

ia.rbre

Mika INISAN

LIRIS

23 octobre 2025

Synthèse : Data pipelines

Chaînes de traitement de données

MÉTROPOLE
GRAND **LYON**

Teles**Coop**

université
LUMIÈRE
LYON 2



 **BANQUE des**
TERRITOIRES |  Caisse
des Dépôts
GROUPE

Sommaire

- 01 Objectifs de la présentation
- 02 Objectifs de la synthèse
- 03 Méthode de recherche employée
- 04 Rappels sur IA.rbre et ses objectifs
- 05 Notions sur les data pipelines
- 06 Meilleures pratiques
- 07 Data pipeline pour IA.rbre



01

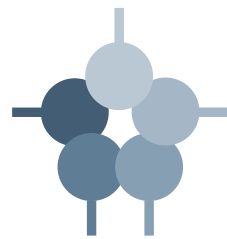
Objectifs de la présentation

01 → 02 → ... → 07

Mika INISAN

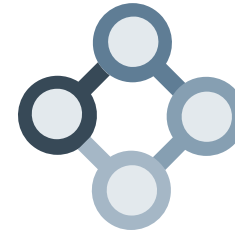
LIRIS

Objectifs de la présentation



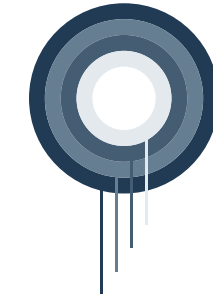
Restituer les principaux enseignements de la synthèse :

approches, défis, tendances,
meilleures pratiques...



Partager une base commune de connaissance :

enjeux, problématiques
techniques et stratégiques...



Être un point de départ favorisant les échanges pour ajuster et affiner le projet :

questions, idées, compléments,
approfondissements,
confrontation avec les attentes...



02

Objectifs de la synthèse

Objectifs de la synthèse



Positionner la solution développée par rapport à l'état de l'art (convergence, divergence, opportunités d'ajustement, axes différenciants)



Tirer parti des travaux existants pour accélérer les développements et éviter de repartir de zéro



Identifier des leviers d'innovation et fonctionnalités, susciter des idées nouvelles



Réduire l'incertitude, éclairer les décisions, valider et justifier les choix a priori et a posteriori



