

## □ Architecture moderne Big Data Analytics

### Nouvelles architectures analytiques

#### ◆ Lambda Architecture

- **Combine batch (Spark SQL, ETL périodique) + streaming (Structured Streaming, Kafka)**
- **Avantage : précision + temps réel**
- **Spark 3+ simplifie la gestion du double flux avec le même code DataFrame**

#### ◆ Kappa Architecture

- **Tout passe par le flux temps réel (Kafka + Spark Streaming)**
- **Le batch est remplacé par la relecture du topic Kafka**
- **Avantage : plus simple, un seul pipeline ; support natif dans Spark 3+**

## □ Ingestion & Intégration des données

Collecter les données depuis les sources (bases, API, fichiers, SaaS, Kafka, etc.)

Catégorie	Outils	Rôle
Ingestion batch & streaming	<b>Airbyte, Fivetran, Kafka, NiFi, Flume</b>	Connecter les sources et charger les données dans un data lake/warehouse
Message broker / streaming	<b>Apache Kafka</b>	Transport temps réel, file d'attente distribuée

## □ Stockage & Data Lake

Lieu où sont stockées les données brutes et transformées.

Catégorie	Outils	Rôle
Object storage	<b>MinIO, Amazon S3, Azure Blob Storage</b>	Stocker de gros volumes de données non structurées (parquet, csv, json...)
Data Lakehouse	<b>Dremio, Databricks, Snowflake, Delta Lake, Iceberg</b>	Combiner lac + entrepôt (SQL + performance)

## □ Orchestration & Workflow

Gérer les dépendances, planifier et exécuter les pipelines de données.

Catégorie	Outils	Rôle
Orchestration	<b>Apache Airflow, Dagster, Prefect</b>	Automatiser les pipelines (exécuter ETL, transformations, contrôles)

## 4 Transformation & Modélisation

Transformer les données brutes en données prêtes pour l'analyse.

Catégorie	Outils	Rôle
Data transformation (ELT)	<b>dbt (Data Build Tool)</b>	Gérer les transformations SQL dans l'entrepôt, versionner, documenter, tester
Compute engine	<b>Spark, Flink, Dremio</b>	Exécuter les transformations à grande échelle

## 5 Data Quality & Testing

Garantir la fiabilité et la cohérence des données.

Catégorie	Outils	Rôle
Contrôle qualité	<b>Great Expectations, dbt tests, Soda, Monte Carlo</b>	Vérifier la cohérence, la fraîcheur, la complétude des données

## 6 Gouvernance & Catalogage

Assurer la conformité, la traçabilité, et la documentation des données.

Catégorie	Outils	Rôle
Data catalog	<b>Collibra, Alation, DataHub, Amundsen</b>	Centraliser la documentation et la découverte des données
Gouvernance	<b>Apache Atlas, OpenMetadata, DataHub</b>	Gérer les politiques, lineage, conformité RGPD, rôles et sécurité

---

## 7 Data Consumption (Consommation)

Mise à disposition des données aux utilisateurs métiers, data scientists, BI, IA.

Catégorie	Outils	Rôle
BI / visualisation	<b>Power BI, Tableau, Superset, Metabase</b>	Créer des dashboards et indicateurs
Machine Learning	<b>Python (pandas, scikit-learn), Spark ML, TensorFlow, Dremio</b>	Exploiter les données pour la prédiction et l'IA

## 8 Sécurité & Conformité (Security Layer)

Protéger les données, contrôler les accès et assurer la conformité réglementaire.

Catégorie	Outils / Technologies	Rôle
Authentification & IAM	<b>LDAP, Kerberos, OAuth2, Azure AD, AWS IAM</b>	Gérer l'identité et l'accès des utilisateurs
Chiffrement & Sécurité des données	<b>Ranger, Knox, Vault, SSL/TLS, KMS</b>	Chiffrement au repos et en transit, gestion des clés
Contrôle d'accès & Audits	<b>Apache Ranger, Atlas, DataHub</b>	Politiques RBAC, ABAC, logs d'accès, audit
Conformité & RGPD	<b>Collibra, DataHub, OneTrust</b>	Masquage des données sensibles, consentement, traçabilité
Monitoring & Sécurité réseau	<b>Prometheus, Grafana, ELK, SIEM</b>	Surveiller les anomalies, intrusion, performance

Étape	Catégorie	Exemples
1. Ingestion	Airbyte, Kafka, NiFi	Collecte des données
2. Stockage	MinIO, Dremio, S3	Data lake / lakehouse
3. Orchestration	Airflow	Planification et exécution
4. Transformation	dbt, Spark	Modélisation analytique
5. Qualité	Great Expectations, Soda	Vérification des données
6. Gouvernance	DataHub, Collibra	Traçabilité, catalogage
7. Consommation	Power BI, Tableau	Visualisation, IA

Étape	Catégorie	Exemples
8. Sécurité	Ranger, Knox, Atlas	Chiffrement, contrôle d'accès, conformité

## Flux visuel (mise à jour avec sécurité)

### Sources (API, fichiers CSV, logs)



**Kafka** (*streaming temps réel*)



**Airbyte / Airflow** (*ingestion + orchestration*)



**MinIO / S3** (*data lake sécurisé – chiffrement & IAM*)



**DBT + Dremio / Spark** (*transformation & analyse*)



**Contrôles qualité & gouvernance** (*Great Expectations + DataHub*)



**Sécurité & conformité** (*Ranger, Knox, Atlas*)



**Power BI / Superset / ML** (*consommation des données fiables et sécurisées*)