/en-us/microsoft-agent-365/developer/reference/observability-schema/>
69. Observability with OpenTelemetry - Gemini CLI, dernier accès : décembre 12, 2025, <https://geminicli.com/docs/cli/telemetry/>
70. A Measurable AGI-Oriented Framework with a Modular Cognitive Architecture and Early Prototype Demonstration An AGI-inspired cognitive system demonstrating autonomous learning, reasoning, and cross-domain transfer in a prototype environment - ResearchGate, dernier accès : décembre 12, 2025,  
[https://www.researchgate.net/publication/397916420\\_A\\_Measurable\\_AGI-Oriented\\_Framework\\_with\\_a\\_Modular\\_Cognitive\\_Architecture\\_and\\_Early\\_Prototype\\_Demonstration\\_An\\_AGI-inspired\\_cognitive\\_system\\_demonstrating\\_autonomous\\_learning\\_reasoning\\_and\\_cross-domai](https://www.researchgate.net/publication/397916420_A_Measurable_AGI-Oriented_Framework_with_a_Modular_Cognitive_Architecture_and_Early_Prototype_Demonstration_An_AGI-inspired_cognitive_system_demonstrating_autonomous_learning_reasoning_and_cross-domai)
71. Utilizing Generative AI and LLMs to Automate Detection Writing | by Dylan - Medium, dernier accès : décembre 12, 2025,  
<https://medium.com/@dylanhwilliams/utilizing-generative-ai-and-llms-to-automate-detection-writing-5e4ea074072e>
72. Think Inside the JSON: Reinforcement Strategy for Strict LLM Schema Adherence - arXiv, dernier accès : décembre 12, 2025, <https://arxiv.org/html/2502.14905v1>
73. Structured model outputs | OpenAI API, dernier accès : décembre 12, 2025,  
<https://platform.openai.com/docs/guides/structured-outputs>
74. Structured Output in LLMs: Why JSON/XML Format Matters | by Tahir | Medium, dernier accès : décembre 12, 2025,  
<https://medium.com/@tahirbalarabe2/%EF%88%8Fstructured-output-in-llms-why-json-xml-format-matters-c644a81cf4f3>