Anticiper les futurs besoins, démontrer l'impact

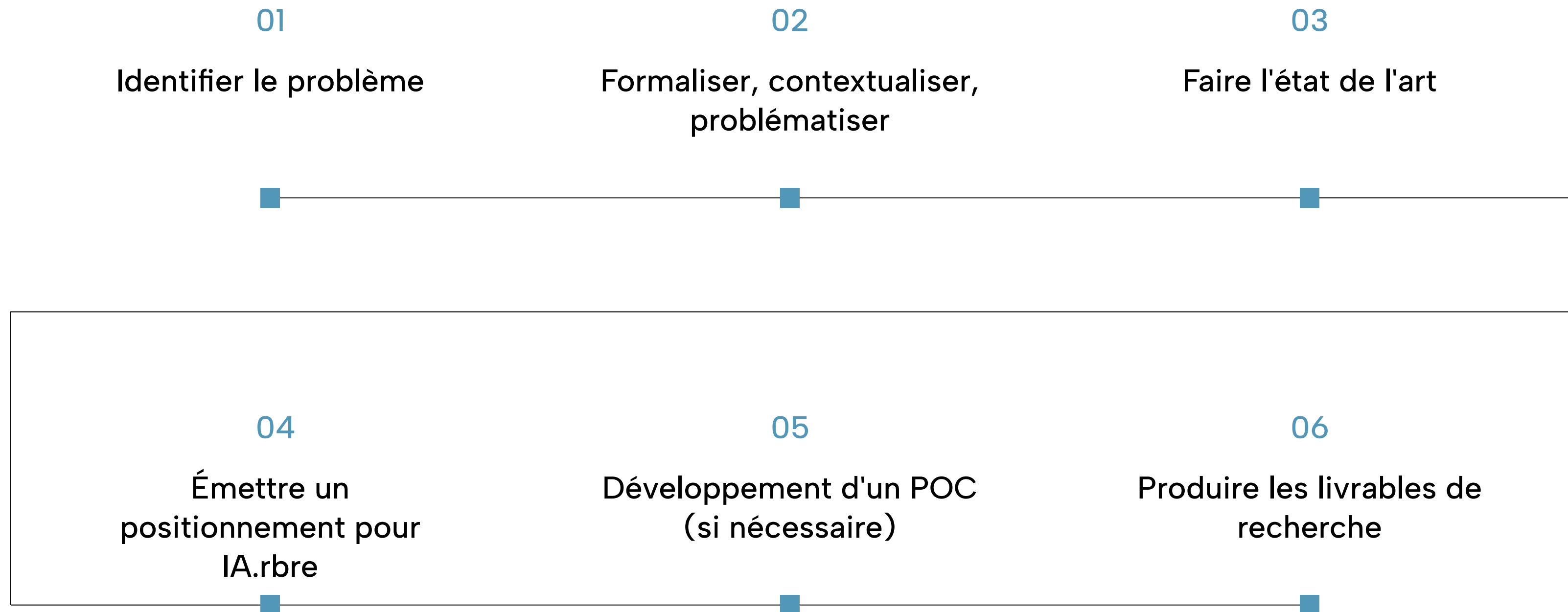


03

Méthode de recherche employée

01 → 02 → **03** → 04 → ... → 07

Méthode de recherche employée





04

Rappels sur l'A.rbre et ses objectifs₇

Rappels sur IA.rbre et ses objectifs

Rappels sur IA.rbre

- **Finalités :**
 - Aider à la décision
 - Servir les services métiers et le grand public grâce aux besoins généralisables pris en compte
- **Propriété intellectuelle :** commun numérique
- **Formes :**
 - Outils de visualisation et de scénarisation
 - Produit de données
- **Adjectifs :** frugalité, innovation, basé sur la science, open data / source / innovation / science, données FAIR

Rappels sur IA.rbre et ses objectifs

Objectifs fonctionnels

- Intégrer les données territoriales disponibles pour agrandir la connaissance et la pertinence de l'analyse :

Sources : acteurs multiples du territoire, publics et privés.

DataGrandLyon, IGN, exploitant de réseaux (TCL, ENEDIS, SFR...)...

Contenu : données de réseaux, de chantiers, d'occupation de sols, inventaire du végétal urbain existant...

Moyen d'acquisition : en ligne, issues d'une démarche DICT, locales à la métropole...

Types : modèles, données territoriales, vérité terrain, données issues de modèles...

Formats : raster, vectoriel, GeoTiff, Shapefile, GeoPackage, CityGML, CSV, HDF5, MLmodel...

- Transférer la connaissance d'une métropole ayant beaucoup de données, sur une métropole en ayant moins

