



Stanford Ribonanza RNA Folding

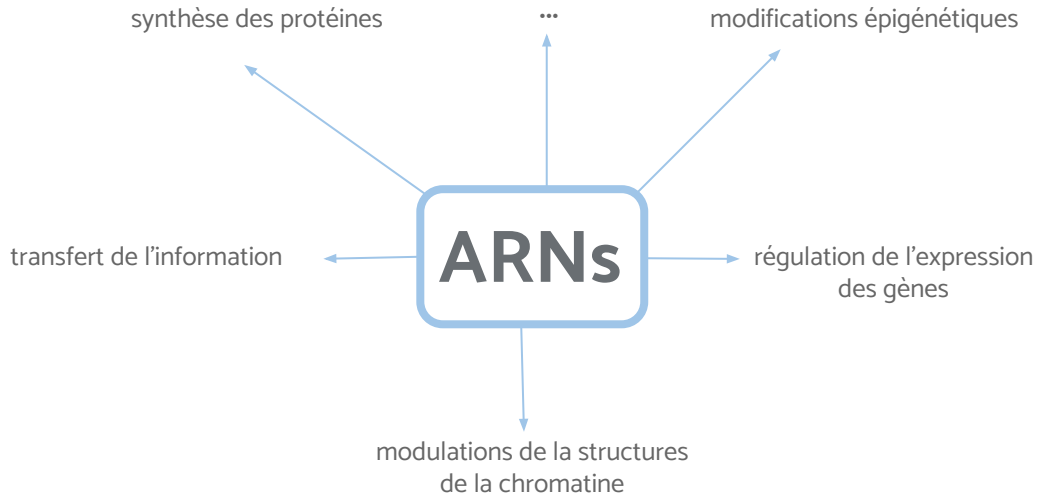
Projet Kaggle : Création d'un modèle de prédiction de la structure de molécules d'ARN

UE Apprentissage, Intelligence artificielle et Optimisation 1 - M2BI

OUADAH Lilia - MERABET Anis - BAILLIF Marine - JEAN-MARIE ROUDE

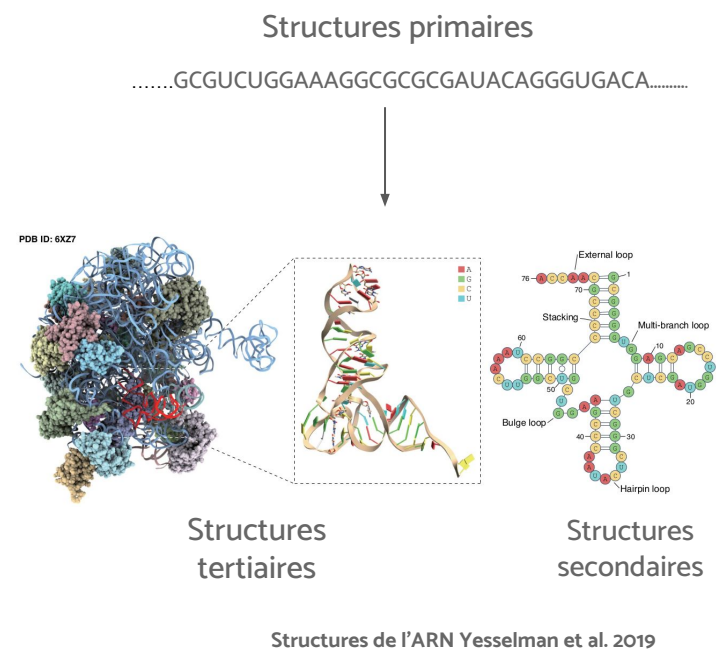
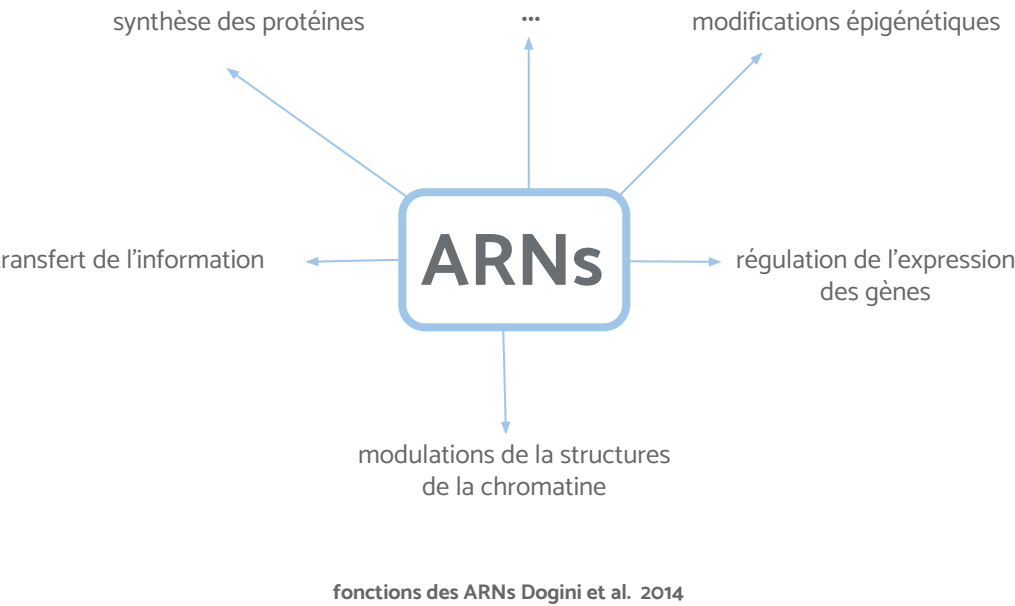


