

Apprentissage fédéré et analyse décentralisée des données pour une recherche collaborative en santé à l'échelle nationale et internationale

Livre blanc du
RLS-Digital Health







Introduction	4
La recherche fondée sur les données dans le domaine de la santé : perspectives et défis	7
Les progrès de la recherche fondée sur les données au sein des institutions et des juridictions	7
Défis de la recherche fondée sur les données à l'échelle nationale et internationale	8
L'apprentissage fédéré et l'analyse décentralisée des données de santé : perspectives et défis	10
En quoi consistent l'apprentissage fédéré et l'analyse décentralisée des données ?	10
Avantages possibles de l'apprentissage fédéré et de l'analyse décentralisée pour la recherche dans le domaine de la santé	14
Défis posés par l'apprentissage fédéré et l'analyse décentralisée, et les exigences nécessaires pour les affronter	15
Focus sur les plateformes d'apprentissage fédéré et d'analyse décentralisée nationales et internationales dans les régions de RLS-Sciences	19
Le projet Beacon de l'Alliance GA4GH	19
La plateforme CODA au Québec (Canada)	20
L'initiative GAIA-X en Europe centrée sur l'Allemagne	22
Le projet BORN en Bavière (Allemagne)	25
Les conditions de réussite de l'apprentissage fédéré et de l'analyse décentralisée pour la recherche fondée sur les données à l'échelle internationale	26
Remerciements	28
Références	29

Introduction

Au cours de la dernière décennie, les données massives et l'intelligence artificielle (IA) ont généré une vague de transformations que nous avons encore de la difficulté à évaluer et à envisager de manière globale. Dans le domaine de la santé comme dans plusieurs autres secteurs, l'augmentation exponentielle des données et de l'information générée, conjuguée à une meilleure capacité de calcul et de stockage, a radicalement changé la façon dont nous effectuons de la recherche, du développement technologique et dont nous prenons des décisions stratégiques (Mehta et collab., 2019). Cette transformation rapide a toutefois suscité des questions et des inquiétudes parmi les chercheurs, les décideurs politiques et la société en général, comme le montre le débat actuel sur la réglementation de l'IA qui vise à maximiser ses avantages et à atténuer ses risques dans l'intérêt du public.

La santé est souvent présentée comme un secteur à la fois prometteur pour l'application des données massives et en retard dans la mise en œuvre de l'IA (Shaw et collab., 2019). Ainsi, les algorithmes d'IA appliqués à l'aide au diagnostic, à la prédiction et à la précision des traitements, dans des disciplines telles que la radiologie et la cardiologie, ont connu des avancées considérables au cours des dernières années (Fornell, 2023). Pour autant, les experts ont souligné l'existence d'un fossé entre le développement de l'IA en médecine et ses applications concrètes au chevet des patients (McCradden et collab., 2019).

Ce fossé s'est certainement réduit grâce aux progrès réalisés pendant la pandémie de COVID-19. Il n'en reste pas moins que la collecte, l'utilisation et le partage de données de santé de qualité posent encore aujourd'hui de nombreux défis pour la mise au point d'applications prédictives utiles pour les patients, les professionnels et les décideurs.

Le lien entre la science et les politiques publiques est particulièrement important dans le domaine de la santé pour favoriser une utilisation responsable des mégadonnées et de l'IA. Dans ce cadre, le réseau scientifique multirégional et international, RLS-Sciences, réunit un groupe de chercheurs qui œuvrent pour la santé numérique. RLS-Sciences « apporte les bénéfices tirés de diverses perspectives culturelles et scientifiques grâce à la collaboration multilatérale entre les scientifiques, les décideurs politiques et les gestionnaires académiques [traduction] »¹. Le réseau fonctionne sous l'égide d'un forum politique, connu sous le nom de Sommet RLS-Sciences. Le forum se réunit tous les deux ans pour favoriser les échanges entre sept gouvernements régionaux : la Bavière (Allemagne), la Géorgie (États-Unis), le Québec (Canada), São Paulo (Brésil), la province de Shandong (Chine), la Haute-Autriche (Autriche) et le Cap-Occidental (Afrique du Sud) sur le thème « Des politiques publiques pour toutes les générations [traduction] »².

Les régions partenaires de RLS-Sciences ont convenu en 2012 de soutenir la collaboration scientifique entre leurs établissements de recherche régionaux, et ils ont d'abord choisi le domaine de l'énergie renouvelable comme champ d'action [traduction]³. Les régions ont alors invité les acteurs de la recherche de leur écosystème à se joindre au premier groupe

1. <https://oecd-opsi.org/innovations/rls-sciences/>

2. <https://www.mdpi.com/2071-1050/13/1/76>

3. <https://www.rls-sciences.org/about.html>

nommé *RLS-Energy Network*, dans lequel les chercheurs ont pu échanger des pratiques exemplaires, des données et des résultats de recherche dans le cadre d'un réseau scientifique de confiance. En 2016, lors de la 8e édition du Sommet RLS-Sciences à Munich, sur la base d'une analyse des forces scientifiques des régions, les membres bavarois de RLS ont encouragé la mise en place de trois autres groupes : *RLS-Small Satellites*, *RLS Global Aerospace Campus* et *RLS-Expert Dialogue on Digitalization*.

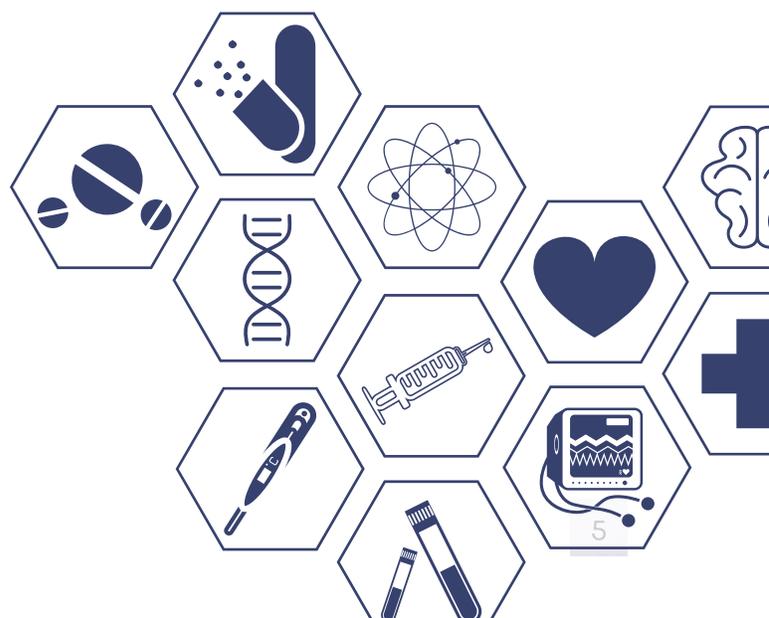
Pour soutenir le réseau RLS-Sciences dans son ensemble, ainsi que les quatre groupes, les gestionnaires académiques régionaux ont instauré des infrastructures de gouvernance et de soutien à la recherche, en nommant notamment des coordonnateurs experts issus des secteurs gouvernementaux, scientifiques et de la gestion des sciences, ainsi que des chefs de file scientifiques de chaque région pour chaque groupe thématique. Ensemble, ces représentants se sont employés activement à créer les conditions d'une collaboration internationale efficace au niveau régional dans le domaine scientifique. Les membres de RLS-Sciences ont participé aux sommets politiques biannuels en 2016, 2018 et 2021, renforçant ainsi le lien entre la science et le champ des politiques publiques. En 2020, sur l'invitation du gouvernement de Haute-Autriche, les représentants des régions partenaires et leurs dirigeants, ainsi que le réseau RLS-Sciences, se sont réunis autour d'une « table ronde virtuelle sur la COVID-19. »

Lors de cette rencontre, les régions partenaires ont échangé leurs expériences, leurs pratiques exemplaires, leurs défis et leurs stratégies pour lutter contre la pandémie de COVID-19. Avant la table ronde, le gouvernement bavarois a proposé la création d'un cinquième groupe

axé sur le domaine de la santé numérique, le *RLS-Digital Health*. Le groupe a alors entrepris une phase préparatoire à partir de 2020 et a été officiellement lancé en tant que cinquième groupe de RLS-Sciences en 2022.

Étant le plus récent groupe de RLS-Sciences, RLS-Digital Health travaille selon une feuille de route élaborée sur la base d'une matrice de méthodes et de cas d'utilisation cliniques pertinents pour les régions partenaires.

Le groupe s'engage à faire progresser les connaissances et à promouvoir la recherche translationnelle portant sur un large éventail de sujets dans le domaine de la santé numérique, et ce de façon à combler le fossé entre la recherche et la mise en pratique des applications. Le groupe a d'ailleurs identifié les méthodes d'apprentissage fédéré et d'analyse décentralisée des données en santé comme étant particulièrement pertinentes pour les futurs efforts de collaboration, notamment parce qu'elles constituent des mécanismes permettant de tirer le meilleur parti de la diversité des données au sein de chacune des régions concernées.



Au moment de la conception et de la rédaction de ce livre blanc, les chercheurs suivants étaient désignés comme scientifiques principaux du RLS-Digital Health:

- > PD Dr Sebastian Bickelhaupt, centre hospitalier universitaire d'Erlangen, Bavière
- > Prof. Jessica Kissinger, Université de Géorgie, Géorgie
- > Prof. Yves Joannette, Université de Montréal, Québec
- > Prof. Agma Traina, Université de São Paulo, São Paulo
- > Prof. Yu Changbin, Shandong First Medical, Université, Shandong
- > Prof. Michael Giretzlehner, RISC Software GmbH, Haute-Autriche

Le présent livre blanc porte sur les progrès en matière de données massives et d'IA en santé, et s'intéresse à une tendance centrale dans le domaine de la santé numérique, à savoir l'apprentissage fédéré et l'analyse décentralisée des données. Ces approches proposent d'exploiter tout le potentiel des données de santé en permettant une utilisation sécurisée de sources multiples sans avoir à les regrouper sur un seul site (AbdulRhaman, 2020). L'apprentissage fédéré et l'analyse décentralisée des données peuvent ainsi être présentés comme une réponse aux défis juridiques, éthiques et techniques actuels qui restreignent le partage des données entre les institutions et les juridictions et réduisent ainsi notre capacité à mener des recherches collaboratives à l'échelle nationale et internationale (Kairouz et collab., 2021). Si l'apprentissage fédéré et l'analyse décentralisée des données offrent de réelles opportunités d'amélioration de la recherche et de l'innovation fondées sur les données massives et l'IA, ces approches soulèvent toutefois plusieurs enjeux

tels que la protection de la vie privée, la fiabilité des données et l'utilisation judicieuse des ressources. Ces enjeux constituent quelques-unes des questions explorées dans le présent document.

