

Approche multi-omique du suivi et de la thérapie personnalisée de l'insuffisance rénale chronique en utilisant des transformers hybrides

Projet de recherche présenté par

Ngoue David

**Étudiant en Master 2 en Intelligence Artificielle et
Big Data à Keyce Informatique**

655103813

ngouedavidrogerlyannick@gmail.com

Abdouraman Bouba Dalil
Tel : 699334228
Encadrant académique

Diffouo Tazo Evarist
Tel : 697481969
Responsable académique

27 octobre 2024

Table des matières

Introduction

Contexte

L'insuffisance rénale chronique (IRC) est une pathologie progressive qui affecte la capacité des reins à filtrer les déchets du sang, entraînant une accumulation de toxines dangereuses dans l'organisme. Elle touche des millions de personnes dans le monde et constitue un problème de santé publique majeur. L'IRC peut être causée par diverses maladies sous-jacentes, telles que le diabète, l'hypertension artérielle ou encore des infections rénales, et elle évolue souvent vers une insuffisance rénale terminale nécessitant une dialyse ou une transplantation rénale.

Un des défis majeurs dans la gestion de l'IRC est l'hétérogénéité des patients. Chaque individu peut réagir différemment au traitement, en raison de facteurs génétiques, environnementaux, ou liés à son style de vie. D'où la nécessité de traitements personnalisés pour ralentir la progression de la maladie et améliorer la qualité de vie des patients. C'est ici qu'intervient l'approche multi-omique, qui englobe l'intégration de différents types de données biologiques comme les génomes (ensemble des gènes), les transcriptomes (ARN messagers), les protéomes (protéines exprimées), et les métabolomes (métabolites).

Cette approche permet une vision plus holistique du patient, en analysant simultanément plusieurs niveaux biologiques pour découvrir des biomarqueurs pertinents, identifier des profils pathologiques uniques et ajuster les traitements en conséquence. Cependant, ces données massives et complexes posent des défis en matière d'analyse et de traitement, ce qui ouvre la porte à l'utilisation des méthodes d'intelligence artificielle (IA), notamment les modèles basés sur les transformers.

Problématique

Actuellement, les méthodes de suivi et de thérapie de l'IRC reposent principalement sur des indicateurs standards comme la clairance de la créatinine ou la protéinurie. Bien que ces mesures offrent une indication globale de la fonction rénale, elles ne tiennent pas compte des spécificités génétiques et moléculaires des patients. De plus, les traitements standardisés manquent de précision et d'efficacité pour certains sous-groupes de patients. Cela aboutit à une gestion sous-optimale de la maladie, avec des résultats cliniques variables, et parfois un retard dans la détection de la progression de la maladie.

Les approches basées uniquement sur un type de données biologique, comme les données génétiques, ont montré leurs limites en termes de précision. Elles ne permettent pas une compréhension globale de la dynamique de la maladie. De plus, l'utilisation de méthodes conventionnelles pour analyser ces données multi-omiques reste souvent insuffisante pour capturer les interactions complexes entre différents niveaux biologiques.

Les modèles IA classiques, bien qu'utilisés dans le domaine de la santé, peinent à intégrer efficacement ces grandes quantités de données hétérogènes et à fournir des prédictions robustes. C'est dans ce contexte que l'utilisation des transformers hybrides apparaît comme une solution innovante. Ils sont capables de traiter des données multi-omiques et d'intégrer des informations provenant de diverses sources, ce qui permettrait de mieux personnaliser le suivi et les thérapies de l'IRC.