Introduction

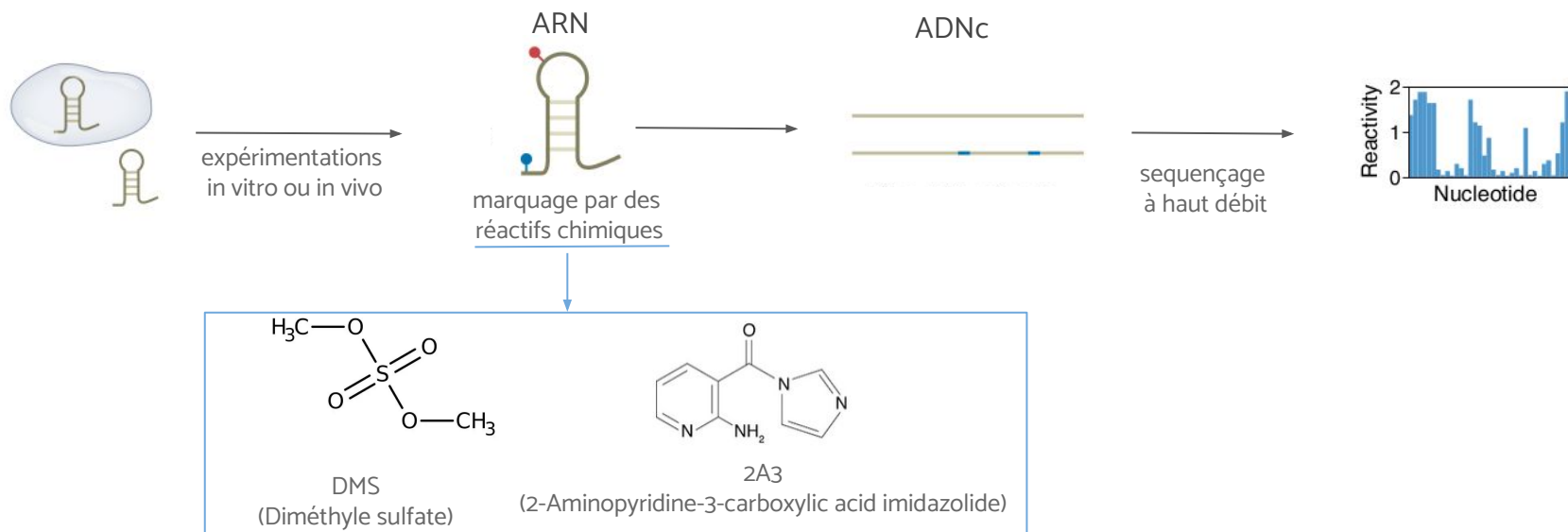


fonctions des ARNs Dogini et al. 2014

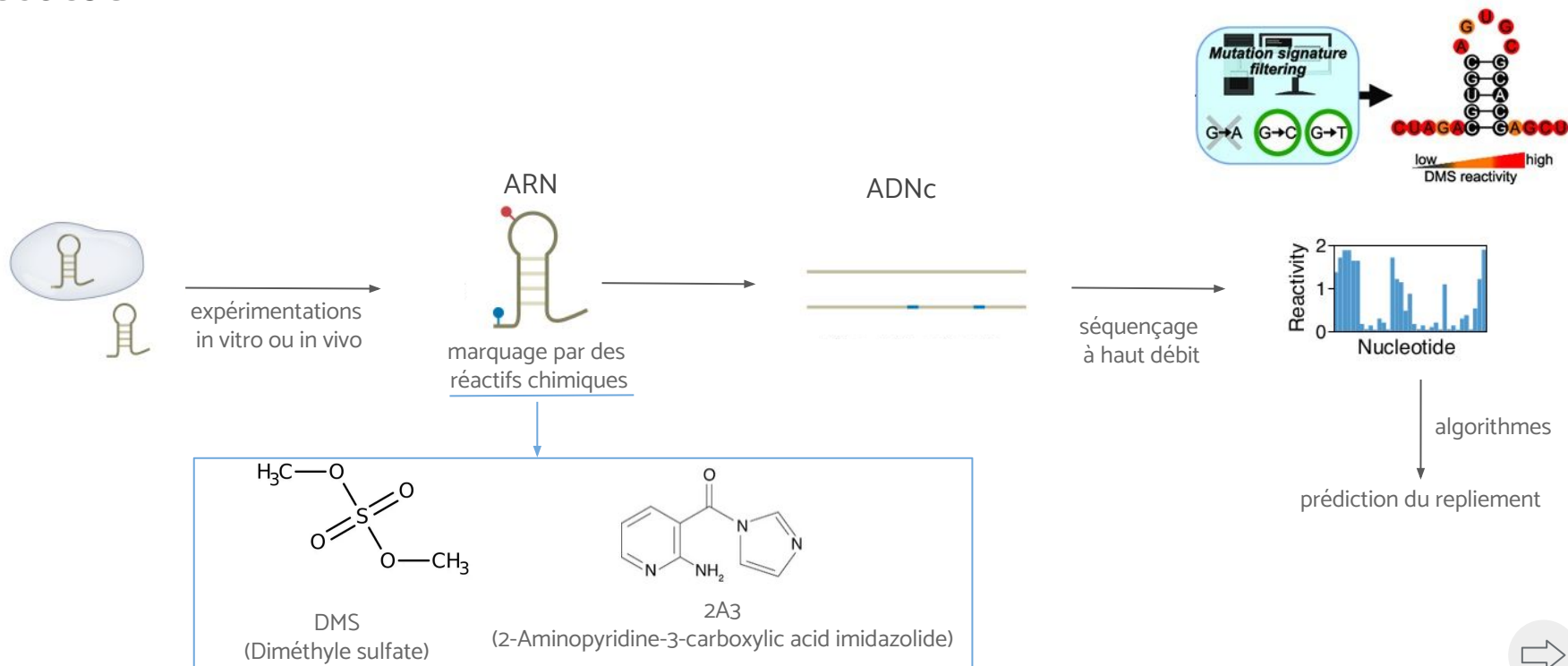
Introduction



Introduction



Introduction



Problématique

Peut on prédire computationnellement le profil de réactivité d'une séquence d'ARN ?

Projet Kaggle

Objectif

Prédire pour une séquence d'ARN la réactivité de chacun de ces nucléotides pour les expériences DMS et 2A3