Ce livre blanc explore ainsi tout le potentiel et les défis de l'apprentissage fédéré et de l'analyse décentralisée des données pour la recherche collaborative dans le domaine de la santé numérique. Il décrit également des plateformes et des technologies robustes qui démontrent comment l'apprentissage fédéré et l'analyse décentralisée peuvent être mis en application dans les systèmes de santé d'aujourd'hui.

La collaboration entre les membres de *RLS-Digital Health* a ainsi permis de mettre en lumière des projets qui constituent des manifestations éloquentes des promesses de l'apprentissage fédéré et de l'analyse décentralisée pour la recherche et l'innovation fondées sur les données de santé. En s'appuyant sur ces initiatives inspirantes, ce livre blanc présente les éléments clés qui pourraient favoriser la mise en place d'infrastructures performantes permettant de connecter et d'analyser des sources de données de haute qualité et en temps réel, tout en garantissant la mise en place des meilleures normes de protection et de standardisation des données. Ces éléments pourraient nous aider à définir et à construire ensemble un modèle pour une recherche collaborative axée sur les données à l'échelle nationale et internationale, qui pourrait profiter aux chercheurs, aux innovateurs, aux décideurs et aux patients.

La recherche fondée sur les données dans le domaine de la santé : perspectives et défis

Les progrès de la recherche fondée sur les données au sein des institutions et des juridictions

Malgré les crises et les tragédies qu'elle a engendrées, la pandémie de COVID-19 a tracé la voie à de véritables progrès en matière de recherche et d'innovation collaboratives fondées sur des données provenant des quatre coins du monde (Bragazzi et collab., 2020). Lorsque le SRAS-CoV-2 a été détecté pour la première fois en Chine en janvier 2020, les chercheurs en savaient très peu sur ce nouveau coronavirus et sur la manière de répondre à la menace qu'il représentait. Cependant, quelques jours plus tard, la séquence complète du génome du virus a été déterminée et communiquée à l'ensemble de la communauté scientifique⁴. À titre de comparaison, lors de l'épidémie de SRAS en 2003, ce même effort avait pris près de trois mois, et avant cela, on croyait faussement que la maladie était causée par la bactérie de la chlamydia.

La pandémie de COVID-19 a donc donné lieu à plusieurs initiatives visant à accélérer le partage des données sur les plans national et international pour la recherche, le développement

de traitements et l'évaluation des politiques publiques.

Les informations partagées durant la crise comprenaient des données moléculaires (des séquences aux cibles médicamenteuses), des données épidémiologiques, des données sur les interventions, ainsi que des données concernant des politiques et stratégies publiques qui ont été essentielles pour faciliter la collaboration à l'échelle internationale et la prise de décisions fondées sur des données probantes en vue de lutter contre le virus (Dagliatti, 2021). De nouveaux dépôts et répertoires de données ont été mis en place pour accélérer la mise en commun et l'utilisation des données sur la COVID-19 pour la recherche, tels le Portail officiel des données européennes et le CoronaNet. De surcroît, des infrastructures, de stockage, d'accès et de traitement des données déjà en place ont été utilisées à leur plein potentiel pour favoriser la collaboration en recherche à l'échelle nationale et internationale. C'est ce qu'illustrent les initiatives relatives à la COVID-19 lancées par des infrastructures de données reconnues comme la biobanque du Royaume-Uni (UK Biobank) et le Consortium international sur les infections respiratoires aiguës sévères et émergentes (ISARIC).

Ces initiatives ont ainsi souligné l'importance d'obtenir l'alignement des cadres de gouvernance des données entre les équipes de recherche et les institutions, ainsi que le rôle clé des standards et des terminologies communes pour l'harmonisation des données (Ros et collab., 2021).

Pour ce qui est des formats de données en particulier, les modèles de données communs, tels le partenariat Observational

4. La Commission européenne (2020). « Covid-19: How unprecedented data sharing has led to faster-than-ever outbreak research ». Dans la revue Horizon : The EU Research & Innovation Magazine. [En ligne]. <https://ec.europa.eu/research-and-innovation/en/horizon-magazine/covid-19-how-unprecedented-data-sharing-has-led-faster-ever-outbreak-research>

Medical Outcomes Partnerships (OMOP) et les standards de communication comme Fast Healthcare Interoperability Resources (FHIR) ont joué un rôle clé dans le partage et l'utilisation des données liées à la COVID-19 à des fins de recherche et d'innovation.

En fin de compte, ces efforts de collaboration ont abouti à une augmentation rapide de la production et de la diffusion des connaissances, comme le montrent les milliers de publications produites pendant la pandémie. Des vaccins ont été mis au point en un temps record et de nouveaux traitements ont été proposés aux patients en l'espace d'un an. Toutefois, cette avancée spectaculaire de la recherche et de la collaboration fondées sur les données a également mis en lumière d'importantes limites dans notre capacité à utiliser et à partager les données pour améliorer la santé de la population et relever les défis des systèmes de santé.

Défis de la recherche fondée sur les données à l'échelle nationale et internationale

La pandémie de COVID-19 a démontré comment, dans un court laps de temps et dans un contexte d'urgence, une communauté mondiale composée de divers acteurs peut se mettre d'accord sur des normes et une terminologie communes pour nommer une maladie, des variants, des gènes, des tests, etc. Ces efforts déployés à l'échelle internationale ont apporté des avantages considérables pour la recherche et l'élaboration de solutions efficaces pour endiguer et combattre la pandémie.

Toutefois, la recherche fondée sur les données se heurte encore à de nombreux obstacles,

à l'intérieur comme à l'extérieur des frontières. Ces obstacles limitent notre capacité à mener des projets de collaboration mobilisant de vastes quantités de données de santé qui soient diverses et de grande qualité (Abdulrahman et collab. 2021; Nguyen et collab. 2022; Ros et collab., 2021). En voici quelques exemples :

- > **Enjeux relatifs à la vie privée :** Les données de santé sont des renseignements personnels sensibles qui nécessitent un degré élevé de protection et d'attention afin de garantir le respect de la vie privée et la confidentialité.
- > **Réglementations restrictives :** Comme les données de santé sont considérées comme des renseignements personnels sensibles, les réglementations existantes (*Règlement général sur la protection des données* (RGPD), *Loi sur la portabilité et la responsabilité des assurances-maladie* (HIPAA), etc.) tendent à restreindre le transfert de données d'un pays à l'autre et même, le plus souvent, entre établissements de santé.
- > **Préoccupations éthiques :** Dans la plupart des pays, les citoyens partagent les mêmes préoccupations quant à l'utilisation de leurs données de santé et exigent des garanties en contrepartie de leur usage : consentement individuel éclairé, informations relatives à l'utilisation des données, possibilité de retirer son consentement, etc.
- > **Manque de standardisation et d'interopérabilité des données :** Les données de santé sont collectées et conservées sous divers formats qui ne respectent pas toujours les standards et la terminologie internationaux, ce qui compromet l'appariement et la réutilisation des données.
- > **Qualité et disponibilité médiocres des données :** Il est bien connu que la qualité des données de santé varie considérablement

et qu'il est difficile de les utiliser en raison du manque de standardisation entre les sources de données, en particulier d'un système de santé à un autre. Certaines données peuvent être trouvées dans un format structuré (données d'imagerie, données sociodémographiques, résultats de laboratoire sous format tabulaire), mais bon nombre d'entre elles ne sont pas structurées (notes en texte libre, documents en papier numérisés, etc.).

> **Manque d'alignement entre les pratiques d'évaluation des projets de recherche :**

En plus de tous ces défis, les comités d'éthique de la recherche au sein des pays et entre eux adoptent des normes et des pratiques différentes pour évaluer les projets de recherche, ce qui peut poser des obstacles à la collaboration internationale.

Ces défis ne sont pas nouveaux et constituent depuis des décennies de sérieuses barrières à la réutilisation des données de santé à des fins de recherche et d'innovation (Price et collab., 2019). Pourtant, avec les progrès de l'analyse des données et l'augmentation rapide de la production d'informations par les systèmes, les capteurs et les applications, ces défis sont désormais considérés comme des menaces qui pourraient mettre un terme au développement et à la mise en œuvre de l'IA en santé (Morley et collab., 2020).

En effet, les algorithmes d'IA, en particulier en particulier ceux qui reposent sur l'apprentissage automatique et l'apprentissage profond, ont besoin d'une grande quantité de données structurées et mises en forme pour l'apprentissage et la validation algorithmiques. Des données de haute qualité, diversifiées et représentatives permettent de favoriser la précision,

la robustesse et la transférabilité des algorithmes d'IA dans les milieux de soins (Peifer et collab., 2020).

Faute de données de qualité, les algorithmes peuvent d'ailleurs donner de très mauvais résultats dans la pratique; c'est précisément ce que signifie le dicton « garbage in, garbage out » (à données inexactes, résultats erronés, en français).

Comme indiqué précédemment, la disponibilité et l'accessibilité des ensembles de données de santé sont très limitées à l'heure actuelle, de sorte que les efforts visant à regrouper des données massives pour la recherche et l'innovation dans le domaine de l'IA peuvent se heurter à des limites importantes. Tout d'abord, bien que le nombre de dépôts de données de nature académique ait augmenté ces dernières années, ceux-ci contiennent principalement des données de recherche qui ne sont pas représentatives des populations de patients et des véritables parcours de soins (c'est-à-dire qu'il s'agit de données d'études cliniques, de données non-observationnelles, etc.). Autrement dit, les données issues des entrepôts de recherche sont utiles, mais elles contiennent des biais qui peuvent limiter les performances en milieu réel des algorithmes qui sont entraînés et calibrés à l'aide de ces informations.



Deuxièmement, les efforts visant à structurer les données de vie réelle (c'est-à-dire les données dossiers électroniques des patients, les données de l'Internet des objets, etc.) pour une utilisation secondaire ont été remarquables ces dernières années, notamment grâce au développement de lacs et d'entrepôts de données pour la recherche et l'analyse à l'échelle organisationnelle ou nationale (Rieke et collab., 2020). Néanmoins, ces infrastructures sont pour la plupart basées sur un modèle centralisé où les données réelles sont entreposées, conservées et utilisées sur un site unique.

Le partage des données de vie réelle au-delà du site où elles ont été collectées et générées demeure un défi de taille dans plusieurs juridictions. Il pose non seulement des enjeux juridiques et éthiques, liés à la protection de la vie privée et de la confidentialité des données, mais aussi des enjeux techniques.

