

# Quel modèle et quelle base de données 16S pour l'analyse du microbiome urinaire bactérien?

Marthe Leoz<sup>1,2</sup>, Johan Brière<sup>1,2</sup>, Marie Leoz<sup>1</sup>, Camille Villenave<sup>1</sup>, Martine Pestel-Caron<sup>3</sup>, and Sandrine Dahyot<sup>3</sup>

<sup>1</sup>Univ Rouen Normandie, Université de Caen Normandie, INSERM, Normandie Univ, DYNAMICURE UMR 1311, F-76000 Rouen, France; <sup>2</sup>DataScientest, 92800 Puteaux, France;

<sup>3</sup>Univ Rouen Normandie, Université de Caen Normandie, INSERM, Normandie Univ, DYNAMICURE UMR 1311, CHU Rouen, Department of Bacteriology, F-76000 Rouen, France

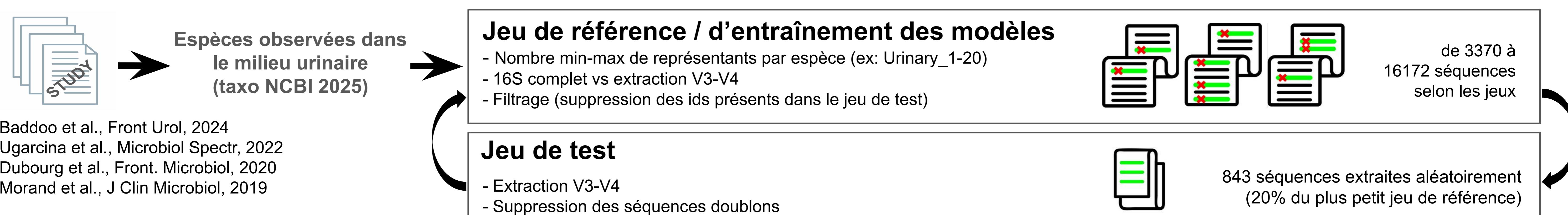
Corresponding Author: marthe.leoz@gmail.com

## Contexte

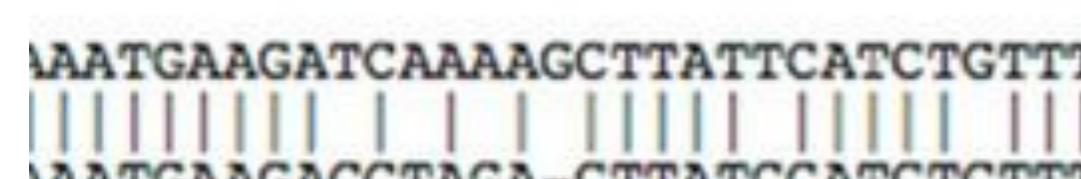
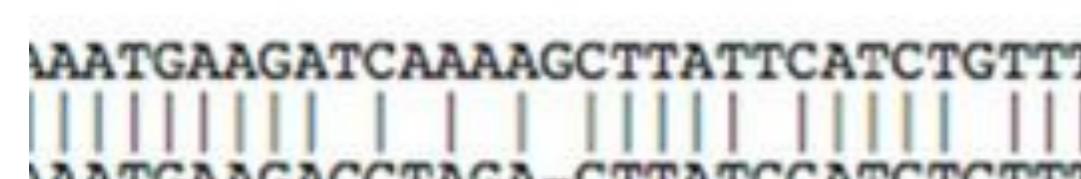
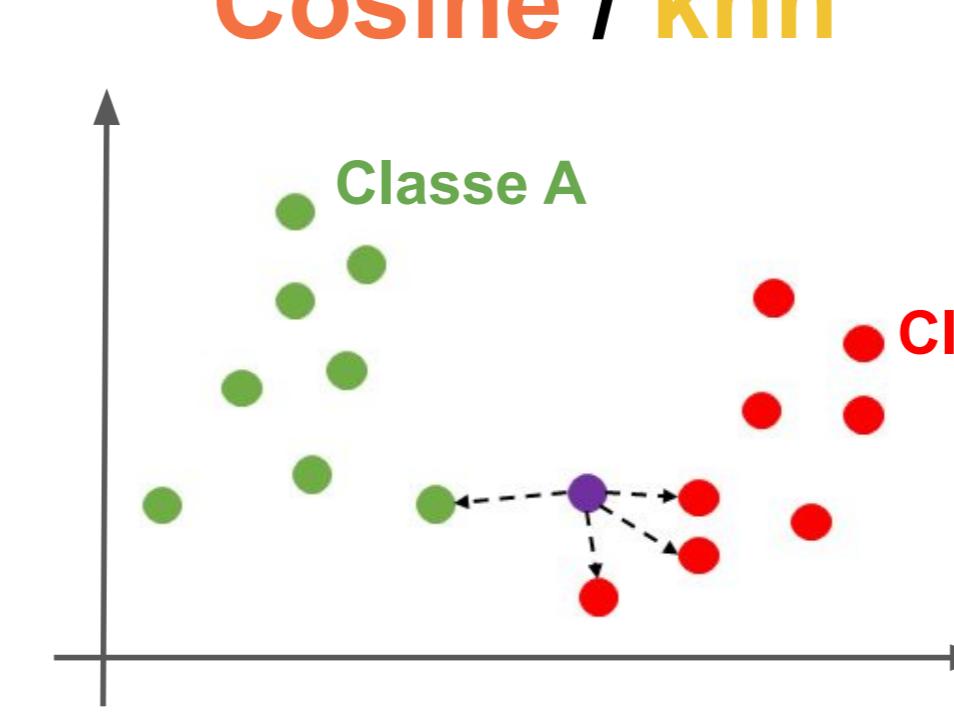
L'analyse du microbiome urinaire par metabarcoding 16S pose des difficultés de spécificité et sensibilité. Les résultats obtenus peuvent varier d'une approche à l'autre, et certains genres / espèces d'intérêt peuvent être mal identifiés. Parmi eux, on note des uropathogènes (*Escherichia coli*, *Klebsiella pneumoniae*...) et des représentants importants du microbiome (dont les différents lactobacilles).