DMS	0.2	0.2	0.4	0.3
2A3	0.1	0.1	0.1	0.4
	A	U	U	A

Projet Kaggle

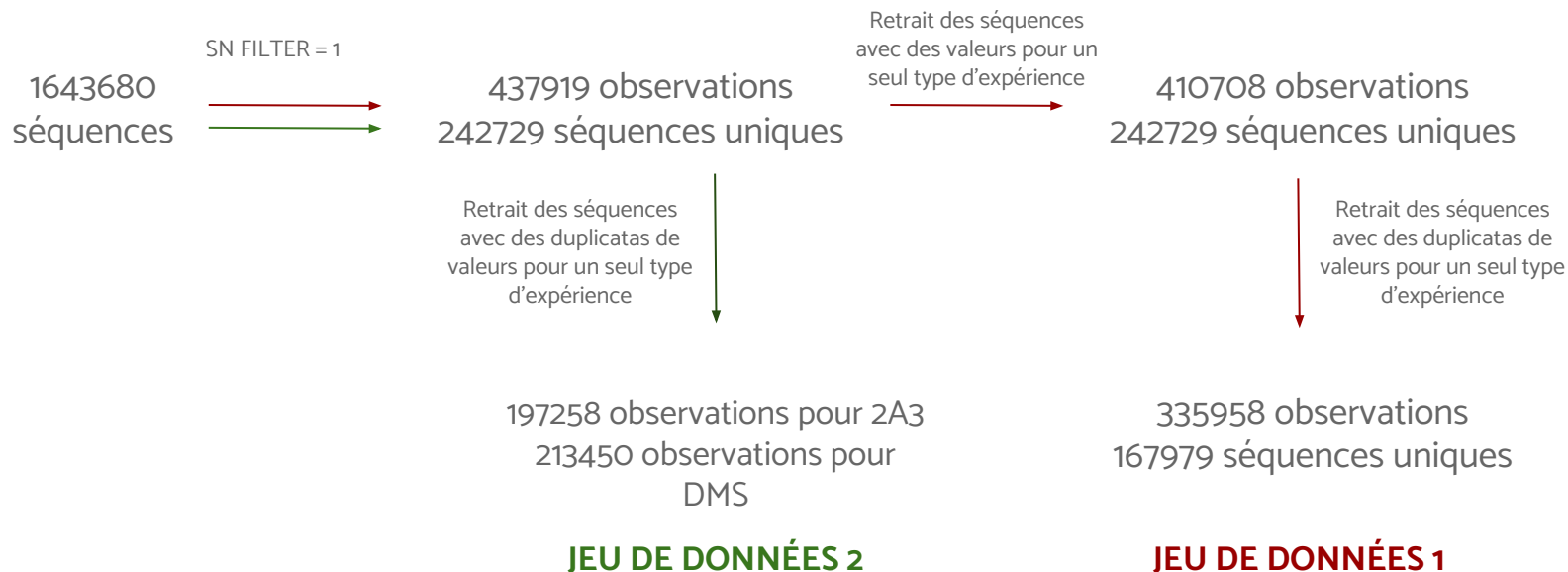
Présentation des données

1643800 observations
806578 séquences uniques

Sequence	Dataset name	Experiment type	Signal to noise	SN Filter	Reactivity 1	...	Reactivity n	Reactivity error 1	...	Reactivity error n
AU...A	...	2A3/DMS	0.02	0/1	0.012	...	0.003	0.02	...	0.001

Preprocessing

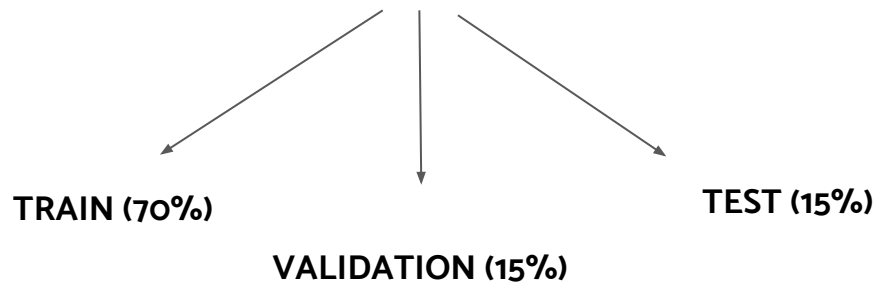
Filtrage



Preprocessing

Filtrage

JEU DE DONNÉES 1 / JEU DE DONNÉES 2



Preprocessing

Encoding

A

1
0
0
0

U

0
1
0
0

C

0
0
1
0

G

0
0
0
1

Preprocessing

Encoding



Preprocessing

Padding

seq 1

A	G	G	C	A	U	G	G	U	A	G	G	C
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 nmax = 457

seq 2

G	G	U	U	C	G	A	A	U	C	C	A
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

 n = x



PADDING →

0	0	0	0	0	0	1	1	0	0	0	1	0	0	0	0
0	0	1	1	0	1	0	0	1	0	0	0	0	0	0	0
0	0	0	0	1	0	0	0	0	1	1	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0

n = 457

Preprocessing

Padding



Preprocessing

Padding

Encoded
sequence

0	0	0	0	0	0	1	1	0	0	0	1	0	0	0	0
0	0	1	1	0	1	0	0	1	0	0	0	0	0	0	0
0	0	0	0	1	0	0	0	0	1	1	0	0	0	0	0
1	1	0	0	0	0	0	0	0	0	0	0	0	0	0	0

