

Prédiction de structures biomoléculaires complexes par AlphaFold 3

Antoine Taly¹, Alexis Verger²

¹Laboratoire de biochimie théorique, UPR 9080 CNRS, Université Paris Cité, Paris, France.

²CNRS EMR 9002 Biologie structurale intégrative, Inserm U1167 – Facteurs de risques et déterminants moléculaires des maladies liées au vieillissement (RID-AGE), Univ. Lille, Centre hospitalo-universitaire de Lille, Institut Pasteur de Lille, Lille, France.

Taly@ibpc.fr

Alexis.Verger@univ-lille.fr

L'année 2021 a été marquée par la formidable concrétisation de la prédiction de structures tridimensionnelles des protéines à partir de leurs séquences d'acides aminés grâce à AlphaFold et RoseTTAFold^[2], des systèmes d'intelligence artificielle (IA) fondés sur l'apprentissage profond par réseaux de neurones. L'IA a d'ailleurs marqué la toute récente édition 2024 des prix Nobel de physique et de chimie et de chimie décernés respectivement pour le développement de l'apprentissage profond et son application dans la prédiction des structures des protéines. Cette innovation est un accélérateur de recherche et un véritable générateur d'hypothèses. Elle transcende les méthodes traditionnelles comme la cristallographie aux rayons X et la cryo-microscopie électronique en fournissant des modèles structuraux précis qui guident et complètent l'interprétation des données expérimentales.

La prédiction de structures biomoléculaires complexes par AlphaFold 3 permet d'obtenir des prédictions plus précises. Le module Pairformer d'AlphaFold 3 remplace le module Evoformer d'AlphaFold 2 (Figure 1). L'innovation accordée par AlphaFold 3 est donc sa capacité à modéliser une large gamme de complexes moléculaires incluant quelques ligands, informations coévolutives extraites des acides nucléiques, quelques modifications post-traductionnelles (multiple sequence alignment, MSA) a été considérablement réduite. Enfin, la représentation module de confiance a été introduit pour évaluer les erreurs provenant des calculs au niveau des atomes et des paires d'acides aminés (Figure 2). Certaines de ces modifications étaient nécessaires pour prédire avec précision (ADN, ARN) et même aux modifications des complexes biomoléculaires, mais pas toutes. Tandis que d'autres sont des mises à jour des approches d'apprentissage automatique qui ont émergé et se sont répandues au cours des trois dernières années, depuis la publication d'AlphaFold 2.

Un nouveau champ des possibles

Le module de structure d'AlphaFold 3 a été optimisé d'AlphaFold 2 et a été remplacé par un module de diffusion AlphaFold 3 à presque 200 millions de paramètres (Figure 3). En tant que méthode générative, la diffusion AlphaFold 3 génère directement la génération de coordonnées atomiques, ce qui permet d'explorer des solutions très diverses. Pour éviter de générer des structures improbables dans les régions mal structurées, qui ne représentent pas la réalité (ce qu'on appelle l'hallucination), la distillation croisée All-Atom sur le benchmark Protein Data Bank (PDB) a été utilisée avec des données d'entraînement d'AlphaFold-RoseTTAFold All-Atom combinées en Multimodal^[6]. Cette stratégie permet de prédire en 3D la meilleure orientation d'un ligand (petite molécule) dans une protéine cible.