Objectifs

L'objectif principal de cette étude est de proposer une approche innovante pour le suivi et la personnalisation des thérapies chez les patients atteints d'insuffisance rénale chronique, en exploitant les données multi-omiques et en utilisant des modèles de transformers hybrides. Plus précisément, nous souhaitons :

- Intégrer des données issues de différentes "omiques" (génomique, transcriptomique, protéomique et métabolomique) pour identifier des biomarqueurs uniques associés à la progression de l'IRC.
- Développer un modèle de transformers hybrides capable de traiter ces données complexes, pour améliorer la prédiction de l'évolution de la maladie et personnaliser les interventions thérapeutiques.
- Comparer les résultats de cette approche avec ceux des méthodes traditionnelles afin d'évaluer ses avantages en termes de précision et d'efficacité clinique.

Chapitre 1

Revue de la littérature

L'insuffisance rénale chronique (IRC) est une pathologie caractérisée par une détérioration progressive et irréversible de la fonction rénale, conduisant à une incapacité des reins à éliminer efficacement les déchets du corps. Parmi ces déchets, on trouve la créatinine, un marqueur clé utilisé pour évaluer la fonction rénale.

1.1 Source de la créatinine

La créatinine est un sous-produit du métabolisme de la créatine, une molécule principalement présente dans les muscles. La créatine est cruciale pour le métabolisme énergétique musculaire, jouant un rôle dans la régénération rapide de l'ATP, qui fournit de l'énergie aux cellules musculaires. Lorsqu'elle est utilisée par les muscles, la créatine est convertie en créatinine. Cette créatinine est ensuite libérée dans la circulation sanguine, d'où elle est filtrée par les reins avant d'être éliminée dans l'urine.

Contrairement à la créatine, qui peut être stockée temporairement dans les muscles, la créatinine n'est pas réutilisée et doit être constamment éliminée par les reins. Ainsi, sa concentration dans le sang est généralement stable, ce qui en fait un marqueur fiable pour évaluer la fonction rénale. Une élévation du taux de créatinine sanguine peut suggérer une réduction de la fonction rénale, en particulier dans le contexte de l'IRC.

1.2 Importance de la créatinine

La créatinine est l'un des principaux marqueurs utilisés pour évaluer la fonction rénale. La mesure de son taux dans le sang (créatininémie) et son élimination dans l'urine permettent de calculer la clairance de la créatinine, qui est une estimation du taux de filtration glomérulaire (TFG), un indicateur clé de la santé rénale. La clairance de la créatinine peut être calculée à partir des concentrations sanguines et urinaires de créatinine, selon la formule suivante :

$$\text{Clairance de la créatinine} = \frac{\text{Créatinine urinaire (mg/dL)} \times \text{Volume urinaire (mL/min)}}{\text{Créatinine plasmatique (mg/dL)}}$$

Le taux de filtration glomérulaire est un outil essentiel pour diagnostiquer et suivre l'évolution de l'insuffisance rénale chronique (IRC). Une baisse du TFG indique une détérioration de la capacité des reins à filtrer le sang, ce qui peut entraîner l'accumulation

de toxines dans l'organisme. La créatinine étant un produit de dégradation musculaire relativement constant, elle est considérée comme un indicateur fiable pour surveiller la fonction rénale.

1.3 Insuffisance rénale chronique et multi-omique

L'utilisation des données multi-omiques dans le cadre de l'insuffisance rénale chronique (IRC) est relativement récente, mais elle s'appuie sur des avancées significatives dans le domaine de la médecine de précision. Les études menées par Eddy et al. (2020) dans l'article intitulé **"Integrated multi-omics approaches to improve classification of chronic kidney disease"** [eddy2020integrated] montrent que l'intégration des données génomiques, transcriptomiques, protéomiques et métabolomiques permet de mieux classer les patients souffrant d'IRC. Cette classification multi-omique permet d'identifier des sous-groupes de patients présentant des caractéristiques spécifiques, facilitant ainsi la personnalisation des traitements. Leur recherche met en lumière l'importance de combiner plusieurs couches d'informations biologiques pour comprendre pleinement la progression de la maladie.