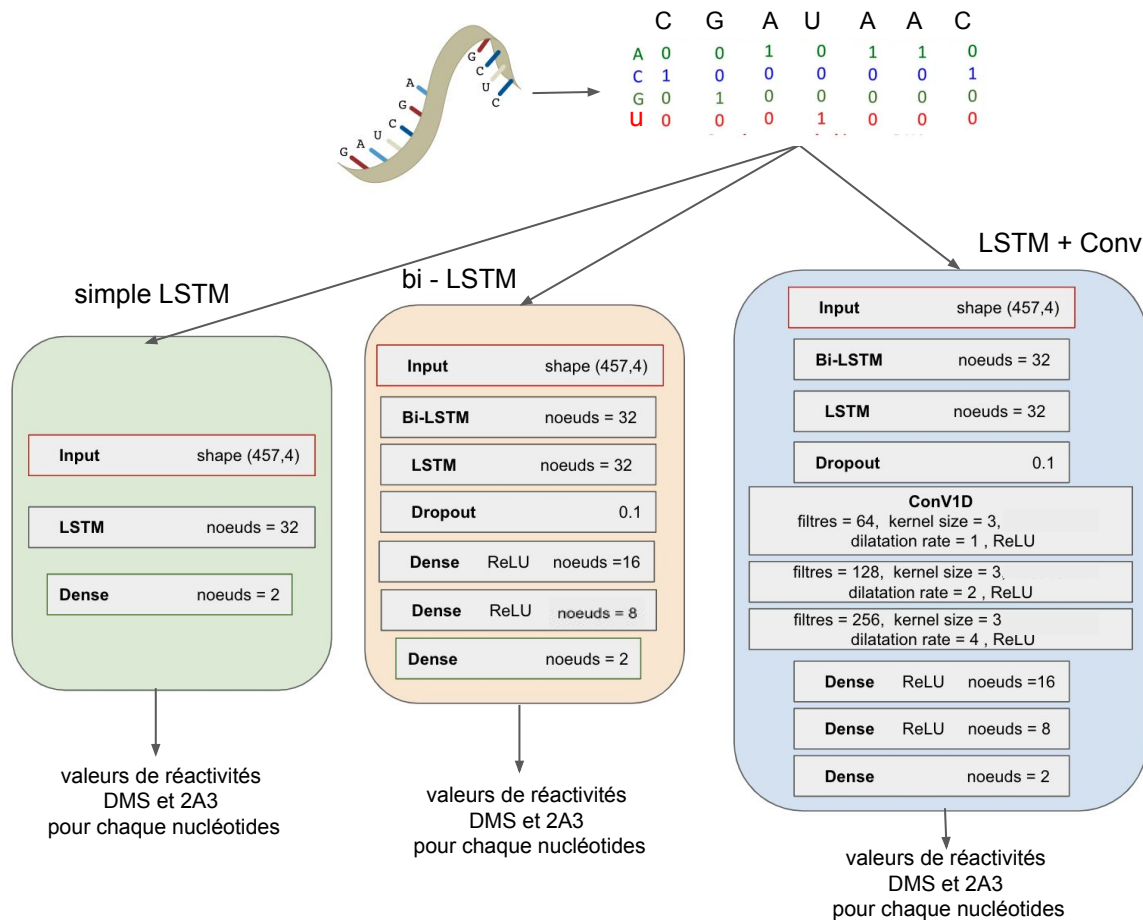
Reactivities

0.25	0.1 1	NA	NA	0.3	0.4 1	0.29	NA	0.31	0.2	0.1	NA	0	0	0	0
------	----------	----	----	-----	----------	------	----	------	---	---	---	---	-----	-----	----	---	---	---	---

Masking

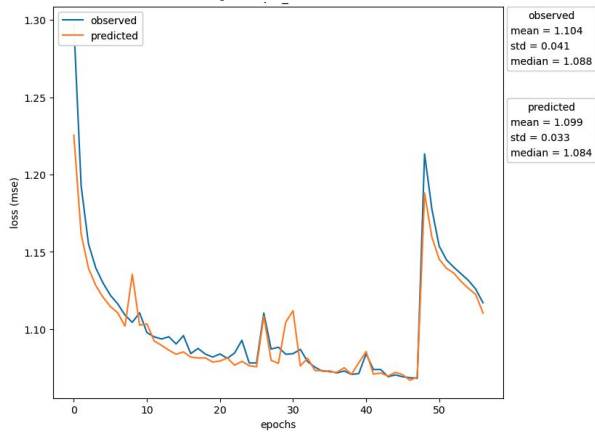
0.25	0.1 1			0.3	0.4 1	0.29		0.31	0.2	0.1					
------	----------	--	--	-----	----------	------	--	------	---	---	---	---	-----	-----	--	--	--	--	--

Modèles simples :

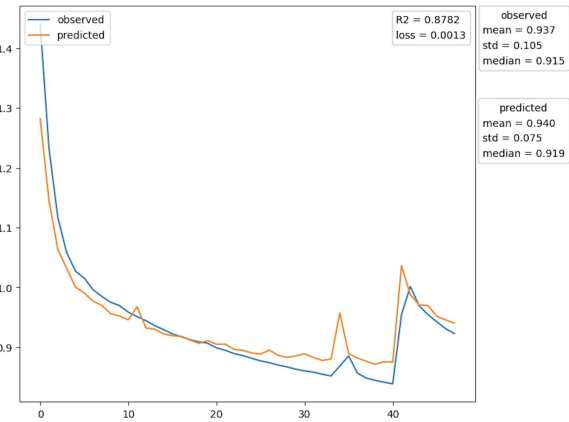


Résultats et Discussions :

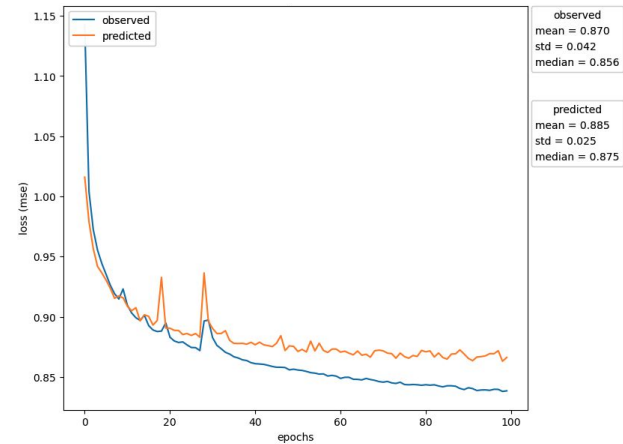
simple LSTM



Bi - LSTM

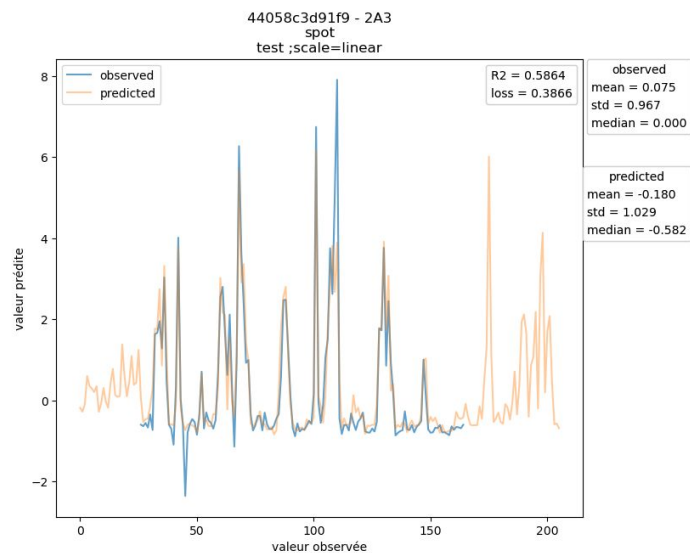
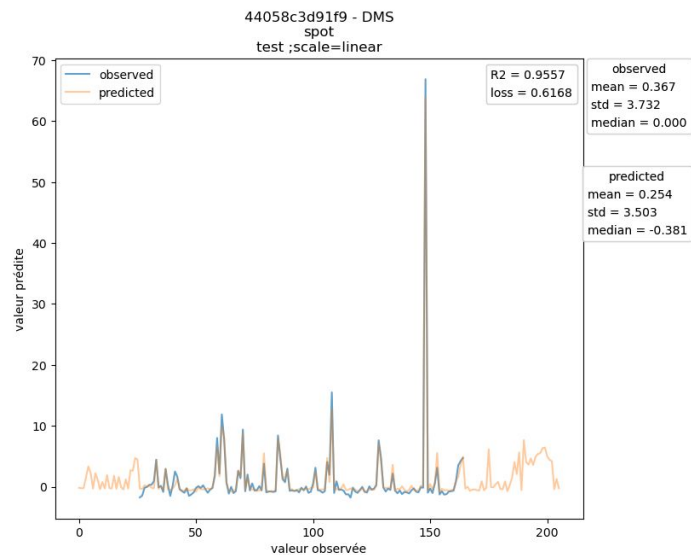