L'objectif de ce travail est d'optimiser un outil d'analyse des séquences des régions V3-V4 du gène 16S obtenues selon le protocole standard Illumina.

## Jeux de référence et de test



## Modèles

Alignment-based models		Machine Learning (ML) models		
<b>Rapid Fuzz</b>  généraliste parcimonie   > Retournent des scores de similarité		<b>Blastn</b>  nucléotide - nucléotide extension de k-mer   > Retournent des scores de similarité		
<b>Data Pre-processing</b> 	<b>Cosine / knn</b>  > Retourne la classe la plus proche (Cosine), ou la classe majoritaire parmi les K plus proches voisins (K-Nearest Neighbors, knn)	<b>Naive Bayes (NB)</b> $P_{\text{classe} kmer} = \frac{P_{kmer \text{classe}} \times P_{\text{classe}}}{P_{kmer}}$ > Retourne la classe la plus probable		

## Résultats

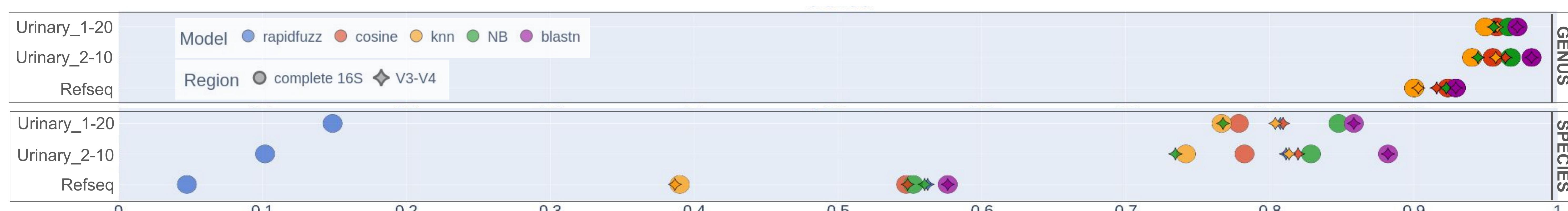


Fig1: F1-score moyen par modèle (couleur), région d'entraînement (symbole), jeu de référence/entraînement et niveau taxonomique.

L'analyse des F1-scores (Fig.1) indique un **bénéfice des jeux de référence / d'entraînement spécifiques du milieu urinaire**, avec peu ou pas d'avantage à (1) équilibrer le nombre de représentants par espèce; (2) utiliser v3-v4 vs 16S complet. A ce stade, RapidFuzz < modèles ML < Blastn.



Fig2: Apport du jeu de référence spécifique du milieu urinaire pour l'identification d'espèces d'intérêt avec Blastn.

## Discussion et perspectives

- Optimisation des modèles ML: classes mieux équilibrées (enrichir les sous-représentées), varier les tailles de k-mers et les jeux de test, le nombre de voisins (knn), l'estimateur alpha (NB), tester le modèle de vectorisation DNA-BERT
- Limite jeu de référence urinaire: pas de possibilité de découvrir de nouvelle espèce

L'assignation taxonomique des **espèces d'intérêt** est améliorée par l'utilisation d'un jeu de référence spécifique du milieu urinaire (Fig.2):

- **Uropathogènes:**
  - Moins de faux+ et faux- (notamment chez *E.coli* et *K.pneumoniae*)
  - *E.faecium* à améliorer
- **Lactobacilles:**
  - Plus aucun faux- sur ce jeu test