L'anonymisation des données, la gestion d'un accès sûr et efficace aux données et le transfert des données sont des activités non triviales qui nécessitent des ressources importantes et une expertise pluridisciplinaire, en particulier dans un contexte où la réglementation en matière de protection des données évolue rapidement (Daglieti, 2021).

C'est pourquoi des modèles alternatifs ont été envisagés et expérimentés comme solution de rechange à la centralisation et au transfert des données. L'apprentissage fédéré et l'analyse décentralisée des données peuvent être un moyen efficace de contourner les obstacles actuels liés au partage des données et d'accélérer la recherche et l'innovation fondées sur des informations réparties entre les sites et entre les frontières. Le présent document a pour but d'explorer les avantages de ce nouveau paradigme pour l'utilisation des données et de présenter certains des nouveaux défis qu'il soulève pour la recherche fondée sur les données dans le domaine de la santé.

L'apprentissage fédéré et l'analyse décentralisée des données de santé : perspectives et défis

En quoi consistent l'apprentissage fédéré et l'analyse décentralisée des données ?

Compte tenu des préoccupations actuelles en matière de protection de la vie privée et des restrictions juridiques au partage des données, un paradigme a récemment émergé dans le domaine de l'analyse des données et de l'apprentissage automatique, que l'on appelle l'apprentissage fédéré. Selon Adbul Rahman et ses collègues (2021), l'apprentissage fédéré « est une approche décentralisée préservant la vie privée, qui conserve les données brutes sur les sites et suppose un apprentissage automatique local qui peut éliminer la surcharge liée au transfert des données [traduction]. » L'apprentissage fédéré présente donc une alternative aux systèmes centralisés pour l'analyse des données et l'entraînement d'algorithmes d'IA.

Dans le cadre de l'apprentissage fédéré, les données sont conservées sur les sites d'origine où elles ont été collectées ou générées (ces sites sont parfois appelés « nœuds »). Les sites s'entendent pour collaborer à l'analyse commune des données ou à l'entraînement d'un modèle algorithmique sous la direction d'un serveur central. Le serveur reçoit des résultats

d'analyse des sites (statistiques, paramètres, gradients, poids, etc.), mais jamais les données brutes elles-mêmes. Les résultats envoyés sont ensuite agrégés sur le serveur central pour obtenir un algorithme amélioré et plus performant, qui a été entraîné sur une base de données beaucoup plus importante. Le modèle ou les statistiques peuvent ensuite être partagés entre les différents sites afin de contribuer à la connaissance collective. À aucun moment, les sites n'ont accès aux données brutes des collaborateurs, uniquement aux résultats analytiques.

Grâce à ce processus décentralisé qui assure la protection de la vie privée, l'apprentissage fédéré offre

la possibilité de contourner les enjeux légaux, éthiques et techniques liés à la mise en commun et au partage de renseignements sensibles et potentiellement identifiables dans le domaine de la santé (Rieke et collab., 2020).

La recherche a d'ailleurs montré que les modèles d'apprentissage automatique entraînés à partir de l'apprentissage fédéré peuvent atteindre des niveaux de performance comparables à ceux entraînés sur des ensembles de données hébergés de manière centralisée. Les modèles fédérés sont aussi supérieurs en performance à ceux entraînés sur des données provenant d'un seul site (Li et collab., 2019; Roy et collab., 2019).

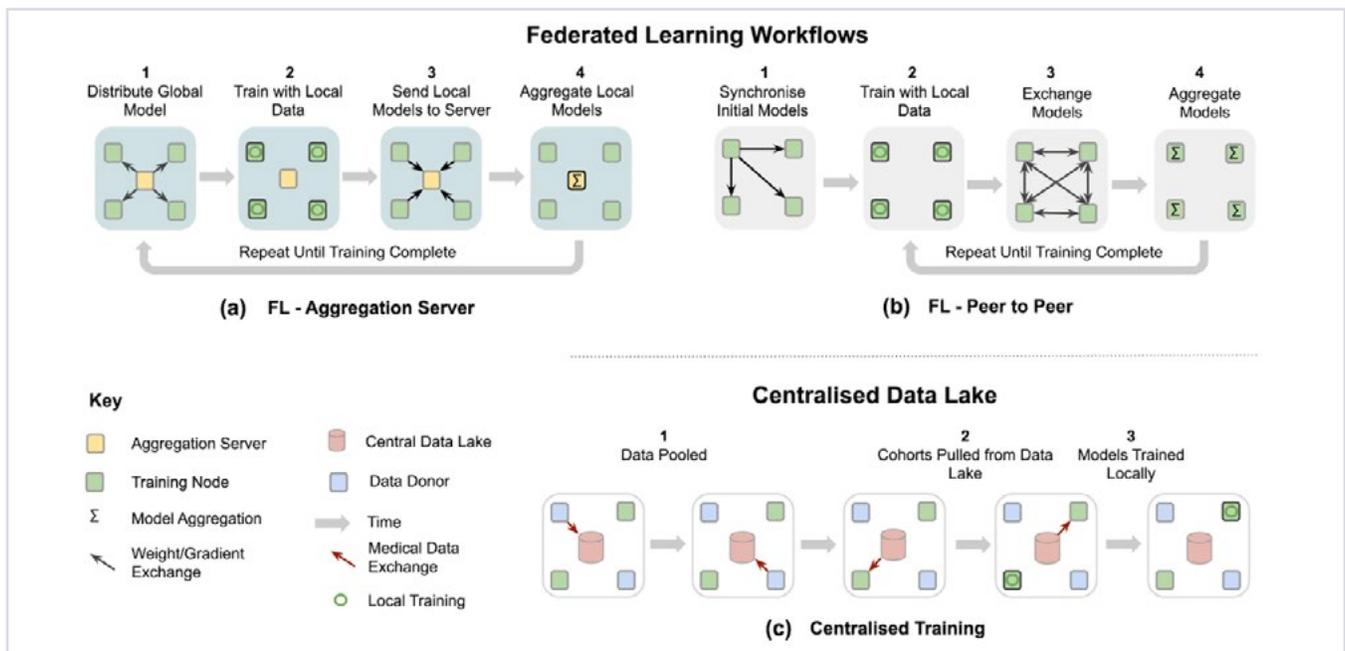


Figure 1 : Exemple de flux de travail de l'apprentissage fédéré et comparaison avec l'apprentissage à partir d'un lac de données centralisées (extrait de Rieke et collab., 2020)

L'apprentissage fédéré n'est toutefois qu'une approche parmi d'autres de l'analyse décentralisée des données, qui recouvre un plus large éventail de méthodologies et de processus.

Ainsi, l'analyse décentralisée consiste à appliquer les méthodes de base de la science des données (statistiques, régression linéaire, etc.) pour le traitement des données, tandis que l'apprentissage fédéré porte sur le perfectionnement des modèles d'apprentissage automatique à distance et sur l'obtention de résultats de prédiction agrégés.

En ce sens, l'apprentissage fédéré peut être considéré comme un sous-ensemble de l'analyse décentralisée, car il comprend un type spécifique d'analyse reposant sur l'apprentissage automatique dans un contexte qui est en fait beaucoup plus large.

Dans ce contexte plus global, l'un des inconvénients de l'apprentissage fédéré réside dans le fait qu'il dépend d'un serveur central, ce qui oblige tous les sites participants à se mettre d'accord pour désigner un organisme central de confiance, dont la défaillance ou l'inefficacité pourrait en fait compromettre le processus d'entraînement algorithmique et d'analyse des

données. Une alternative à l'apprentissage fédéré repose sur une approche décentralisée qui ne nécessite pas de coordination centrale. Cette approche a déjà été envisagée et mise en place au moyen d'une infrastructure telle que la *Decentralized Zero-Trust IoT Data Fabric* de l'Université de Géorgie (États-Unis).⁵

Ce projet de recherche consiste à établir une architecture de données décentralisée fondée sur les principes du Web 3.0 et de la chaîne de blocs. Le système permettrait à chaque propriétaire de données d'exercer une autorité totale sur ses données et d'accorder ou de révoquer l'accès aux données à n'importe quel utilisateur sans avoir besoin de passer par un intermédiaire.

En outre, toutes les modifications apportées aux données et la gestion des accès seraient identifiables et enregistrées, et feraient l'objet d'un audit, ce qui garantit la transparence et la responsabilité. Ce modèle est appelé « Zero-Trust » parce qu'il n'exige pas des propriétaires ou des fiduciaires de données qu'ils fassent confiance à un organisme central de coordination pour décider de l'accès et de l'utilisation des données selon des modalités qui sont acceptables et autorisées.

5. University of Georgia College of Engineering. « Decentralized Zero-Trust IoT Data Fabric ». [En ligne]. <https://sensorweb.engr.uga.edu/index.php/wns/>

Types de modèles d'apprentissage fédéré

Il existe plusieurs types de modèles d'apprentissage fédéré qui peuvent être mis en œuvre dans le domaine de la santé pour faciliter l'analyse décentralisée des données et l'entraînement de modèles d'IA (Joshi et collab., 2022; Mammen, 2021) :

Apprentissage fédéré vertical — Ce type d'apprentissage peut être utilisé, par exemple, lorsque différentes organisations disposent de données sur le même groupe

de patients, mais avec des caractéristiques différentes. L'apprentissage fédéré vertical permet de construire un modèle d'IA à partir d'un ensemble de données plus complet.

Apprentissage fédéré horizontal –

Ce type d'apprentissage peut être utilisé lorsque différentes organisations disposent de données présentant les mêmes caractéristiques, mais concernant des groupes de patients différents. L'apprentissage fédéré horizontal peut servir à entraîner un modèle sur un ensemble de données plus vaste contenant un plus grand nombre de patients et une plus grande variabilité.

Apprentissage fédéré par transfert –

Ce type d'apprentissage consiste à ajouter une nouvelle caractéristique à un modèle préalablement entraîné, à l'instar de l'apprentissage automatique traditionnel. Il s'agit par exemple d'étendre l'apprentissage fédéré vertical de manière à inclure des échantillons supplémentaires qui peuvent ne pas exister au sein de toutes les organisations participant à l'apprentissage.

Apprentissage fédéré intersilos –

Ce type d'apprentissage est utilisé lorsque les nœuds/sites participants sont moins nombreux et non disponibles pour tous les cycles de l'apprentissage d'un modèle. Ce type d'apprentissage est normalement appliqué aux organisations qui disposent d'une grande quantité de données, mais qui ne peuvent pas les partager. Ce modèle peut s'ajouter à l'apprentissage fédéré vertical ou horizontal.