01 → ... → 03 → **04** → 05 → ... → 07

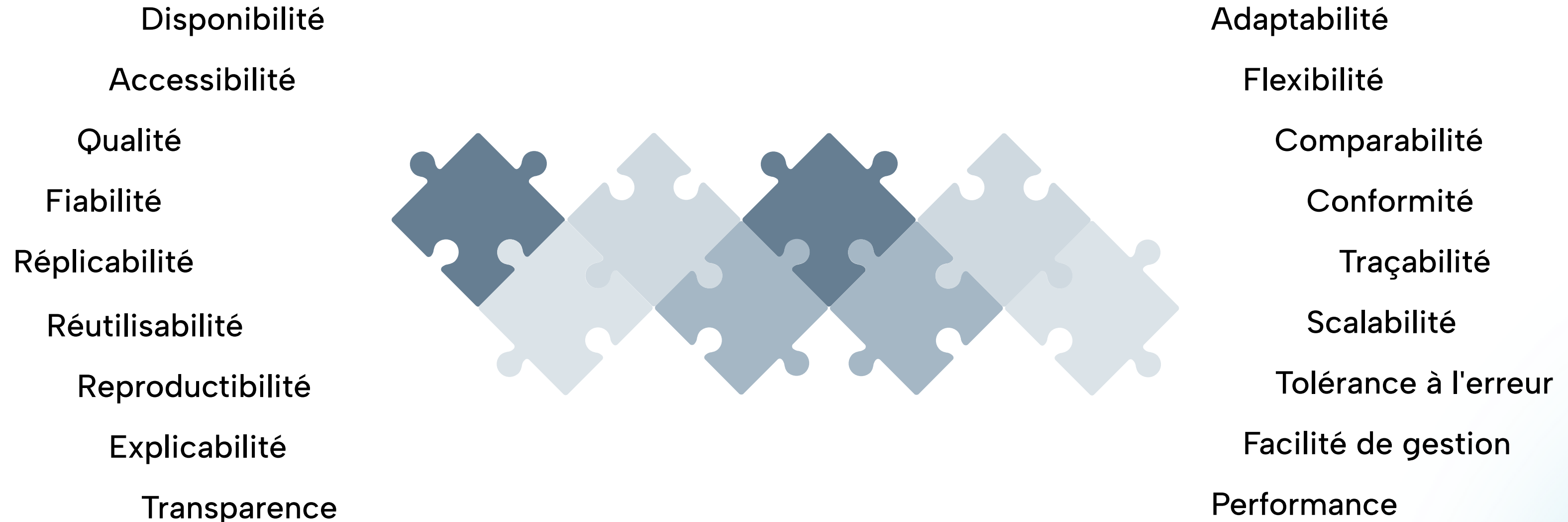
Rappels sur IA.rbre et ses objectifs

Objectifs fonctionnels

- Faire des analyses complexes dont ML, AI :
 - Amélioration des données par croisement
 - Génération de nouvelles données à partir de celles existantes
 - Réunification, réconciliation des données
 - Croisement des données
 - Substituer les données pour le calcul des facteurs
 - Calcul des facteurs de faisabilité
 - Transformer les données ingérées en indices intermédiaires
 - Pondération possiblement automatique des facteurs
 - Génération des calques basiques (plantabilité) et thématiques
 - Exploration de scénario, données dans le temps, évolution temporelle des villes
 - A/B testing de modèles

Rappels sur IA.rbre et ses objectifs

Objectifs non fonctionnels





05

Notions sur les data pipelines

01 → ... → 04 → **05** → 06 → 07

Notions sur les data pipelines

Introduction

Structure logicielle qui permet le déplacement et la manipulation systématiques de données provenant de sources potentiellement multiples et hétérogènes, vers des destinations variées.

Se compose d'un ensemble d'étapes, possiblement automatisées, interconnectées à travers lesquelles les données passent de manière séquentielle, la sortie d'une étape servant d'entrée pour la suivante.