LSTM + ConV



Courbes Loss des modèles simples

Résultats et Discussions :

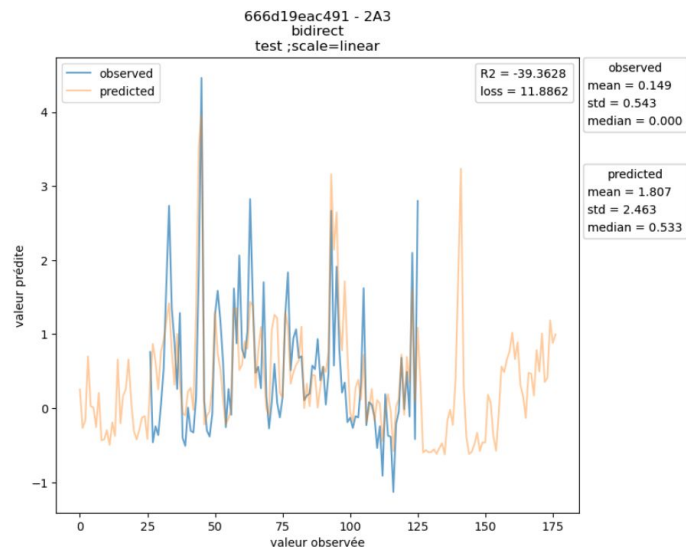


Comparaison des valeurs observées et prédites de la réactivité par le modèle LSTM Convolutionnel pour une séquence (id: 44058c3d91f9)

meilleures prédictions des modèles simples pour le jeu test

Résultats et Discussions :

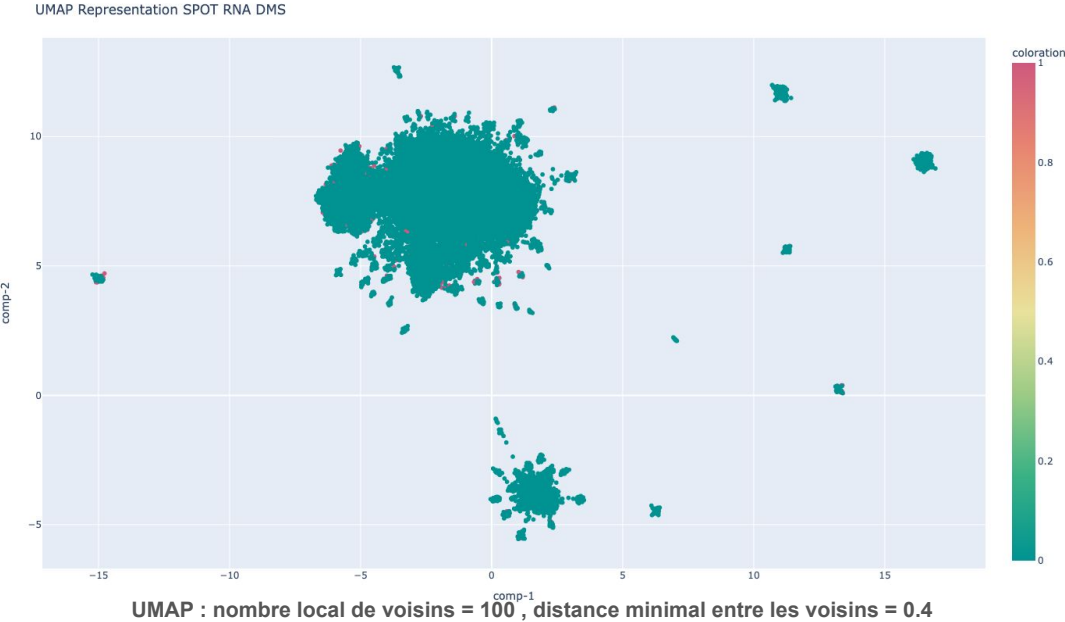
Bi - LSTM



Comparaison des valeurs observées et prédites de la réactivité par le modèle Bi- LSTM pour une séquence (id: 666d19eac491)