Apprentissage fédéré interappareils –

Ce modèle est approprié lorsque le nombre de nœuds participants est élevé et que chaque nœud dispose de petites quantités de données. Par conséquent, l'approche interappareils permet de développer des modèles algorithmiques lorsque les données sont réparties à grande échelle tout en faisant usage d'une même application. Cela peut être le cas lors de l'entraînement d'un modèle sur des appareils mobiles et des applications de l'Internet des Objets.

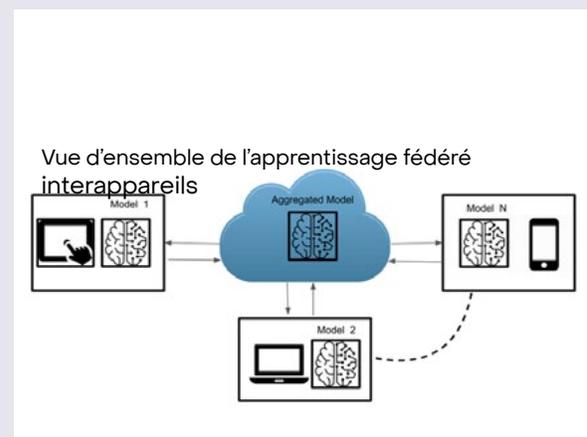
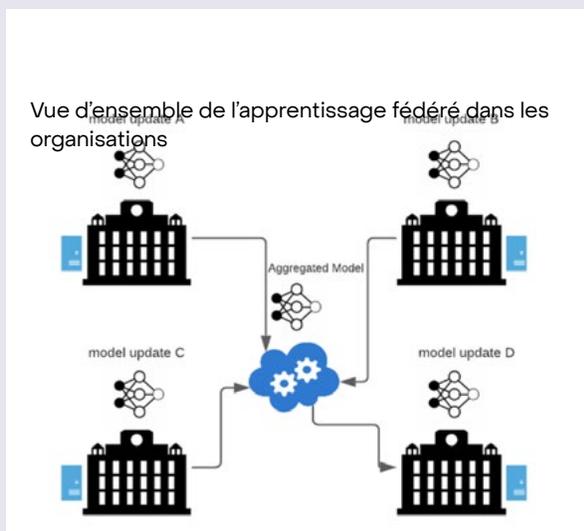


Figure 2 : L'apprentissage fédéré interappareils et interorganisations (extrait de Mammen., 2021)

Avantages possibles de l'apprentissage fédéré et de l'analyse décentralisée pour la recherche dans le domaine de la santé

En comparaison avec les approches centralisées, l'apprentissage fédéré et l'analyse décentralisée sont susceptibles d'offrir plusieurs avantages pour la recherche et à la collaboration fondés sur les données de santé. Nous présentons ci-dessous les avantages les plus fréquemment relevés par les experts consultés et par la littérature (Mammen, 2022; Rieke et collab., 2020; Sheller et collab., 2020; Xu et collab., 2020) :

> **Confidentialité des données :** L'apprentissage fédéré et l'analyse décentralisée pourraient renforcer la confidentialité des données, notamment par l'entraînement de modèles d'apprentissage automatique sur des appareils et des nœuds locaux, sans avoir à partager les données brutes (uniquement les données agrégées, les poids et les gradients) avec un serveur central. Cela permettrait de protéger les données confidentielles des patients et d'assurer la conformité avec les réglementations en vigueur, comme HIPAA, la Loi sur la protection des renseignements personnels et les documents électroniques (LPRPDE) et le RGPD, réduisant ainsi le risque de fuite de données.

- > **Contrôle des données locales :** L'apprentissage fédéré et l'analyse décentralisée pourraient permettre aux établissements de santé de conserver le contrôle de leurs données locales tout en continuant à produire des connaissances et à effectuer des analyses statistiques destinées à être partagées de manière collaborative. Cela permettrait d'assurer l'autonomie des organisations et le contrôle des données, ce qui peut s'avérer nécessaire pour répondre aux exigences des autorités réglementaires ou des organisations.
- > **Diversité et mise à l'échelle des données :** L'apprentissage fédéré et l'analyse décentralisée permettent de regrouper diverses données provenant de plusieurs établissements de santé, ce qui faciliterait l'entraînement de modèles d'IA plus robustes et transférables. Cela pourrait améliorer la performance des modèles et fournir de meilleures analyses prédictives, en exploitant le potentiel d'ensembles de données plus complets et diversifiés tout en éliminant le besoin de transfert ou de centralisation des données.
- > **Mises à jour en temps réel :** L'apprentissage fédéré et l'analyse décentralisée pourraient permettre d'apporter des mises à jour en continu, au fur et à mesure que les modèles algorithmiques sont raffinés localement sur les serveurs ou les appareils présents sur site, permettant ainsi d'adapter et d'améliorer le modèle global en temps réel. Cette fonctionnalité est particulièrement utile dans les établissements de santé où la distribution et les caractéristiques des données peuvent changer au fil du temps.
- > **Rapport coût-efficacité :** L'apprentissage fédéré et l'analyse décentralisée ont le potentiel d'être rentables, car ils peuvent



éliminer le besoin de transfert ou de stockage central des données, réduisant ainsi les frais de communication et les coûts qui y sont associés. Cela permettrait également aux établissements de santé de tirer parti de leur infrastructure et de leurs ressources existantes pour effectuer des analyses de données localement.

> **Collaboration et partage des connaissances :**

L'apprentissage fédéré et l'analyse décentralisée peuvent promouvoir la collaboration entre les établissements de santé en leur permettant de travailler ensemble à l'élaboration d'un modèle ou d'un projet global, tout en préservant la sécurité et le contrôle des données sur le plan local. Cela peut favoriser l'échange de connaissances et de savoir-faire ainsi que des efforts de recherche conjoints, ce qui a le potentiel d'améliorer les pratiques de santé et les résultats pour les patients.

Défis posés par l'apprentissage fédéré et l'analyse décentralisée, et les exigences nécessaires pour les affronter

En dépit des avantages potentiels de l'apprentissage fédéré et de l'analyse décentralisée des données pour la recherche et l'innovation en santé, plusieurs préoccupations et défis soulevés par les experts et la littérature limitent encore leur mise en œuvre au sein et entre les établissements de santé (Kairouz et collab., 2021; Liu et collab., 2022). C'est pourquoi il est nécessaire de mettre en œuvre des techniques et des processus spécifiques pour

relever ces défis et faciliter le déploiement de l'apprentissage fédéré et de l'analyse décentralisée dans les établissements de santé. Parmi les défis à relever et les exigences nécessaires pour les affronter, il convient de souligner les points suivants :

Hétérogénéité des données : Il est reconnu que les données de santé présentent une variabilité à plusieurs égards, du fait de leurs formats, leur qualité, leur distribution et leur représentation différenciés. Par conséquent, les ensembles de données entre les nœuds et les établissements participants à une infrastructure fédérée peuvent présenter des caractéristiques différentes qui sont susceptibles de limiter notre capacité à les analyser de manière décentralisée (Kairouz et collab., 2021; Liu et collab., 2022).

Exigences : La standardisation des données et l'interopérabilité entre les systèmes d'information sont essentielles pour garantir le succès d'une infrastructure d'apprentissage fédéré ou d'analyse centralisée. L'harmonisation entre les nœuds et les sites peut être facilitée par l'application de normes, comme le standard Fast Healthcare Interoperability Resources (FHIR) pour les données des dossiers électroniques des patients et le standard Digital Imaging and Communications in Medicine (DICOM) pour les données d'imagerie. En outre, des corrections doivent être appliquées pour garantir que des ensembles de données non-identiquement distribués puissent être utilisés conjointement pour l'entraînement d'un modèle algorithmique.

Biais dans les données : Cet enjeu fait référence à la présence d'erreurs systématiques ou de postulats biaisés dans les données présentes dans les nœuds contribuant au

processus d'apprentissage fédéré ou d'analyse décentralisée. Cela peut se produire, par exemple, en raison des choix effectués par chaque établissement participant lors de la collecte et du traitement des données (Kairouz et collab., 2021). Différents appareils peuvent également être utilisés pour collecter, analyser et stocker les données, et ces dernières peuvent être recueillies auprès de groupes populationnels différents d'un établissement à l'autre, ce qui peut conduire à des résultats biaisés dans le modèle global.

Exigences : Au sein des sites et entre ceux-ci, les vérifications et les corrections des biais doivent être constantes pendant le processus d'apprentissage fédéré ou d'analyse décentralisée. En outre, certaines techniques peuvent être utilisées pour contrôler ce risque et garantir une plus grande équité dans le processus; il en est ainsi de l'apprentissage fédéré agnostique (Mohri et collab., 2019). Dans ce scénario, le serveur central peut appliquer des mécanismes de pondération ou d'équité lors de l'agrégation des résultats d'analyses provenant de différents nœuds. Cela permet une représentation et une prise en compte égales des nœuds dont les caractéristiques des données varient.

Sécurité et confidentialité des données : Des études récentes ont montré que l'apprentissage fédéré et l'analyse décentralisée ne garantissent pas suffisamment le respect de la vie privée, car des renseignements confidentiels peuvent être révélés au cours du processus d'analyse et d'entraînement des modèles (Mothukuri et collab., 2021; Yin et collab., 2021). En effet, au cours du processus, les nœuds envoient des informations à un serveur central, ce qui rend l'infrastructure vulnérable à plusieurs types d'attaques comme les attaques par inférence de

e propriétés, les attaques par reconstruction et les attaques par inférence d'appartenance (Hu et collab., 2021). Ces attaques visent à identifier si un individu était présent ou non dans les ensembles de données d'apprentissage, ce qui présente des risques pour la vie privée.

Exigences : Plusieurs techniques de protection de la vie privée doivent être mises en œuvre afin de limiter le risque de réidentification des individus au cours du processus d'entraînement des modèles (voir l'encadré ci-dessous). Cependant, le risque ne peut jamais être réduit à zéro, car il est nécessaire de trouver un équilibre entre l'exactitude des informations transmises au serveur central et le degré de confidentialité que les participants souhaitent maintenir au sein de l'infrastructure d'apprentissage fédéré ou d'analyse décentralisée.

Empoisonnement des données : Il s'agit d'un cas où un établissement malveillant cherche à empoisonner le modèle global en envoyant des mises à jour du modèle dérivées de données étiquetées de façon erronée. De telles attaques par empoisonnement des données peuvent provoquer une baisse substantielle de l'exactitude des informations transmises et de la précision du modèle global, même si le pourcentage d'établissements aux intentions malveillantes est faible (Tolpegin et collab. 2020).