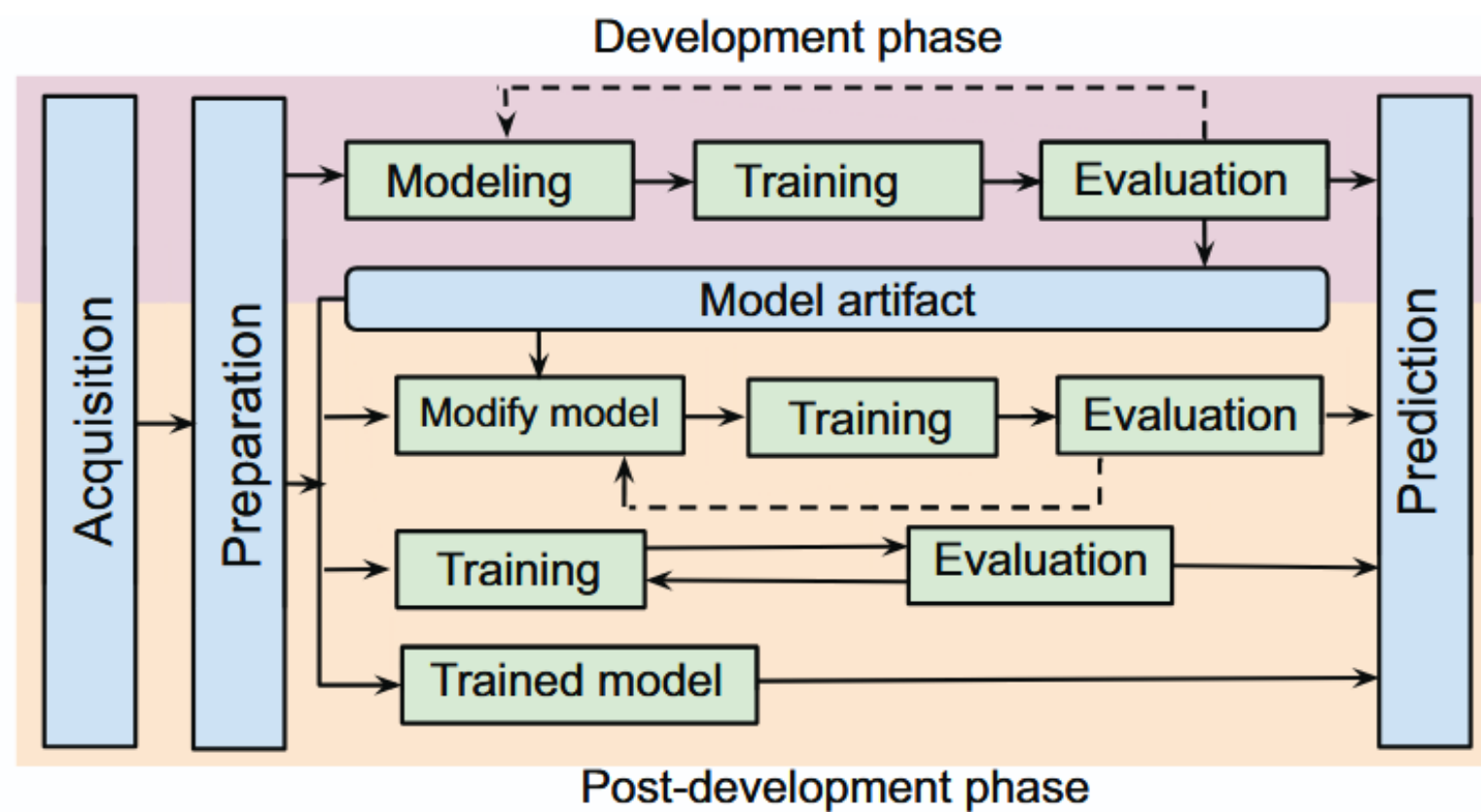
Pas de standardisation et de terminologie largement acceptée dans la littérature scientifique [1]

[1] S. Biswas, M. Wardat, and H. Rajan, "The art and practice of data science pipelines," *Proceedings of the 44th International Conference on Software Engineering*, pp. 2091–2103, May 2022

01 → ... → 04 → **05** → 06 → 07

Notions sur les data pipelines

Exemple et représentation



Représentation : graphe orienté acyclique
(Directed acyclic graph, DAG)

Exemple que nous détaillerons par la suite

Notions sur les data pipelines

Caractéristiques

Formats des données

Non structurés, semi-structurés, structurés

Sources des données

Locales, distantes ; centralisées, distribuées ; API, fichiers, base de données, origines d'un crawling ou scraping...

Destinations

Stockage, applications (outils de visualisation...), autre data pipelines, modèles d'apprentissage automatique...

Modes de déclenchement

Manuel, programmé ; ponctuel, récurrent, en réponse à des stimuli basés sur des événements

Modes d'ingestion de données

Par lots (« batch »), continu (« streaming »), hybride

Modes d'exécution

Séquentiel, parallèle ; centralisé, distribué

Notions sur les data pipelines

Étapes

Extraction	Transformation	Chargement	Filtrage	Fusion
Sélection	Ingestion	Acquisition	Exploration	Agrégation
Validation	Enrichissement	Stockage	Visualisation	...