De plus, l'étude menée par Hill et al. (2022) intitulée **"Harnessing the Full Potential of Multi-Omic Analyses to Advance the Study and Treatment of Chronic Kidney Disease"** [hill2022harnessing] aborde l'importance de l'intégration des données multi-omiques dans l'étude de la progression de l'IRC. Les auteurs soulignent que chaque "omique" capture un aspect différent de la maladie, comme les altérations transcriptomiques ou les réponses protéomiques spécifiques. Ils discutent également des défis liés à l'analyse de ces données complexes, notamment la gestion de la dynamique temporelle et la nécessité de nouvelles méthodes pour l'intégration des différents types de données.

Enfin, Ramírez Medina et al. (2023), dans leur article intitulé **"Proteomic signature associated with chronic kidney disease (CKD) progression identified by data-independent acquisition mass spectrometry"** [medina2023proteomic], identifient des signatures protéomiques spécifiques associées à la progression de l'IRC. Ces signatures permettent de mieux comprendre les variations individuelles dans les réponses aux traitements, et ouvrent la voie à des stratégies thérapeutiques plus ciblées. Les auteurs montrent également comment la spectrométrie de masse est utilisée pour identifier ces biomarqueurs protéomiques de manière plus précise.

1.4 Transformers dans la bio-informatique

Les transformers sont des modèles de deep learning qui ont révolutionné le traitement des données séquentielles, notamment dans les domaines du traitement du langage naturel (NLP). Ils se basent sur des mécanismes d'attention qui permettent de capturer des relations complexes entre différents éléments d'une séquence, sans se limiter aux dépendances locales comme le font d'autres modèles de type LSTM (Long Short-Term Memory) ou

CNN (Convolutional Neural Networks).

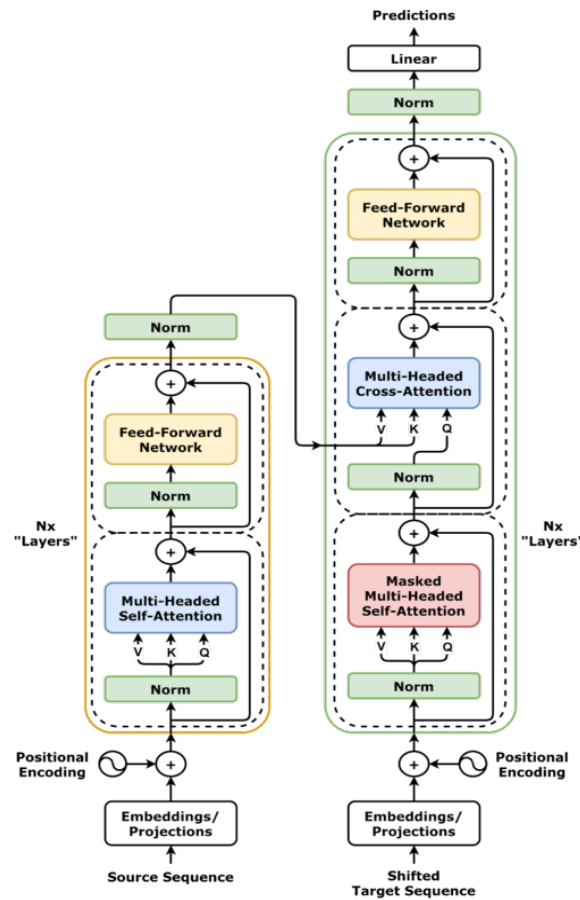


FIGURE 1.1 – Architecture d'un transformer

Dans la bio-informatique, ces modèles ont récemment été adaptés pour traiter des données génétiques et autres données biologiques séquentielles, en particulier grâce à leur capacité à intégrer et à analyser simultanément des données de différentes échelles. Par exemple, des transformers ont été utilisés pour prédire l'expression des gènes, détecter des mutations dans les séquences ADN, ou encore modéliser des interactions protéine-protéine. Grâce à leur flexibilité, les transformers sont bien adaptés à la nature complexe et multi-échelle des données biologiques, telles que celles produites par les approches multi-omiques.