pires prédictions des modèles simples pour le jeu test

Prédiction du jeu de données test

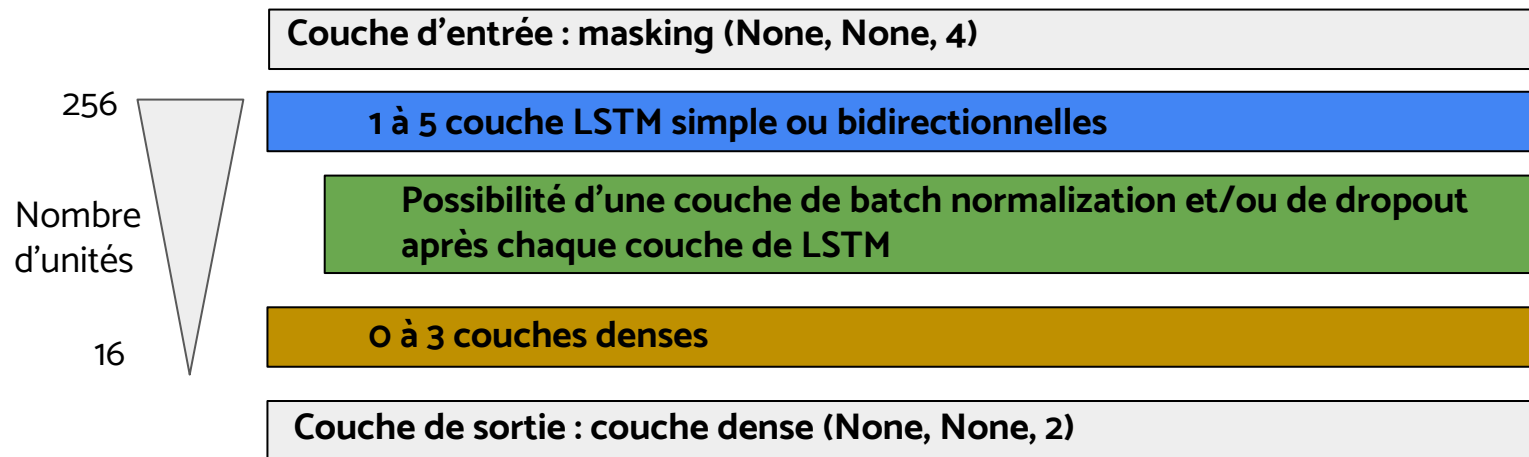


Modèles	$R^2 > 0.5$
Simple LSTM 2A3	3
Simple LSTM DMS	42
Bi -LSTM 2A3	19
Bi -LSTM DMS	610
LSTM + Conv 2A3	385
LSTM + Conv DMS	1312

Nombre de prédictions associé à un $R \geq 0.5$ pour chaque modèles

Optimisation des hyperparamètres du modèle

Approche bayésienne visant réduire le nombre d'essais nécessaires pour identifier les hyperparamètres optimaux.



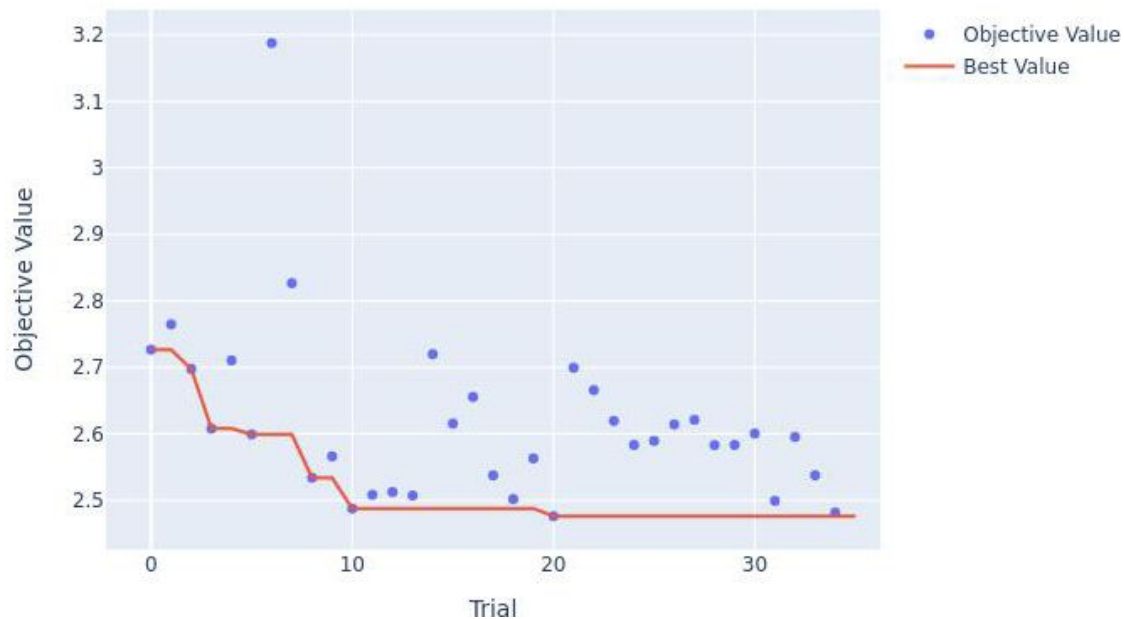
Paramètres fixes : taille des batches : 64. Nombre d'époques : 10. Optimiseur : Adam. fonction d'activation : "relu" pour les couches dense et "linear" pour la couche de sortie.

Optimisation des hyperparamètres du modèle

Un total de 35 itérations.

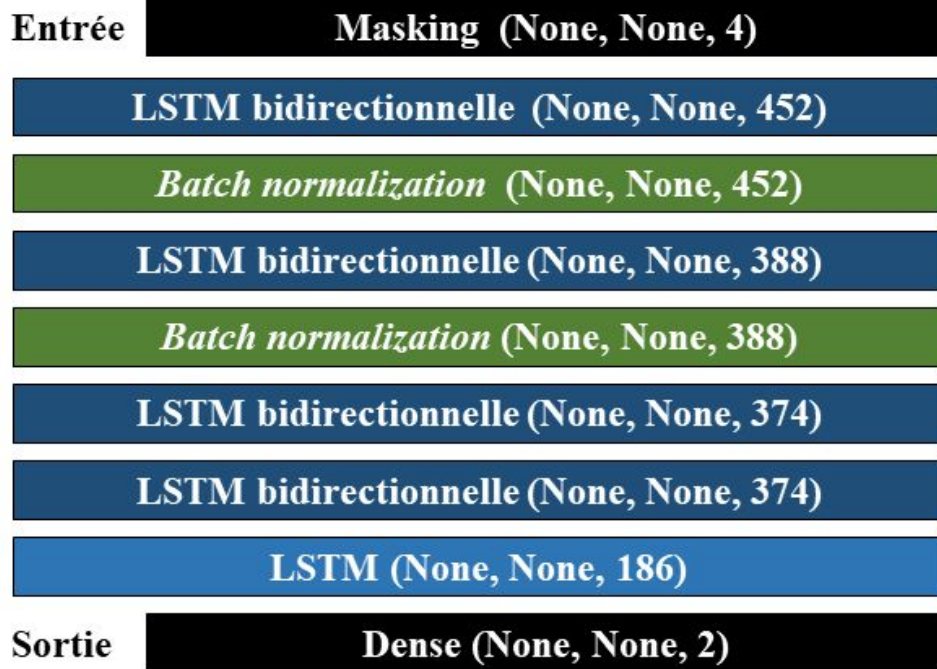
Le meilleur modèle a été enregistré à l'**itération 20**