Étapes spécifiques à l'analyse (ML, AI...) :

Feature engineering	Modélisation	Entraînement	Évaluation	...
---------------------	--------------	--------------	------------	-----

Notions sur les data pipelines

Rôles et objectifs

Les chaînes de traitement de données, ou data pipelines, occupent une place de plus en plus centrale au sein des systèmes de gestion de données contemporains grâce à leur **rôle structurant**. Elles constituent la fondation des activités de traitement, d'analyse et de prise de décision.

Elles permettent de **contrôler toutes les opérations** liées aux données et d'**orchestrer l'ensemble du flux** (traitement, transfert, stockage) de manière rationalisée de la source à la destination.

En **remplaçant des opérations manuelles** susceptibles d'introduire des erreurs **par des processus systématiques** et reproductibles.

Elles **décomposent les analyses complexes** de grands ensembles de données **en une série de tâches plus simples**.

Les propriétés des implémentations des différentes étapes visent à encourager la réutilisation, la composition flexible et la configurabilité pour des usages spécifiques.

Notions sur les data pipelines

Rôles et objectifs

- Atténuer les goulots d'étranglement et les délais
- Renforcer la scalabilité (notamment face au volume de données)
- Simplifier la conception et le déploiement des services de traitement des données
- Reproductibilité, traçabilité, fiabilité, tolérance aux pannes
- Améliorer la cohérence des systèmes d'information
- Augmenter la vitesse de bout en bout, l'efficacité
- Réduire la latence dans le développement des produits de données
- Accès temps réel à l'information actualisée
- Une meilleure prise de décision



06

Meilleures pratiques

Meilleures pratiques



- Contrôle de la dette technique
- Évolutivité, flexibilité, robustesse
- Faciliter maintenance, entretien, dépannage, documentation...
- Meilleure testabilité, sécurisabilité, reproductibilité, réutilisabilité



Concevoir de manière modulaire (séparer, isoler clairement les responsabilités au sein des étapes, de minimiser les redondances)

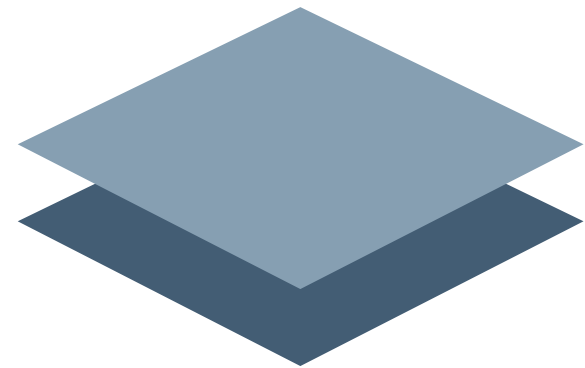



Une attention particulière doit être portée sur l'interface entre les étapes



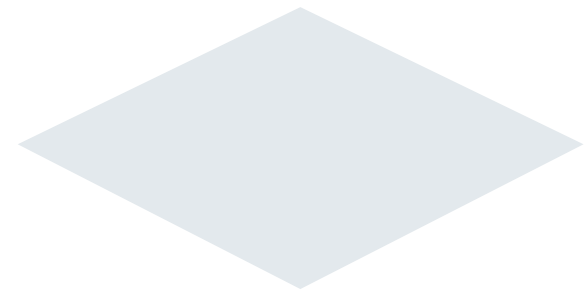
L'utilisation d'un métamodèle peut simplifier la phase de conception

Meilleures pratiques



- Recentrer l'effort sur des missions à plus forte valeur ajoutée que celles de gestion de données
 - Minimiser l'erreur humaine, la charge de travail et les coûts liés au volume, la vitesse et la variété des données
-
-  Maximiser l'automatisation des étapes et tâches