Dans leur travail, Jumper et al. (2021) [jumper2021alphafold] ont montré l'efficacité des transformers dans la prédiction de la structure des protéines, notamment avec l'outil AlphaFold. Cette capacité à prédire des structures biologiques à partir de données séquentielles a ouvert de nouvelles perspectives pour la modélisation des interactions moléculaires, un domaine clé dans la compréhension des maladies chroniques telles que l'IRC. Les transformers permettent d'analyser simultanément plusieurs couches d'information biologique, ce qui est essentiel dans le cadre d'une approche multi-omique.

1.5 Transformers hybrides

L'utilisation de transformers hybrides consiste à combiner ces modèles avec d'autres architectures neuronales ou algorithmes, afin de tirer parti des avantages complémentaires de chaque approche. Dans le cadre de cette étude, les transformers hybrides sont envisagés pour intégrer plusieurs couches de données omiques (génomique, protéomique, transcritomique, etc.) et fournir des prédictions globales sur l'évolution de l'IRC.

Par exemple, un modèle hybride pourrait combiner un transformer pour le traitement des données séquentielles (comme les séquences d'ADN) avec un réseau de neurones convolutionnels (CNN) pour analyser les images de biopsie rénale ou d'autres types de données structurales. De plus, la combinaison avec des méthodes bayésiennes pourrait permettre de mieux gérer l'incertitude liée aux petites tailles d'échantillon et aux données hétérogènes souvent rencontrées en biologie.

Cette approche permet de capturer à la fois les relations complexes au sein de chaque type de données (grâce au transformer) et les interactions entre ces données (grâce à l'hybridation avec d'autres techniques). Cela est particulièrement pertinent dans le cas de l'IRC, où les relations entre les différents biomarqueurs peuvent être non linéaires et difficiles à modéliser avec des techniques classiques.

En conclusion, les transformers hybrides représentent une opportunité pour repousser les limites des méthodes actuelles d'analyse multi-omique dans l'insuffisance rénale chronique. Ils offrent la possibilité de fournir des prédictions plus robustes, d'identifier des biomarqueurs plus précis, et de personnaliser les traitements pour chaque patient, contribuant ainsi à une gestion plus efficace de cette maladie complexe.

Chapitre 2

Objectifs spécifiques

Les objectifs spécifiques de cette recherche sont formulés pour répondre aux défis identifiés dans la problématique, tout en assurant une évaluation rigoureuse de l'efficacité de l'approche multi-omique associée à l'utilisation de transformers hybrides. Voici les principaux objectifs :

- **Analyser les données multi-omiques des patients atteints d'IRC.**

L'un des objectifs clés est de collecter et d'analyser des données provenant de différentes omiques (génomiques, transcriptomiques, protéomiques, métabolomiques, etc.) chez des patients atteints d'insuffisance rénale chronique. Ces données permettront de mieux comprendre la complexité moléculaire de l'IRC, d'identifier des biomarqueurs pertinents et d'explorer comment ces données interagissent pour influencer la progression de la maladie.

Cette étape implique non seulement la collecte de ces données, mais aussi leur pré-traitement pour assurer leur qualité et leur intégrité. Elle inclura des étapes comme la normalisation des données, la gestion des valeurs manquantes et la réduction de la dimensionnalité si nécessaire.

- **Développer et ajuster un modèle de transformers hybrides pour intégrer ces données.**

L'objectif principal est de concevoir un modèle de transformers hybrides capable d'intégrer les différentes couches de données omiques et de générer des prédictions fiables concernant l'évolution de l'IRC. Ce modèle hybride combinera des transformers avec d'autres types de réseaux neuronaux, ou éventuellement des méthodes bayésiennes ou de l'apprentissage automatique traditionnel, pour tirer parti des complémentarités entre ces approches.

Il sera également nécessaire d'ajuster et de paramétrer ce modèle en fonction des données spécifiques à l'IRC, de tester différentes architectures et d'optimiser les hyperparamètres pour obtenir les meilleures performances possibles.