Optimization History Plot



Optimisation des hyperparamètres du modèle

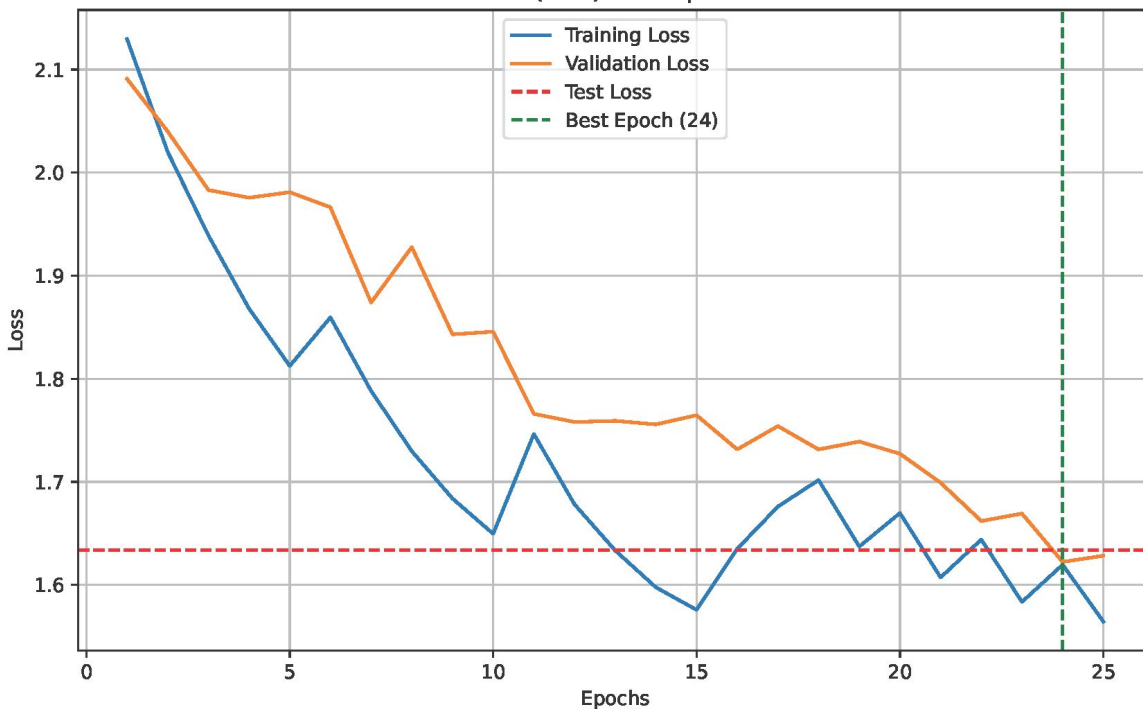
Approche bayésienne visant réduire le nombre d'essais nécessaires pour identifier les hyperparamètres optimaux.