Exigences : Le risque d'empoisonnement des données peut être réduit s'il existe des moyens de contrôler les établissements qui peuvent contribuer au processus d'apprentissage fédéré ou d'analyse décentralisée. Dans le cadre d'un réseau universitaire, par exemple, le risque d'attaques malveillantes sur des modèles provenant des établissements est

considérablement réduit, car les conditions de participation à l'entraînement des modèles peuvent être très rigoureuses et respecter des lois et des lignes directrices éthiques contraignantes (se référer à la partie 3 pour obtenir des exemples de tels réseaux universitaires).

Défis en matière de communication et de calcul : L'apprentissage fédéré et l'analyse fédérée supposent l'analyse et l'échange d'informations entre des modèles ou des serveurs locaux, ce qui peut se traduire par une augmentation des coûts de communication et de calcul informatique. Cela peut poser des défis en matière de bande passante, de latence et de ressources informatiques, en particulier lorsque l'infrastructure fédérée connecte des établissements à grande échelle ou quand les sources de données sont éparpillées.

Exigences : L'apprentissage fédéré et l'analyse fédérée nécessitent des ressources importantes, notamment en matière de puissance de calcul, de stockage et d'expertise, sur les sites locaux. Il peut être difficile de s'assurer que tous les participants disposent des ressources nécessaires pour participer activement au processus d'analyse et d'entraînement, en particulier dans les petits établissements de santé ou dans les contextes où les ressources sont limitées. Cette difficulté peut s'expliquer par le fait que l'expertise humaine en matière d'ingénierie et de science des données tend à se faire rare dans les établissements de santé financés par des fonds publics.

Gouvernance et considérations juridiques et éthiques : L'apprentissage fédéré et l'analyse décentralisée des données reposent en grande partie sur la confiance, la collaboration et l'alignement entre de multiples parties prenantes, y compris les établissements de santé, les propriétaires de données, les comités d'éthique, entre autres acteurs. L'élaboration de modèles de gouvernance efficaces et la prise en compte des questions juridiques et éthiques, comme le consentement du patient, la propriété intellectuelle et la responsabilité, peuvent être complexes et exigeantes, et elles nécessitent un examen minutieux et la mise en œuvre de mesures appropriées.

Exigences : Des mécanismes de gouvernance et des cadres juridiques robustes pour assurer une bonne coordination, ainsi que des accords de partage des données et le strict respect des réglementations et des lignes directrices éthiques doivent être mis en œuvre et acceptés par toutes les entités participant à l'infrastructure fédérée. Ces mécanismes et ces cadres peuvent prendre la forme d'une documentation partagée fondée sur un formulaire de consentement unifié, un accord de partage des données et une voie d'accès aux données balisée. De la sorte, tous ces éléments peuvent être élaborés en partenariat et mis en œuvre à travers la création d'un comité représentant tous les acteurs participant à l'infrastructure fédérée.



Techniques de protection de la vie privée dans le cadre de l'apprentissage fédéré – définitions simplifiées

Confidentialité différentielle : Cette technique est utilisée pour protéger la vie privée des personnes en ajoutant du bruit ou des éléments aléatoires aux données de manière à préserver la précision statistique tout en réduisant la possibilité que les personnes soient à nouveau identifiées. Cette technique fournit ainsi un cadre mathématique pour quantifier et contrôler les risques pour la vie privée associés au partage ou à l'analyse de données de nature sensible, comme les renseignements personnels ou confidentiels, tout en permettant l'analyse des données.

Chiffrement homomorphique : Le chiffrement homomorphique est une technique cryptographique qui permet de chiffrer des données de manière à ce qu'elles puissent être traitées par les utilisateurs sans avoir à être décryptées. Autrement dit, cette technique permet d'effectuer des calculs directement

sur les données chiffrées, sans qu'il soit nécessaire de les déchiffrer au préalable. Cela permet d'atteindre un niveau élevé de confidentialité et de sécurité des données, car les données originales demeurent chiffrées tout au long du processus de calcul, y compris lorsqu'elles sont sauvegardées et transférées.

Calcul sécurisé multi-parties (CSMP) : Également connue sous le nom de calcul sécurisé ou d'évaluation sécurisée des fonctions, cette technique de chiffrement permet à plusieurs personnes de calculer conjointement une fonction sur leurs données, sans les révéler les unes aux autres. Elle permet aux participants de collaborer à des analyses tout en gardant leurs données individuelles confidentielles. Le CSMP garantit que les entrées et les résultats de l'analyse restent privés et sécurisés, même lorsque le calcul est effectué sur des systèmes ou des réseaux non fiables ou potentiellement malveillants.

Focus sur les plateformes d'apprentissage fédéré et d'analyse décentralisée nationales et internationales dans les régions de RLS-Sciences

Le projet Beacon de l'Alliance GA4GH

Sommaire

L'un des principaux défis auxquels est confrontée la recherche en génomique humaine est la pénurie de données. Pour surmonter cet obstacle, l'Alliance mondiale pour la génomique et la santé (GA4GH) a lancé le projet Beacon en 2014. Cette initiative vise à faciliter le partage des données génomiques et cliniques entre des réseaux fédérés (Fiume et collab., 2019). Les données génomiques étant particulièrement sensibles, le projet vise à fournir des lignes directrices en matière de réglementation,

d'éthique et de sécurité afin de garantir la mise en place de mesures appropriées pour l'analyse et le partage des données, conformément au « Cadre d'échange responsable de données liées à la génomique et à la santé » de l'Alliance GA4GH (Knoppers, 2014). Le projet Beacon permet ainsi l'intégration de données génomiques provenant de diverses sources dans le monde entier grâce à un protocole de requête commun.

En « beaconisant » leur ensemble de données omiques, les hôpitaux ou les instituts de recherche peuvent contribuer aux efforts scientifiques conjoints visant à accélérer la recherche en génomique et la médecine de précision, sans compromettre la confidentialité ou le contrôle local des données.

L'API qui sous-tend le projet Beacon a donc été conçue pour les chercheurs et les spécialistes afin de permettre l'interrogation des variants génomiques et de leurs informations connexes. Grâce à une infrastructure de données solide et à des pratiques responsables, le partage des données génomiques dans le cadre du projet Beacon permet d'obtenir de précieux renseignements sur les maladies, les pronostics et les variations génomiques liées aux habitudes de vie.

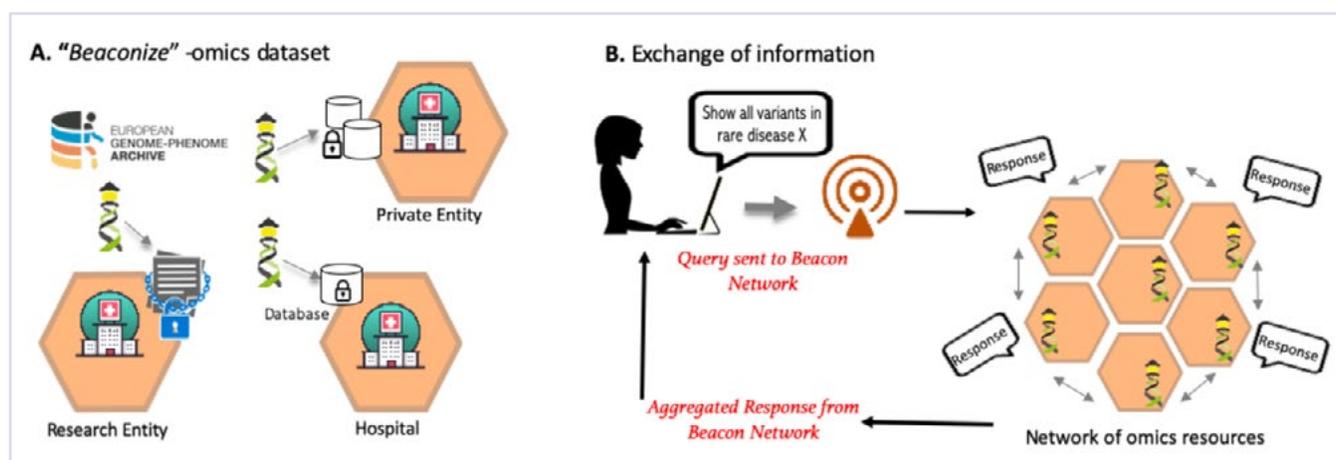


Figure 3 : Illustration de l'infrastructure et du fonctionnement de la balise (extrait de : <https://beacon-project.io/>)

Précisions supplémentaires

Les versions initiales du projet Beacon (versions 0 et 1) se limitent à indiquer la présence ou l'absence d'une mutation génomique spécifique dans un ensemble de données provenant de patients atteints d'une maladie particulière ou de la population générale. L'API de Beacon permet ainsi aux chercheurs de demander des renseignements sur un allèle spécifique. Par exemple, une requête peut consister à demander si un nucléotide, comme le C, a été observé à un emplacement génomique spécifique, comme la position 32936732 sur le chromosome 13. La réponse de Beacon est soit « oui, » soit « non ». L'API du projet Beacon permet donc d'accéder à distance à des renseignements alléliques pertinents sans qu'il soit nécessaire d'identifier un patient ou un échantillon particulier, ce qui réduit les risques en matière de protection de la vie privée.

La dernière version du projet Beacon (version 2) offre aux chercheurs une plus grande souplesse dans la recherche de variants génomiques et permet d'inclure des questions supplémentaires sur l'ensemble des données et les attributs des participants.

Dans un cadre sécurisé, les utilisateurs autorisés peuvent relier les résultats du projet Beacon à des données protégées, notamment le dossier électronique des patients, et de les associer à l'annotation des variants faite par des experts. Par ailleurs, les chercheurs peuvent demander l'accès à un ensemble de données figurant dans les résultats de leur recherche, et la version 2 du projet Beacon peut fournir des informations concernant les restrictions d'utilisation des données pour faciliter le processus d'accès. L'objectif consiste à permettre aux chercheurs d'étudier et de partager des données sur les

variants génomiques ainsi que des métadonnées essentielles comme des données cliniques et phénotypiques. Cette fonctionnalité pourrait permettre aux chercheurs d'étudier davantage de questions liées aux maladies rares et complexes.⁶

La plateforme CODA au Québec (Canada)

Sommaire

La pandémie de COVID-19 a mis en évidence les limites de la mise en commun des données en un lieu unique pour des fins d'analyse et de recherche. Ceci a été particulièrement mis en évidence lorsqu'il a été nécessaire de fournir des données en temps opportun sur une crise de santé publique qui évoluait très rapidement. Pour répondre à ce défi, l'hôpital de l'Université de Montréal au Québec (Canada) a créé la plateforme Collaborative Data Analysis (CODA) (ou « plateforme d'analyse de données collaborative » en français). En janvier 2020, cette plateforme a été conçue avec l'aide de parties prenantes issues de différents domaines tels que la recherche, la clinique, la gestion, les pouvoirs publics, l'éthique et le droit. Les principales exigences comprenaient la capacité d'effectuer des analyses de données décentralisées et de l'apprentissage automatique fédéré, de prendre en charge les normes et la terminologie médicales communes, de mettre en œuvre des mesures visant à minimiser la divulgation des données individuelles des patients et de garantir un déploiement de logiciels libres et non commerciaux.