Meilleures pratiques



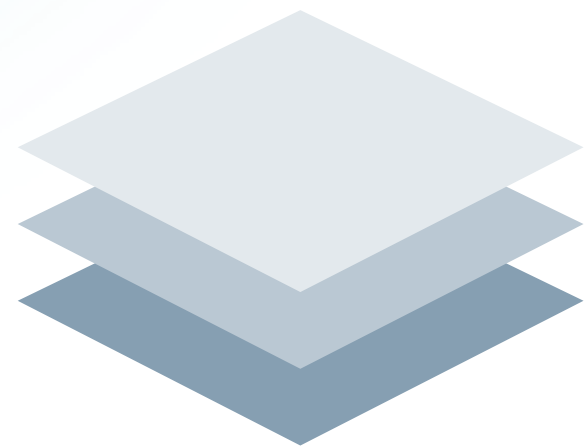
- Maximiser la tolérance aux défaillances et le recouvrement (matérielles, algorithmiques, erreurs métier)



Mettre en place :

- Un monitoring (latence, vitesse de transfert, taux d'erreur ; logs d'erreurs, rapports, tableaux de bord...)
- Une détection de défaillance
- Des règles de validation et de réjection entière ou partielle (exemple : espace de staging)
- Des stratégies d'atténuation (exemple : versionnement)
- La levée d'alerte et d'événements automatiques

Meilleures pratiques



- Reproductibilité, traçabilité
- Pouvoir d'appliquer des traitements différents à l'avenir (e.g., comparaisons)
- Éviter de refaire l'intégralité des traitements en cas d'erreur ou de problème de qualité intermédiaire



Stocker les données initiales et les données intermédiaires issues de transformation en plus des données finales

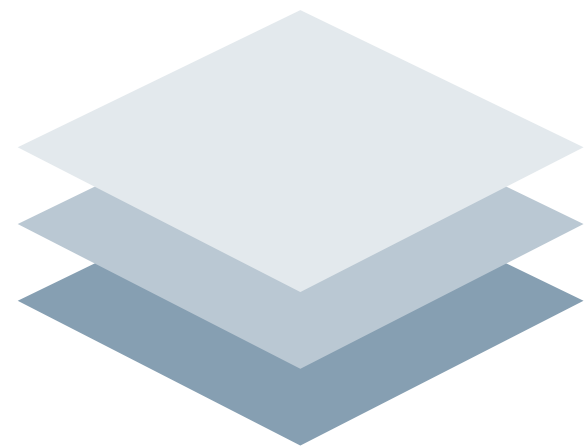


Séparer les stockages initiaux, intermédiaires et finaux



Minimiser la duplication qui amène à des complexités de gestion et une incertitude sur l'état des données comme l'actualisation

Meilleures pratiques



- Minimiser les coûts d'ingestion et maximiser la frugalité
- Traçabilité
- Éviter l'explosion des temps d'exécution, de l'espace de stockage pour des données qui sont potentiellement inutilisables

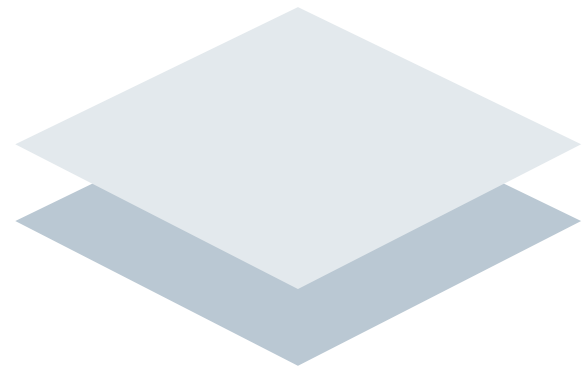


Faire un suivi rigoureux de l'origine des données (e.g, via des métadonnées)



Faire du data profiling pour choisir précisément les données à intégrer, par exemple, les données les plus récentes par rapport à celles déjà intégrées (intégration incrémentale)

Meilleures pratiques



- Être flexible à la variété des données (formats, structures...)
- Supporter l'évolution



Choisir une architecture qui privilégie une compatibilité et une évolution automatique du schéma



Choisir des outils pérennes dans le temps, faciles d'utilisation, faciles à l'intégration, offrant un bon monitoring, sécurisés et scalables