- **Comparer les résultats obtenus avec des méthodes traditionnelles.**

Pour valider l'efficacité du modèle, il est indispensable de comparer ses résultats à ceux des méthodes classiques actuellement utilisées pour le suivi de l'IRC. Cela inclura des comparaisons avec des modèles statistiques (comme les régressions linéaires ou logistiques) ainsi qu'avec d'autres techniques de machine learning non

basées sur les transformers (comme les réseaux de neurones traditionnels, les forêts aléatoires, etc.).

L'objectif est de démontrer en quoi les transformers hybrides offrent une meilleure performance, notamment en termes de précision, de sensibilité et de capacité à intégrer des données hétérogènes.

- **Évaluer l'impact de l'approche sur la personnalisation des traitements.**
Enfin, il s'agira d'évaluer l'impact de l'approche multi-omique sur la personnalisation des traitements de l'IRC. Grâce à l'intégration des données omiques, nous espérons mieux personnaliser les thérapies en fonction des biomarqueurs identifiés et des caractéristiques moléculaires des patients, dans le but d'améliorer les résultats cliniques, notamment en ralentissant la progression de la maladie.

Chapitre 3

Hypothèses de recherche

3.1 Hypothèse 1

Les données multi-omiques permettent une meilleure prédiction de l'évolution de l'IRC par rapport aux méthodes traditionnelles.

Cette hypothèse suggère que l'intégration de plusieurs couches de données biologiques (génomiques, protéomiques, métabolomiques, etc.) fournira des informations plus complètes et plus précises sur la progression de l'insuffisance rénale chronique que les approches basées sur une seule dimension des données (comme l'analyse génomique seule). Il est anticipé que cette richesse en données conduira à une meilleure identification des biomarqueurs de progression et à des prédictions plus fines de l'évolution de la maladie.

3.2 Hypothèse 2

Les transformers hybrides améliorent la précision des prédictions de la progression de l'IRC par rapport aux modèles classiques d'intelligence artificielle.

Cette deuxième hypothèse propose que l'utilisation de transformers hybrides, capables de gérer des données séquentielles et de capturer des relations complexes entre les différents types de données omiques, permettra une amélioration significative des performances de prédiction par rapport aux méthodes classiques, telles que les modèles de régression linéaire ou les algorithmes de machine learning traditionnels. Il est attendu que les transformers hybrides exploitent mieux la diversité des données et permettent des prédictions plus robustes et personnalisées.

Chapitre 4

Méthodologie

4.1 Type d'étude

Cette recherche s'inscrit dans le cadre d'une étude rétrospective d'analyse de données. Les données multi-omiques proviendront de bases de données publiques et/ou de cohortes cliniques déjà existantes de patients atteints d'insuffisance rénale chronique. Une étude rétrospective est justifiée car elle permet de travailler sur des données déjà collectées, ce qui offre une économie de temps tout en assurant la diversité et la richesse des informations disponibles pour l'analyse.

4.2 Population

Les patients inclus dans cette étude seront des individus diagnostiqués avec une insuffisance rénale chronique à différents stades de la maladie. Les critères de sélection incluent :

- **Âge** : Patients âgés de 18 ans et plus.
- **Stades de l'IRC** : Tous les stades de la maladie seront inclus, allant de la phase précoce à l'insuffisance rénale terminale.
- **Disponibilité des données multi-omiques** : Une évaluation sera effectuée sur l'ensemble des patients, y compris ceux ayant des données partiellement manquantes, afin de ne pas exclure une partie importante de la population étudiée. Différentes techniques d'imputation seront utilisées pour combler ces lacunes.
- **Sexe** : Hommes et femmes seront inclus pour évaluer d'éventuelles différences biologiques dans la progression de l'IRC.

Un sous-groupe de patients avec des comorbidités telles que le diabète ou l'hypertension artérielle sera également étudié pour explorer l'impact de ces conditions sur la progression de l'IRC.