La viabilité de la plateforme CODA a fait l'objet d'études dans huit hôpitaux au Canada, où des patients atteints de la COVID-19 (cas soupçonnés ou confirmés) ont été inclus dans l'analyse durant une période de trois ans. Les

6. Global Alliance for Genomics & Health. « New release of GA4GH Beacon expands genomic and clinical data access ». [En ligne]. <https://www.ga4gh.org/news/new-release-of-ga4gh-beacon-expands-genomic-and-clinical-data-access/>

capacités en matière d'apprentissage fédéré ont été examinées sur des ensembles de données cliniques et d'imagerie de référence provenant de patients gravement malades.

La plateforme CODA a été déployée avec succès et les résultats de l'étude de viabilité du déploiement seront publiés prochainement dans le Lancet Journal. Le code de logiciel, la documentation et les documents techniques de la plateforme CODA ont été publiés sous une licence en source ouverte. Désormais, la plateforme sera utilisée pour élaborer et valider prospectivement des modèles d'évaluation des risques, de surveillance proactive et de prévision de l'utilisation des ressources chez les patients hospitalisés et ambulatoires dans le cadre des efforts de validation en cours.

Précisions supplémentaires

La plateforme CODA consiste en un ensemble de microservices qui

fonctionnent conjointement pour permettre l'analyse décentralisée des données de santé (comme le montre la figure 4).

Elle comprend divers services qui effectuent la collecte et le traitement des données sur les différents sites hospitaliers (les nœuds), ainsi qu'un mécanisme de consolidation des analyses locales pour l'exécution des tâches distribuées (centre de coordination) et des composantes préliminaires (applications de tableau de bord et de prises de notes) qui facilitent les requêtes analytiques personnalisées, la visualisation des données et l'entraînement des modèles d'apprentissage automatique. Avant d'être ingérées dans la plateforme CODA, les données sont dépersonnalisées. Tous les canaux de communication entre les composantes de la plateforme sont sécurisés à l'aide des protocoles SSL/TLS (*Secure Sockets Layer/Transport Layers Security*).

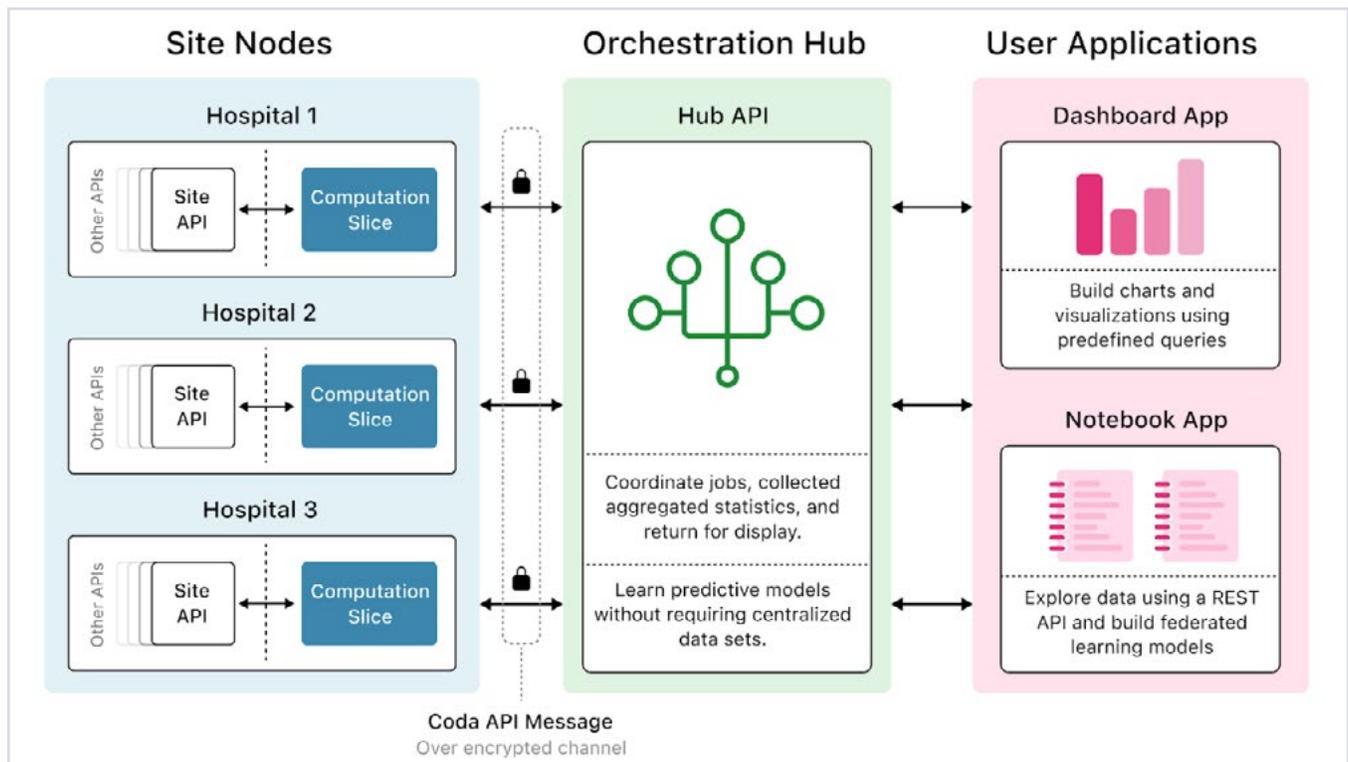


Figure 4 : Vue d'ensemble de l'infrastructure de la plateforme CODA (source : CITADEL)

Les nœuds des sites, situés à l'intérieur des pare-feux institutionnels des établissements de santé participants, comprennent des composantes permettant de sauvegarder et d'extraire des données dépersonnalisées de dossiers électroniques des patients (DME), des données d'imagerie et des données sur les formes d'onde. Les données des DME sont sauvegardées conformément à la norme Fast Healthcare Interoperability Resources (FHIR), alors que les données d'imagerie et de formes d'onde sont sauvegardées conformément à la norme Digital Imaging and Communications in Medicine (DICOM).

L'étude de viabilité de CODA visait à déployer la plateforme dans neuf hôpitaux publics du Québec (Canada), comme l'illustre la figure 5.

Un cadre de gouvernance a été établi pour formaliser les modalités juridiques et éthiques de la collaboration entre les établissements participants.⁷

Huit des neuf sites inscrits ont déployé avec succès la plateforme localement et sont connectés à l'infrastructure CODA, alors qu'un site a abandonné en raison de ressources informatiques limitées, et deux sites n'ont pas encore fourni de données sur les patients. Au moment de la publication de ce document, la cohorte de l'étude de viabilité de CODA comprend 1091540 patients, avec un total de 46181904 objets de la norme FHIR et 377716 études d'imagerie.

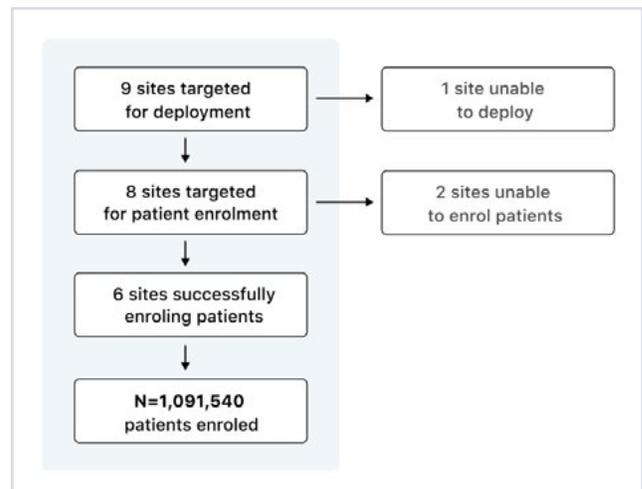


Figure 5 : Organigramme sur le recrutement des sites et des patients dans l'étude de viabilité (source : CITADEL)

Le code de logiciel, la documentation et les documents techniques concernant la plateforme CODA ont été publiés sous licence GPL v3 (www.coda-platform.com). Un ensemble de modèles du standard FHIR ont été élaborés pour aider les utilisateurs à migrer les données à partir des anciens formats de stockage. Une spécification de référence de l'API a été élaborée pour orienter la mise en œuvre des différentes composantes de la plateforme. Un guide de déploiement a été conçu pour faciliter la création d'environnements de test (*sandbox*). Un Cadre de sécurité des données a été créé pour régir les pratiques de mise en œuvre relatives à l'authentification et à l'autorisation des utilisateurs, ainsi qu'à la protection des données.

L'initiative GAIA-X en Europe centrée sur l'Allemagne

Sommaire

GAIA-X est une initiative européenne lancée en 2019 par l'ancien ministre allemand de l'Économie, Peter Altmaier, et son homologue

7. <https://github.com/coda-platform/guides-and-policies/tree/main/policies/governance>

français, Bruno Le Maire. Parallèlement aux efforts visant à créer un espace européen de données de santé, la collaboration franco-allemande cherchait à approfondir les partenariats dans le domaine du partage des données et de l'entraînement de l'IA grâce à une infrastructure de données ouverte et sécurisée susceptible de préserver la souveraineté numérique de l'Europe. Dès 2022, la mise en œuvre de GAIA-X a commencé par la création d'espaces de données et de services associés, comme l'espace commun des données pour la mobilité. Le ministère fédéral de l'Économie et du Climat en l'Allemagne a facilité ces progrès grâce à son offre de financement intitulée « Innovative and Practical Applications and Data Spaces in the Gaia-X Digital Ecosystem » (ou Applications et espaces de données novateurs et pratiques dans l'écosystème numérique GAIA-X, en français). Pour mettre en œuvre l'initiative GAIA-X et concevoir des solutions numériques novatrices dotées d'un important potentiel de marché, 11 projets ont été sélectionnés.

GAIA-X jette les bases d'une infrastructure de données autonome, unifiée et transparente qui respecte les principes de l'Union européenne. Il s'agit d'une approche stratégique qui favorise la collaboration entre diverses parties prenantes pour construire un environnement de données conforme à la réglementation européenne et favorisant la confiance (ministère fédéral de l'Économie et du Climat, 2021; Otto et Burnmann, 2021).