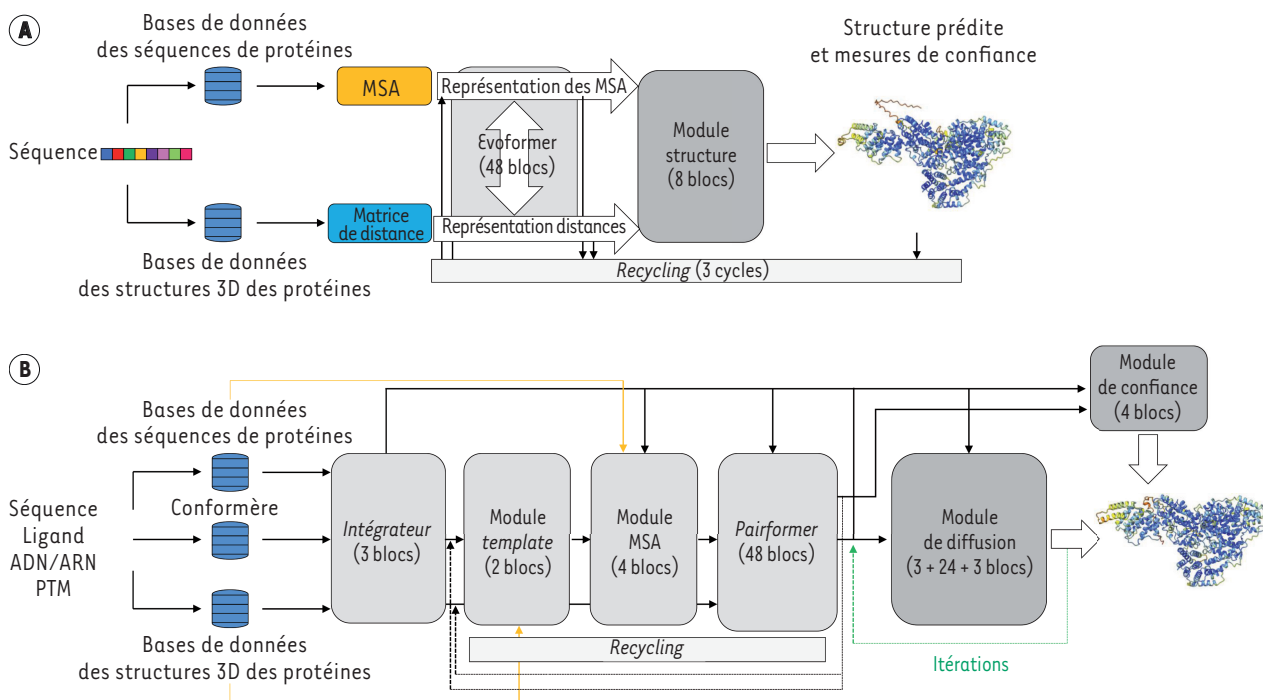


Figure 1. Architecture simplifiée de l'intelligence artificielle de AlphaFold 2 et AlphaFold 3. A. L'architecture d'AlphaFold 2 se décompose en trois modules principaux. Les données de séquences des protéines et les données expérimentales des structures 3D sont traitées par le module d'entrée, puis analysées par le module Evoformer, dont les résultats sont ensuite transmis au module de structure. B. L'architecture d'AlphaFold 3 se décompose en six modules principaux. Le module Pairformer d'AlphaFold 3 remplace le module Evoformer d'AlphaFold 2, et le module de diffusion remplace le module de structure. La représentation des alignements de séquences multiples (multiple sequence alignment, MSA) d'AlphaFold 2 (indiquée par la double flèche dans la figure 1A) n'a pas été conservée, et toutes les informations passent maintenant par la représentation matricielle par paires d'atomes ou d'acides aminés. Enfin, un module de confiance a été ajouté dans AlphaFold 3. Figure adaptée de [1, 6].

de la PDB publiées avant 2021. AlphaFold 3 est également plus précis dans la prédiction des complexes protéine-acide nucléique et des structures ARN que RoseTTAFoldNA [10] et Alchemy_RNA. En outre, AlphaFold 3 prédit plus précisément l'effet de ligands ou de certaines modifications post-traductionnelles des acides aminés et des bases d'acide nucléique (Figure 2). Cependant, aucune comparaison avec d'autres outils n'est rapportée. Le nombre de prédictions correctes, mesurées par un RMSD (root mean square deviation ; mesure utilisée pour comparer la qualité de la superposition entre la structure expérimentale et la structure prédite par AlphaFold 3) inférieure à 2 Å, varie entre 40 % pour les résidus modifiés de l'ARN et près de 80 % pour les ligands [6]. AlphaFold 3 améliore également la prédiction

des complexes protéine-protéine par rapport à AlphaFold-Multimer, notamment celle des interfaces anticorps-protéine (performance de 63 % pour AlphaFold 3, contre 30 % pour AlphaFold-Multimer).

Des possibilités malheureusement limitées

Trois ans après sa sortie, l'IA d'AlphaFold 2 a été largement adoptée et maintes fois validée expérimentalement par la communauté scientifique comme un outil révolutionnaire de prédiction de structure. AlphaFold 3 est une belle optimisation d'AlphaFold 2, mais ce nouvel outil est limité à un usage scientifique non commercial. AlphaFold 3 obtient de très bons résultats dans toutes les catégories testées, mais il convient d'attendre

les évaluations indépendantes des autres équipes scientifiques pour déterminer sa performance globale. Malheureusement, le code/modèle d'AlphaFold 3 n'est pas accessible comme l'était celui d'AlphaFold 2, l'accès étant limité à un serveur web qui impose des contraintes sur le nombre de modélisations par jour (20) et la taille des séquences d'entrée. Les modèles obtenus ne pourront pas être utilisés par la communauté académique pour faire du criblage virtuel par docking moléculaire ou entraîner d'autres modèles d'apprentissage automatique. Les autres restrictions concernent une liste très limitée de ligands (Figure 2). Compte tenu de toutes ces restrictions, il est très difficile, pour le moment, de détermi-