4.3 Données multi-omiques

Les types de données qui seront utilisés dans cette étude incluent :

- **Génomiques** : Séquences d'ADN, mutations génétiques, SNPs (single nucleotide polymorphisms), etc.

- **Transcriptomiques** : Données d’expression des gènes à partir de l’ARN messager.
- **Protéomiques** : Profils protéiques dérivés des techniques comme la spectrométrie de masse.
- **Métabolomiques** : Données sur les métabolites produits par les cellules, obtenues par chromatographie en phase liquide ou d’autres méthodes.

Ces données seront obtenues via des bases de données publiques, telles que TCGA (The Cancer Genome Atlas), ou des collaborations avec des hôpitaux ou des centres de recherche possédant des cohortes cliniques bien établies sur l’IRC.

4.4 Modélisation avec transformers hybrides

Le modèle utilisé sera un transformer hybride. Ce modèle comportera :

- Un transformer pour traiter les données séquentielles, comme les séquences génomiques ou transcriptomiques, qui nécessitent une gestion fine des dépendances entre les positions des éléments dans les séquences.
- Des réseaux neuronaux traditionnels ou convolutifs pour intégrer d’autres types de données moins séquentielles, comme les données métabolomiques ou les profils protéiques. Le modèle hybride permettra d’intégrer ces différentes sources d’information de manière fluide, en exploitant les points forts de chaque sous-modèle pour capturer les relations entre les différentes omiques.

Les paramètres du modèle, comme la profondeur des couches du transformer ou la taille des réseaux de neurones convolutifs, seront ajustés en fonction des résultats préliminaires sur les données d’entraînement. Une validation croisée sera effectuée pour éviter le surapprentissage.

4.5 Validation

Pour évaluer la performance du modèle de transformers hybrides, plusieurs métriques seront utilisées :

- **Précision (Accuracy)** : Pour mesurer la proportion de prédictions correctes.
- **Rappel (Recall)** : Pour évaluer la capacité du modèle à détecter correctement les patients atteints de la maladie.
- **F1-score** : Une mesure qui combine la précision et le rappel.
- **Courbe ROC et AUC (Area Under the Curve)** : Pour analyser la capacité discriminante du modèle entre les différents stades de l’IRC.

Ces mesures seront comparées à celles obtenues avec des modèles classiques pour démontrer l’efficacité de l’approche hybride. Un ensemble de test indépendant sera utilisé pour valider les résultats finaux et garantir la robustesse des prédictions.

Chapitre 5

Résultats attendus

5.1 Amélioration de la prédiction des stades de l’IRC

L’intégration des données multi-omiques devrait permettre d’améliorer la prédiction des stades de l’IRC, offrant une précision accrue par rapport aux méthodes classiques. Cependant, chaque patient réagissant différemment, il n’est pas possible d’anticiper de manière uniforme les résultats pour tous les individus. Les résultats varieront en fonction des caractéristiques omiques et cliniques de chaque patient.

5.2 Personnalisation des traitements

L’approche multi-omique combinée aux transformers hybrides permettra une personnalisation plus fine des traitements, mais cette personnalisation ne doit pas être interprétée comme garantissant des résultats cliniques toujours positifs. Chaque patient a une réponse unique, et les traitements proposés devront être continuellement ajustés selon l’évolution de la maladie et les données cliniques.

Chapitre 6

Plan d'analyse des données

L'analyse des données sera effectuée en plusieurs étapes :

6.1 Analyse statistique initiale

Dans un premier temps, une analyse descriptive des données sera réalisée pour obtenir une vue d'ensemble de la distribution des variables multi-omiques et des caractéristiques des patients (âge, sexe, stade de l'IRC, etc.). Cela inclut des statistiques de base comme la moyenne, la médiane, l'écart-type, et les corrélations entre les différentes variables omiques.