Plus précisément, l'espace des données de santé englobe des capacités partagées et des infrastructures fédérées, où l'on peut accéder aux données de manière granulaire et sélective. À l'avenir, ces espaces de données fédérés seront mis en œuvre à des échelons régionaux,

nationaux et européens. L'objectif est de veiller à ce que l'espace contribue à l'amélioration des soins aux patients tout en permettant l'utilisation secondaire des données à l'échelle d'une cohorte ou d'une population. Cela permettra d'établir une chaîne de valeur des données reliant les propriétaires et les utilisateurs de données au sein d'un écosystème européen de la santé à la fois vaste et complexe.

Les principaux éléments de l'infrastructure de GAIA-X se composent notamment :

- > de services d'identité et de maintien de la confiance pour une gestion fédérée de l'identité des individus et des organisations;
- > d'un catalogue fédéré pour publier les données concernant les bases enregistrées, les consentements et les requêtes de données;
- > de services d'échange de données souverains qui servent à gérer les captures de données, le consentement, les services infonuagiques et en périphérie (cloud and edge services), ainsi que les services de recherche et d'accès aux données.
- > de services de mise en conformité pour la gestion des droits, de l'intégration de nouvelles applications et de la certification des technologies.

La structure organisationnelle de GAIA-X repose sur trois piliers fondamentaux : 1) l'Association GAIA-X pour l'infonuagique et l'infrastructure (AISBL), 2) le Centre national GAIA-X, et 3) la Communauté GAIA-X. Le centre allemand de GAIA-X sert de point de contact principal pour les entreprises, les organisations et les individus en Allemagne qui cherchent à obtenir plus de renseignements sur l'initiative ou à s'engager dans la communauté du libre accès.

Pour en savoir plus :

<https://www.data-infrastructure.eu/GAIA-X/Navigation/EN/Home/home.html>

https://gaia-x.eu/wp-content/uploads/2022/05/Gaia-X-Event-Report_Health-Data-Space-Event-4_4_2022.pdf

https://healthmanagement.org/uploads/article_attachment/gaia-x-federated-data-infrastructure-the-future-of-data-management.pdf

Cas d'utilisation en Allemagne : Health-X dataLOFT

Health-X, soutenu par le ministère fédéral de l'Économie et du Climat en Allemagne, vise à établir une plateforme validée, transparente et interconnectée pour les données de santé, appelée dataLOFT. Rassemblant plusieurs acteurs publics et privés dans le domaine des données de santé en Allemagne, **Health-X dataLOFT vise à assurer la conformité avec les normes GAIA-X tout en améliorant l'accessibilité des données de santé.** L'objectif principal du projet est d'utiliser des données provenant de deux

secteurs clés de la santé : les établissements de soins primaires comme les hôpitaux et les cliniques médicales, ainsi que les institutions privées qui peuvent développer des applications et des capteurs de données. Le projet vise également à placer les citoyens au cœur de la démarche, en soulignant leur pouvoir dans les décisions concernant la collecte, l'utilisation et le contrôle de leurs données de santé. Les citoyens disposent ainsi de l'autonomie nécessaire pour décider quelles données peuvent être rendues accessibles pour leur santé personnelle, les soins médicaux et la recherche, sur la base de leurs propres valeurs et préférences.

À travers quatre cas d'utilisation, le projet HEALTH-X dataLOFT construit un espace de données propre au domaine de la santé. Cet espace vise à servir de base à ces cas d'utilisation, en abordant des questions importantes liées à l'autonomisation des citoyens, à la promotion de la santé et la prévention des maladies, à l'accompagnement du vieillissement en bonne santé et à l'amélioration des soins cliniques.

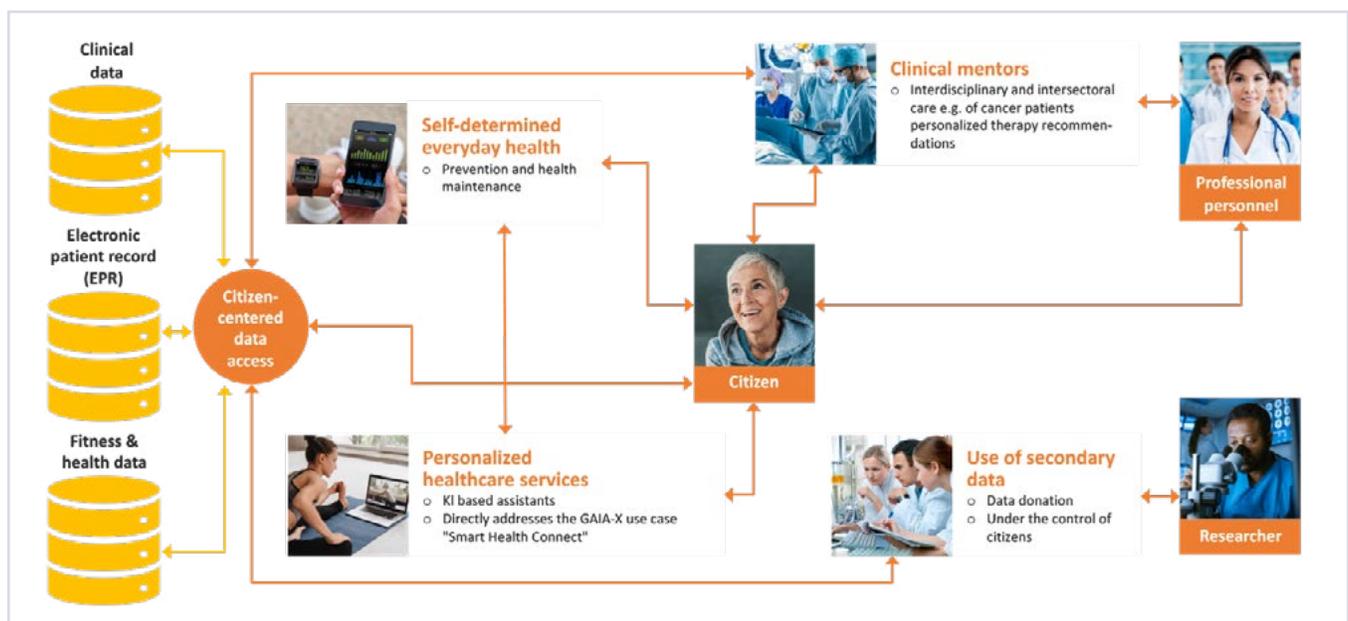


Figure 6 : Cas d'utilisation dans le cadre pour lesquels la plateforme dataLOFT fera l'objet d'études (source : [Health-X 2023](#))

Pour en savoir plus :

<https://www.computer.org/csdl/magazine/mu/2022/01/09770011/1D830WmTsDS>

<https://www.isst.fraunhofer.de/en/business-units/healthcare/projects/HEALTH-X-dataLOFT.html>

Le projet BORN en Bavière (Allemagne)

Lancé en 2022, le projet BORN est le fruit d'une collaboration entre le Centre bavarois de recherche sur le cancer (BZKF) et les instituts de radiologie des six hôpitaux universitaires de Bavière. Il a été lancé avec la participation de Markus Blume, ministre d'État de Bavière pour les sciences et les arts, et de Klaus Holetschek, ministre d'État de Bavière pour la santé et les soins. Mint Medical et Brainlab travaillent en étroite collaboration avec les hôpitaux universitaires et le Centre BZKF pour mettre en place un système cohérent et structuré permettant de connecter les rapports d'imagerie oncologique entre les établissements. Ils élaboreront également une infrastructure informatique sécurisée pour faciliter la saisie et l'échange de données.

Dans un premier temps, des modèles standardisés seront élaborés pour six entités différentes afin d'assurer une documentation uniforme des cas de cancer dans tous les hôpitaux universitaires bavarois. Après une évaluation clinique dans ces hôpitaux, la collecte et l'évaluation standardisées des données d'imagerie pourront être étendues à d'autres centres et pratiques au bénéfice des patients de toute la Bavière. Ceci permettra de constituer une base de données sans précédent, couvrant une population de 13 millions d'habitants, qui pourra être utilisée pour des études prospectives et rétrospectives.

Le but est ainsi de créer un ensemble de données unique pour le développement de biomarqueurs basés sur l'imagerie et de techniques d'apprentissage automatique. Cela permettra de recueillir et de réutiliser des données de santé structurées et exhaustives tout en préservant leur confidentialité, ainsi que d'améliorer l'utilisation des données de santé pour le traitement des patients, la recherche, l'élaboration et la mise en œuvre des politiques publiques.

Dans le cadre du projet BORN, la collecte des données est d'ailleurs décentralisée et est confiée à des cliniques et hôpitaux, ce qui garantit la confidentialité et la bonne gestion des données. Les centres demeurent ainsi propriétaires des données qu'ils fournissent, ce qui leur permet de choisir leur modalité de partage, allant de l'utilisation générale pour la recherche à des projets de collaboration précis.



Les conditions de réussite de l'apprentissage fédéré et de l'analyse décentralisée pour la recherche fondée sur les données à l'échelle internationale

À la lumière des projets présentés dans la section précédente et de l'analyse des avantages et des défis relatifs à l'apprentissage fédéré et à l'analyse décentralisée, plusieurs conditions doivent être réunies pour que ces démarches viennent consolider et faciliter la recherche fondée sur les données de santé à l'échelle internationale :

> **Principes et valeurs fondamentaux :**

Ils doivent être partagés par tous les participants d'une infrastructure d'apprentissage fédéré ou d'analyse décentralisée afin de garantir que les cadres et les processus techniques et de gouvernance reflètent les principaux objectifs et intérêts des participants, qu'il s'agisse d'organisations ou d'individus. Voici quelques exemples de ces principes et valeurs : la science ouverte, les principes FAIR (pour des données facilement trouvables, accessibles, interopérables et réutilisables) et la priorité accordée à la recherche d'intérêt public.