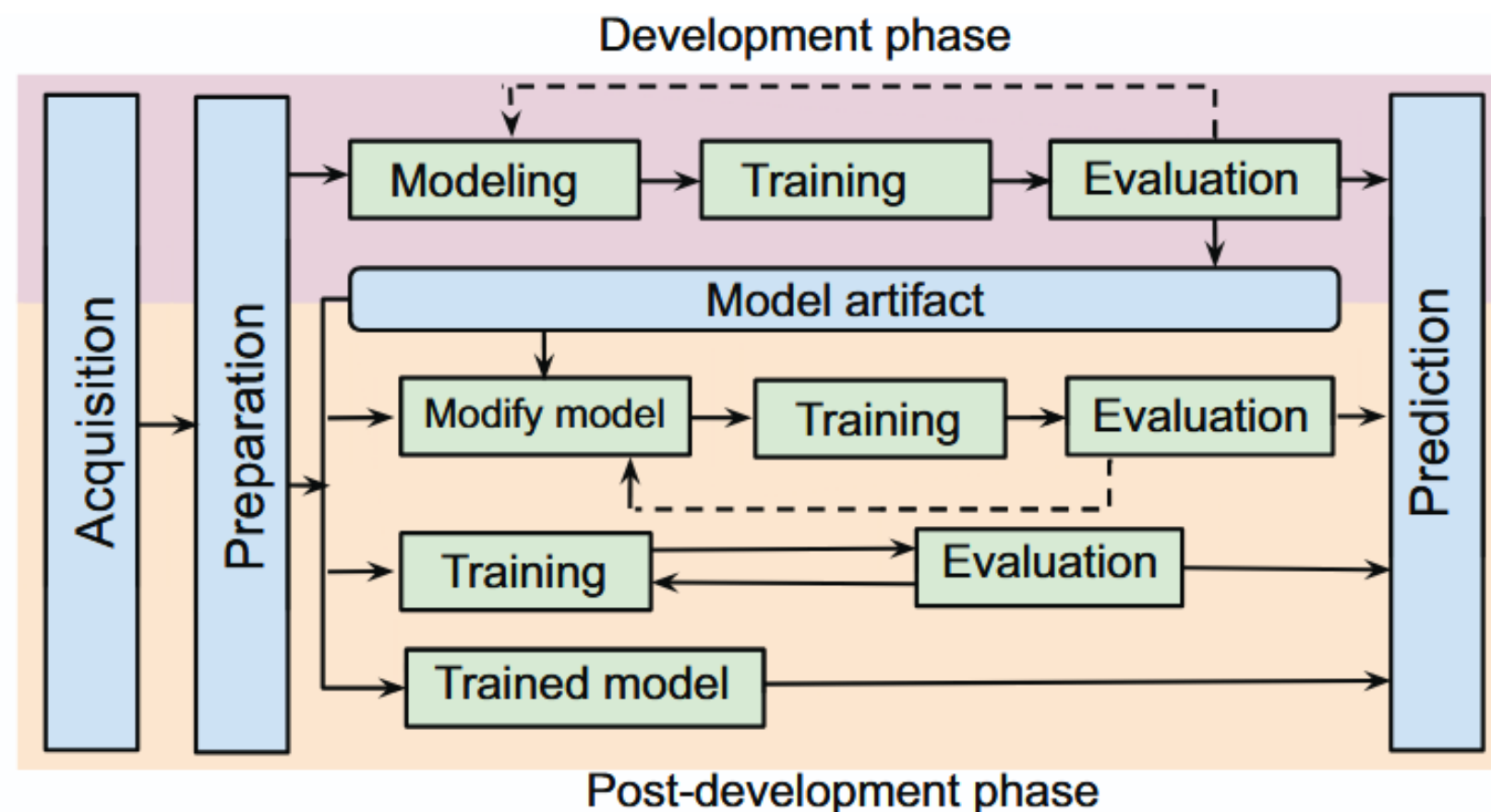


07

Data pipeline pour IA.rbre

Data pipeline pour IA.rbre

Architecture pour la répliquabilité



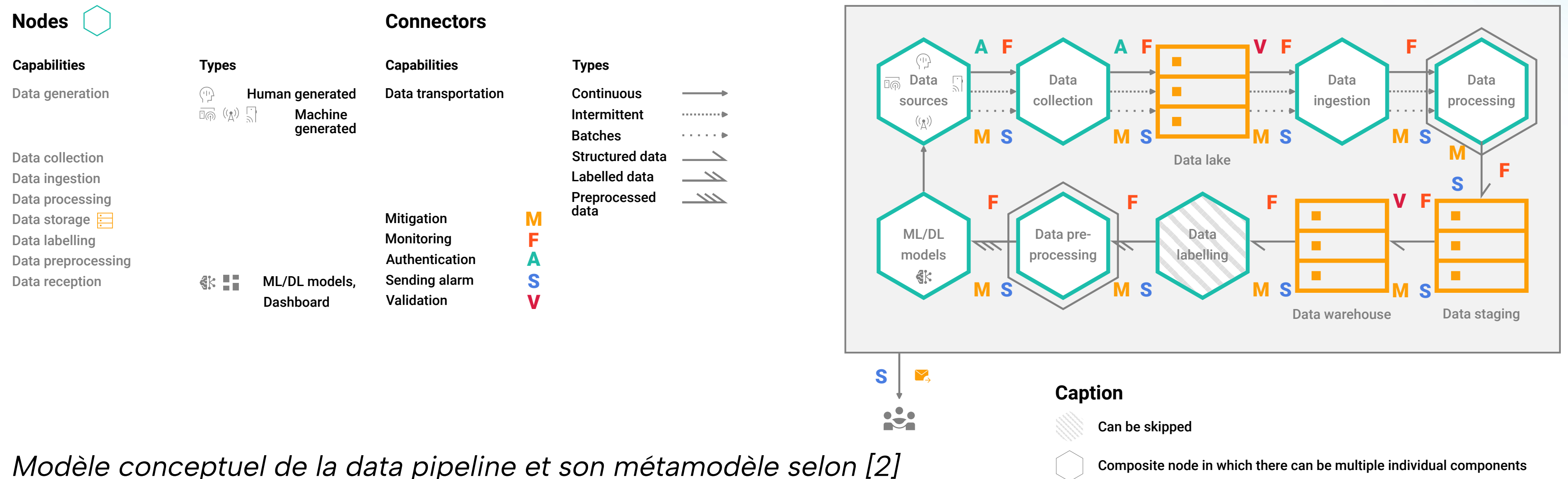
- **Cas de données abondantes :** entraîner les modèles sur des données prétraitées, stocker les artifacts afin de prédire les facteurs à partir de ceux-ci
- **Cas de données de faible qualité, insuffisantes :** réutiliser directement des modèles déjà entraînés sur une autre ville pour effectuer des prédictions
- **Cas intermédiaire :** surentraîner des modèles existants sur les données locales, puis effectuer la prédiction

Architecture représentative des projets d'envergure [1]

[1] S. Biswas, M. Wardat, and H. Rajan, "The art and practice of data science pipelines," *Proceedings of the 44th International Conference on Software Engineering*, pp. 2091–2103, May 2022

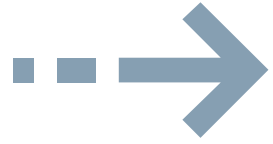
Data pipeline pour IA.rbre

Architecture pour la tolérance aux failles



Modèle conceptuel de la data pipeline et son métamodèle selon [2]

[2] A. Raj, J. Bosch, H. H. Olsson, and T. J. Wang, "Modelling Data Pipelines," 46th Euromicro Conference on Software Engineering and Advanced Applications (SEAA), pp. 13–20, Aug. 2020



Des questions ? Des idées ?

- Synthèse complète sur le *GitHub*



VCityTeam/UD-IArbre-Research/data-pipelines

- Partagez vos retours directement via les Issues *GitHub*

Merci !

License (inclut les conditions relatives aux images) :
[VCityTeam/UD-IArbre-Research/data-pipelines/LICENSE.md](#)

CREDITS: This presentation template was created by **Slidesgo**,
has been recreated and modified in Inkscape by Mika Inisan and
includes icons by **Flaticon** and infographics & images by **Freepik**.