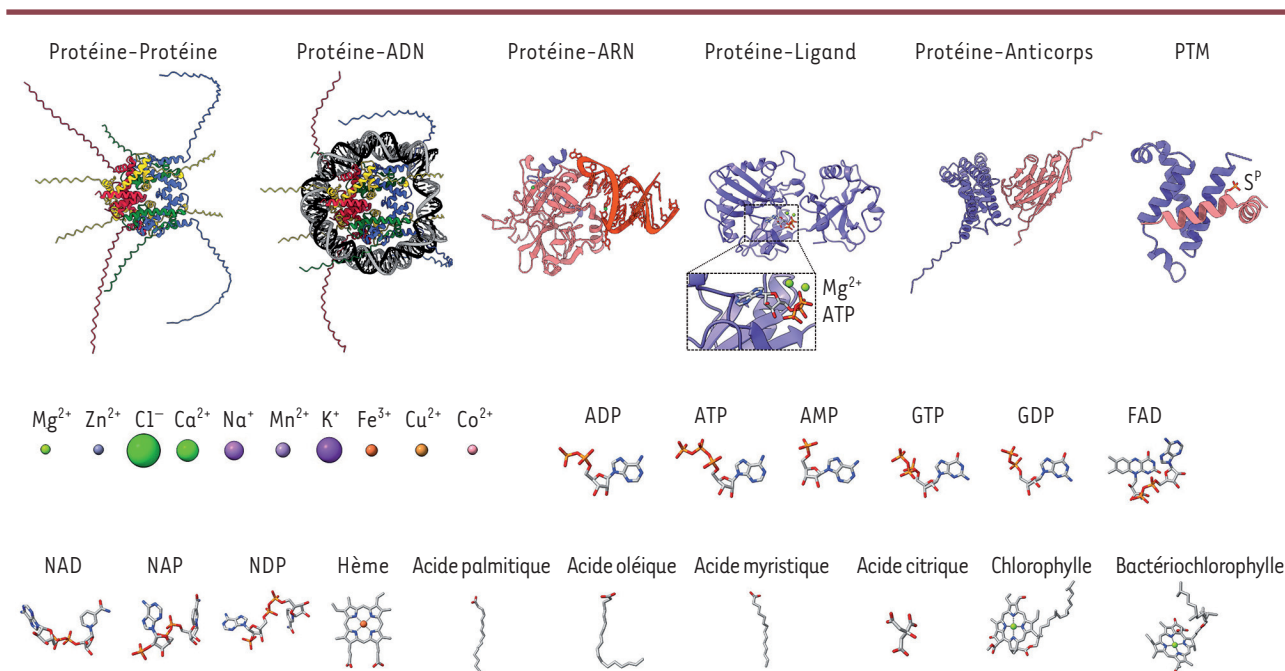


Figure 2. Exemples de complexes biomoléculaires modélisés avec AlphaFold 3. Les modèles présentés ont été réalisés sur le serveur web AlphaFold 3 (<https://alphafoldserver.com/>), et les figures ont été réalisées avec ChimeraX. PTM : post-translational modification ; S^P : résidu sérine phosphorylé. Les quelques ions et ligands disponibles sur le serveur web sont indiqués.

ner les capacités de généralisation d'AlphaFold 3, en particulier pour les autres biomolécules qui ne sont pas contenues dans la liste prédéfinie. \diamond

Prediction of complex biomolecular structures by AlphaFold 3

REMERCIEMENTS

Les auteurs souhaitent remercier Xavier Hanouille ainsi que tous les membres de la section 20 du comité national de la recherche scientifique (CoNRS) pour les discussions fructueuses autour d'AlphaFold.

LIENS D'INTÉRÊT

Les auteurs déclarent n'avoir aucun lien d'intérêt concernant les données publiées dans cet article.

RÉFÉRENCES

1. Jumper J, Evans R, Pritzel A, et al. Highly accurate protein structure prediction with AlphaFold. *Nature* 2021 ; 596 : 583-9.
2. Baek M, DiMaio F, Anishchenko I, et al. Accurate prediction of protein structures and interactions using a three-track neural network. *Science* 2021 ; 373 : 871-6.
3. Graille M, Sacquin-Mora S, Taly A. Best practices of using AI-based models in crystallography and their impact in structural biology. *J Chem Inf Model* 2023 ; 63 : 3637-46.
4. Varadi M, Bertoni D, Magana P, et al. AlphaFold protein structure database in 2024: providing structure coverage for over 214 million protein sequences. *Nucleic Acids Res* 2024 ; 52 : D368-D375.
5. Versini R, Sritharan S, Aykac Fas B, et al. A perspective on the prospective use of AI in protein structure prediction. *J Chem Inf Model* 2024 ; 64 : 26-41.
6. Abramson J, Adler J, Dunger J, et al. Accurate structure prediction of biomolecular interactions with AlphaFold 3. *Nature* 2024 ; 630 : 493-500.
7. Krishna R, Wang J, Ahern W, et al. Generalized biomolecular modeling and design with RoseTTAFold All-Atom. *Science* 2024 ; 384 : ead12528.
8. Eberhardt J, Santos-Martins D, Tillack AF, et al. AutoDock Vina 1.2.0: New docking methods, expanded force field, and Python bindings. *J Chem Inf Model* 2021 ; 61 : 3891-8.
9. Buttenschon M, Morris GM, Deane CM. PoseBusters: AI-based docking methods fail to generate physically valid poses or generalise to novel sequences. *Chem Sci* 2024 ; 15 : 3130-9.
10. Baek M, McHugh R, Anishchenko I, et al. Accurate prediction of protein-nucleic acid complexes using RoseTTAFoldNA. *Nat Methods* 2024 ; 21 : 117-21.