- > **Conformité à la réglementation et alignement juridique et éthique :** Les plateformes d'apprentissage fédéré et d'analyse décentralisée doivent se conformer aux exigences juridiques, notamment aux lois sur la protection des données (par exemple, le RGPD, la LPRPDE, HIPAA), aux réglementations en matière d'IA (par exemple, le Règlement sur l'IA de l'UE, le projet de loi C-27 au Canada), ainsi qu'aux lignes directrices plus vastes relatives à l'éthique de la recherche impliquant des sujets humains. Ces réglementations et ces lignes directrices peuvent toutefois varier d'une juridiction à une autre; c'est pourquoi un niveau minimum d'alignement sur les plans juridique et éthique entre les pays et les régions est nécessaire pour garantir que les plateformes d'apprentissage fédéré et d'analyse décentralisée puissent être mises en œuvre à l'échelle internationale.
- > **Priorité donnée à l'interopérabilité et à la standardisation des données :** La qualité, l'accessibilité et la facilité d'utilisation des données entre les participants d'une plateforme d'apprentissage fédéré ou d'analyse décentralisée dépendent fortement de l'existence de formats de données, de normes, de terminologies et de protocoles partagés par tous. L'interopérabilité entre les systèmes d'information et la standardisation des données sont essentielles pour garantir la compatibilité et l'harmonisation entre les sources de données. En outre, la documentation éthique doit être standardisée au moyen de protocoles et de formulaires partagés décrivant ce qu'est un consentement éclairé, les procédures d'acceptation et de refus du partage de données, les conditions d'utilisation, l'utilisation des données faite par la recherche, etc.

> **Cadres et processus de sécurité rigoureux :**

Une plateforme d'apprentissage fédéré ou d'analyse décentralisée doit intégrer des mécanismes robustes de protection de la vie privée afin de préserver la confidentialité des patients, de se conformer aux réglementations en matière de protection des données et d'empêcher les accès non autorisés ou les fuites de données. Ces mécanismes doivent être en harmonie avec les politiques de protection de la vie privée et s'appuyer sur les standards en vigueur les plus rigoureux pour vérifier l'identité et le rôle des utilisateurs de données, et contrôler l'utilisation de celles-ci.

> **Interfaces de programmation d'applications (API) standardisées :**

Les API peuvent simplifier l'accès aux données stockées dans différents systèmes qui sont associés à une plateforme d'apprentissage fédéré ou d'analyse décentralisée. Ces API permettent aux chercheurs et aux autres utilisateurs d'interroger et d'analyser des ensembles de données distribuées, et aux fiduciaires de données de gérer efficacement les données qui leur sont confiées.

> **Amélioration de l'accessibilité aux données et du contrôle par les patients :**

La confiance à l'égard d'une plateforme d'apprentissage fédéré ou d'analyse décentralisée repose en grande partie sur l'assurance que les patients accordent à la sécurité des données et aux avantages qu'elle offre pour l'amélioration des soins et les progrès de la recherche. Tout comme les professionnels et les chercheurs, les patients doivent avoir la possibilité d'accéder aux données et aux résultats des analyses et de la recherche, et doivent pouvoir conserver le contrôle sur l'utilisation des données. Cet objectif peut être atteint grâce à des processus de consentement numérique, à des interfaces accessibles par des personnes issues de langues et de cultures différentes et à des pratiques dynamiques qui permettent aux patients de contrôler la manière dont leurs données sont utilisées tout en pouvant découvrir les avantages qui découlent de ces utilisations.



Remerciements

Ce travail a été mené par le *RLS-Digital Health* ainsi qu'avec d'autres experts de leurs régions. Leur expertise a été sollicitée en vue de la rédaction et de la révision du présent livre blanc, afin d'obtenir des informations pertinentes, des exemples et des références utiles concernant le développement et la mise en œuvre de FL/FA dans les soins de santé.

Vous trouverez ci-après le nom des experts qui ont été consultés lors des entretiens qualitatifs semi-structurés qui ont eu lieu de février à avril 2023 :

- > Michael Chassé, Centre de recherche du Centre hospitalier de Montréal, Québec, Canada
- > Philippe Després, Université Laval, Québec, Canada
- > Vincent Ferretti, Centre hospitalier universitaire Sainte-Justine, Québec, Canada
- > Louis Mullie, Centre de recherche du Centre hospitalier de Montréal, Québec, Canada
- > Jessica Kissinger, University of Georgia, Géorgie, États-Unis
- > Jaewood Lee, University of Georgia, Géorgie, États-Unis
- > WenZhan Song, University of Georgia, Géorgie, États-Unis
- > Peter Zinterhof, LRZ et Salzburg Federal Hospital, Bavière et Haute-Autriche, Allemagne et Autriche
- > Bjoern Eskofier, Friedrich-Alexander-Universität, Bavière, Allemagne
- > Christian Wachinger, Technical University of Munich, Bavière, Allemagne
- > Michael Giretzlehner, RISC Software GmbH, Haute-Autriche
- > Mohit Kumar, Software Competence Center Hagenberg, Haute-Autriche
- > Bernhard Moser, Software Competence Center Hagenberg, Haute-Autriche
- > Agma Traina, University of Sao Paulo, São Paulo, Brésil
- > Changbin Yu, Shandong First Medical University, Jinan (Shandong), Chine

Le groupe RLS-Digital Health tient à souligner le soutien financier du Fonds de recherche du Québec – Nature et technologies et le Consortium Santé Numérique de l'Université de Montréal qui a permis la réalisation du livre blanc.

Le groupe remercie également Stéphanie Hauschild pour la conception graphique et Cécile Petitgand, PDG de Data Lama, pour son travail remarquable dans la coordination de la production de ce livre blanc.

Références

SAWSAN, Abdulrahman, et collab., « A survey on federated learning: The journey from centralized to distributed on-site learning and beyond », IEEE Internet of Things Journal, vol. 8, no 7, 2020, p. 5476–5497.

BRAGAZZI, Nicola Luigi, et collab., « How big data and artificial intelligence can help better manage the COVID-19 pandemic », International journal of environmental research and public health, vol. 17, no 9, 2020, 3176 p.

DAGLIATI, Arianna, et collab., « Health informatics and EHR to support clinical research in the COVID-19 pandemic: an overview », Briefings in bioinformatics, vol. 22, no 2, 2021, p. 812–822.

Ministère fédéral de l'Économie et du Climat « GAIA-X: Eine vernetzte Datenstruktur für ein europäisches digitales Ökosystem », 2021
En ligne : <https://www.bmwk.de/Redaktion/DE/Dossier/gaia-x.html>

FIUME, Marc, et collab., « Federated discovery and sharing of genomic data using Beacons » dans la revue Nature biotechnology, vol. 37, no 3, 2019, p. 220–224.

FORNELL, Dave, « FDA has now cleared more than 500 healthcare AI algorithms », 2023,
En ligne : <https://healthexec.com/topics/artificial-intelligence/fda-has-now-cleared-more-500-healthcare-ai-algorithms>

HONGSHENG, Hu, et collab., « Source inference attacks in federated learning », Conférence 2021 IEEE International Conference on Data Mining (ICDM), 2021, p. 1102–1107

MADHURA, Joshi, Ankit PAL et Malaikannan SANKARASUBBU, « Federated learning for healthcare domain–Pipeline, applications and challenges », ACM Transactions on Computing for Healthcare, vol. 3, no 4, 2022, p. 1–36.

KAIROUZ, Peter, et collab., « Advances and open problems in federated learning », Foundations and Trends® in Machine Learning, vol. 4, no 1, 2021, p. 1–210.

KNOPPERS, Bartha Maria, « Framework for responsible sharing of genomic and health-related data », The HUGO journal, vol. 8, no 1, décembre 2014.

WENQI, Li, et collab., « Privacy-preserving federated brain tumour segmentation », Atelier Machine Learning in Medical Imaging: 10th International Workshop, MLMI 2019, Colloque 10, 13 octobre 2019, Springer International Publishing.

MAMMEN, Priyanka Mary, « Federated learning: Opportunities and challenges », arXiv preprint arXiv:2101.05428, 2021.

MCCRADDEN, Melissa D., Elizabeth A. STEPHENSON et James A. ANDERSON, « Clinical research underlies ethical integration of healthcare artificial intelligence », Nature Medicine, vol. 26, no 9, 2020, p. 1325–1326.

MEHTA, Nishita, Anil PANDIT et Sharvari SHUKLA, « Transforming healthcare with big data analytics and artificial intelligence: A systematic mapping study », Journal of biomedical informatics, vol. 100, décembre 2019, 103311.

MOHRI, Mohri, Gary SIVEK et Ananda Theertha SURESH, « Agnostic federated learning », International Conference on Machine Learning, mai 2019, 30 p.

MORLEY, Jessica, et collab., « The ethics of AI in health care: a mapping review », *Social Science & Medicine*, vol. 260, septembre 2020, 113172.

MOTHUKURI, V., et collab., « A survey on security and privacy of federated learning », *Future Generation Computer Systems*, vol. 115, février 2021, p. 619–640.

OTTO, B., et A. BURMANN, « European data infrastructures: approaches and tools for using data for the benefit of individuals and communities », *Computer Science Spectrum*, vol. 44, no 4, 2021, p. 283–291.

PRICE II, W. Nicholson, et I. Glenn COHEN, « Privacy in the age of medical big data », *Nature medicine*, vol. 25, no 1, 2019, p. 37–43.

RIEKE, Nicola, et collab., « The future of digital health with federated learning », *NPJ digital medicine*, vol. 3 no 1, article 119, 2020.

ROS, Francisco, et collab., « Addressing the Covid19 pandemic and future public health challenges through global collaboration and a datadriven systems approach », *Learn Health Syst.*, vol. 5, no 1, décembre 2021.

ROY, Abhijit Guha, et collab., « Braintorrent: A peer-to-peer environment for decentralized federated learning », *arXiv preprint arXiv:1905.06731*, 2019.

SHAW, James, et collab., « Artificial intelligence and the implementation challenge », *Journal of medical Internet research*, vol. 21, no 7, 2019, e13659.

SHELLER, Micah J., et collab., « Federated learning in medicine: facilitating multi-institutional collaborations without sharing

patient data », *Scientific Reports*, vol. 10, no 1, 2020, p. 1–12.

TOLPEGIN, Vale, Stacey TRUEX, Mehmet Emre GURSOY et Liu LING, « Data poisoning attacks against federated learning systems », colloque *Computer Security–ESORICS 2020: 25th European Symposium on Research in Computer Security*, ESORICS 2020, Proceedings, Part I: 25th, septembre 2020, Guildford, Royaume-Uni, Springer International Publishing, p. 480–501.

XU, Jie, et collab., « Federated learning for healthcare informatics », *Journal of Healthcare Informatics Research*, vol. 5, 2021, p. 1–19.

XUEFEI, Yin, Yanming ZHU et Jiankun HU, « A comprehensive survey of privacy-preserving federated learning: A taxonomy, review, and future directions », *ACM Computing Surveys (CSUR)*, vol. 54, no 6, 2021, p. 1–36.

Apprentissage fédéré et analyse décentralisée des données pour une recherche collaborative en santé à l'échelle nationale et internationale

Livre blanc du
RLS-Digital Health