6.2 Techniques de machine learning

L'étape principale consistera à entraîner le modèle de transformers hybrides sur les données multi-omiques des patients atteints d'IRC. Le modèle sera ajusté en fonction des données, et des techniques de validation croisée seront employées pour éviter le surapprentissage. Une fois le modèle formé, il sera testé sur des ensembles de données de test indépendants.

6.3 Métriques d'évaluation de la performance

Les performances du modèle seront évaluées à l'aide des métriques suivantes :

- **Précision (Accuracy)** : Pour mesurer la proportion de prédictions correctes.
- **Rappel (Recall)** : Pour évaluer la capacité du modèle à identifier correctement les cas de progression de l'IRC.
- **F1-score** : Une mesure combinée de la précision et du rappel.
- **Courbe ROC et AUC (Area Under the Curve)** : Pour analyser la performance globale du modèle en termes de capacité discriminante.

Ces résultats seront comparés à ceux d'autres méthodes (régressions linéaires, SVM, réseaux de neurones traditionnels, etc.) afin de prouver l'efficacité des transformers hybrides.

Chapitre 7

Calendrier de recherche

Phase 1 : Revue de la littérature et définition du cadre de recherche (1-2 mois)

- Lecture d'articles scientifiques
- Identification des bases de données omiques et des cohortes de patients
- Définition des critères de sélection des données

Phase 2 : Collecte et prétraitement des données (2 mois)

- Collecte des données multi-omiques
- Nettoyage et normalisation des données
- Préparation des données pour la modélisation

Phase 3 : Développement du modèle de transformers hybrides (2 mois)

- Construction et entraînement du modèle
- Ajustement des hyperparamètres et optimisation
- Tests préliminaires sur des ensembles de validation croisée

Phase 4 : Validation du modèle et analyse des résultats (1 mois)

- Validation du modèle sur des ensembles de test
- Comparaison avec d'autres méthodes
- Analyse des performances et interprétation des résultats

Phase 5 : Rédaction et diffusion des résultats (2 mois)

- Rédaction du mémoire
- Publication des résultats (si applicable)

Le calendrier total devrait donc s'étendre sur une durée d'environ 9 à 10 mois.

Chapitre 8

Contraintes et limites

8.1 Taille de l'échantillon

La taille de l'échantillon disponible pour l'analyse des données multi-omiques pourrait être limitée, surtout si les données proviennent de bases de données publiques ou d'études cliniques spécifiques. Une petite taille d'échantillon pourrait réduire la capacité du modèle à généraliser les prédictions à une population plus large.

8.2 Difficultés liées à l'intégration des données omiques

L'intégration de différentes couches de données omiques (génomiques, protéomiques, transcriptomiques, etc.) représente un défi technique important. La nature hétérogène de ces données et leur dimensionnalité élevée peuvent compliquer la modélisation et entraîner des problèmes de surapprentissage.

8.3 Biais potentiels

Les biais dans les données, comme les biais liés à la sélection des patients ou à la collecte des données, peuvent affecter les résultats. Par exemple, des patients sélectionnés dans des bases de données spécifiques pourraient ne pas être représentatifs de la population générale des patients atteints d'IRC, limitant ainsi la généralisation des résultats.

Conclusion

Cette recherche pourrait apporter une contribution significative à la gestion et au suivi de l'insuffisance rénale chronique en intégrant des approches multi-omiques avec des transformers hybrides. En permettant une meilleure personnalisation des traitements, elle pourrait améliorer les résultats cliniques et offrir une nouvelle approche pour prédire la progression de l'IRC.

L'impact potentiel de cette étude sur la recherche médicale est considérable, car elle offre une nouvelle manière d'aborder le suivi des maladies chroniques complexes comme l'IRC. Le succès de cette étude pourrait aussi ouvrir la voie à l'application des transformers hybrides dans d'autres domaines de la médecine personnalisée, tels que le cancer ou les maladies cardiovasculaires.