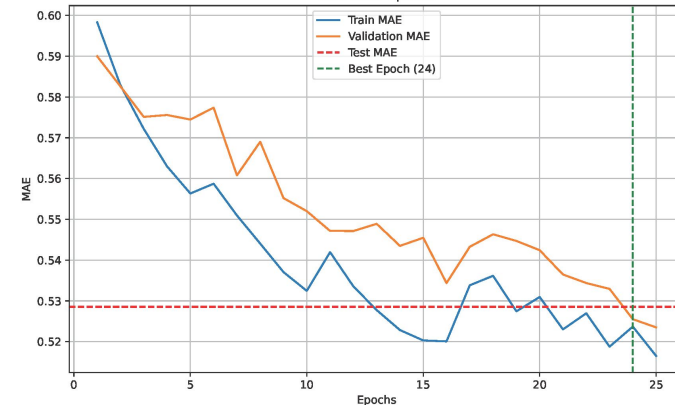
Nombre de paramètres : **3 537 886**

Entrainement du meilleur modèle

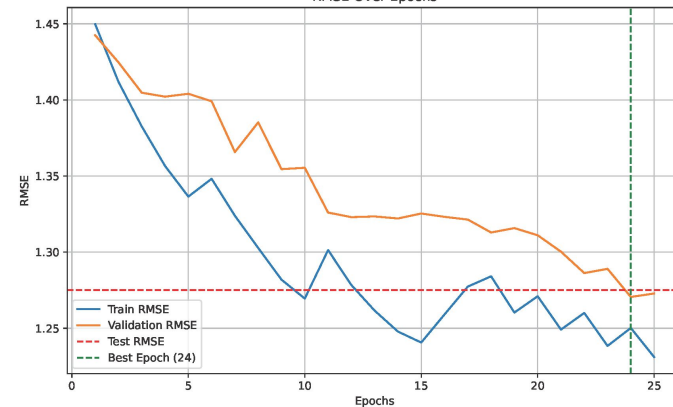
Loss (MSE) Over Epochs



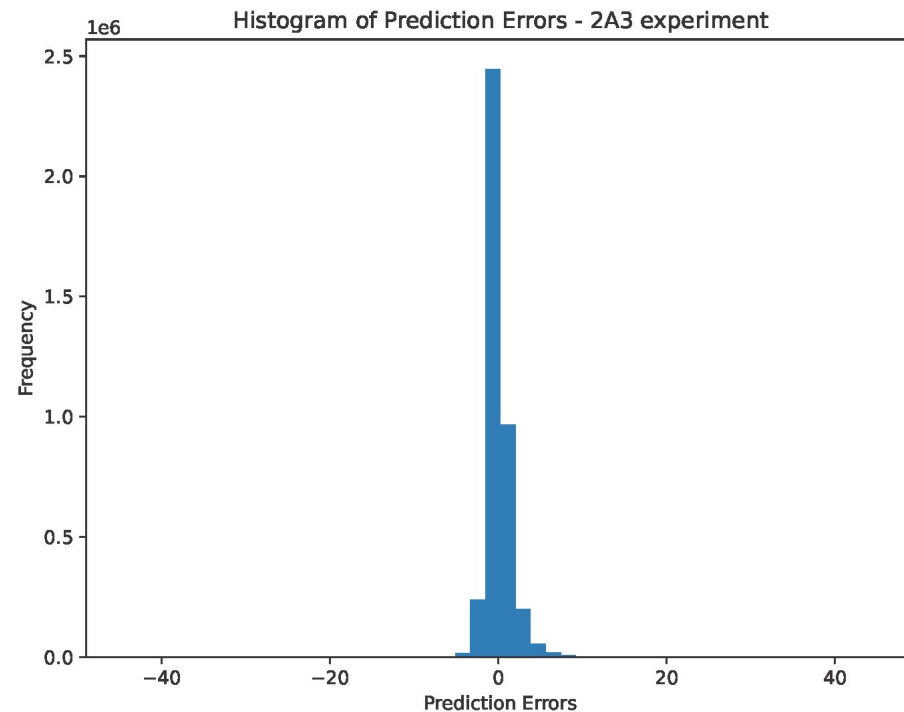
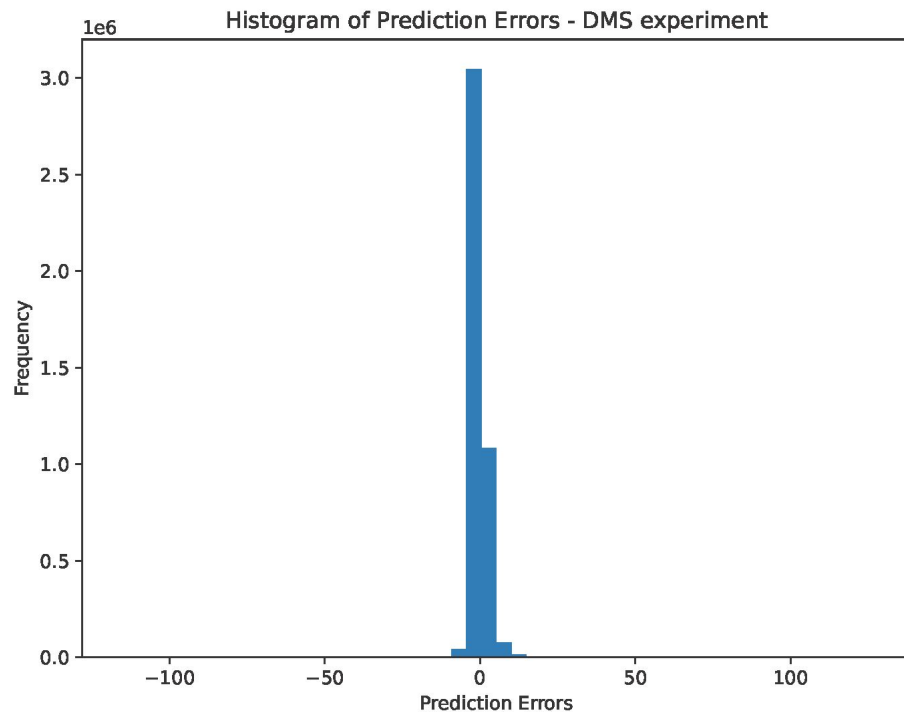
MAE Over Epochs



RMSE Over Epochs

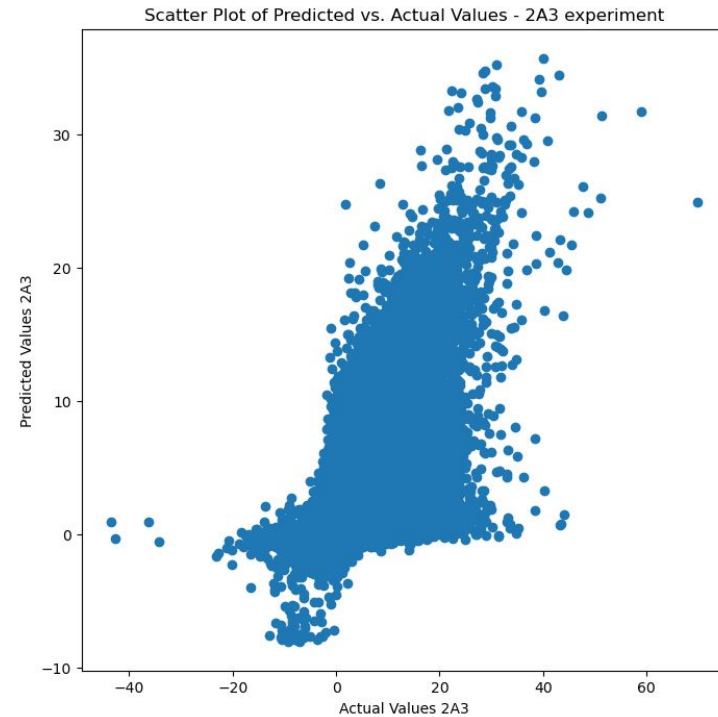
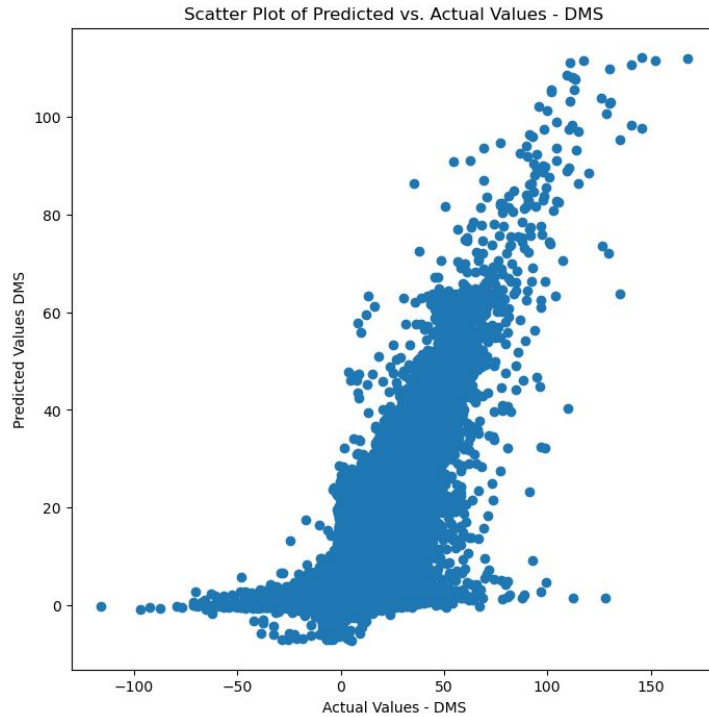


Erreurs de prédiction du meilleur modèle



Les valeurs de réactivité dans l'expérience de **DMS** sont mieux prédites que celles dans l'expérience de **2A3**

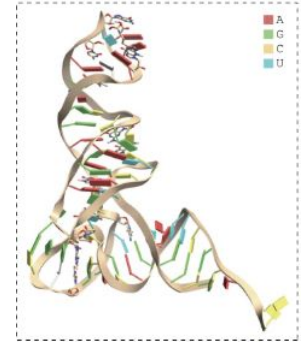
Homogénéité des prédictions des valeurs de réactivité



Le modèle arrive à prédire aussi bien les valeurs proches de zéro que les valeurs extrêmes

Conclusion :

- Performances différentielles entre les différents modèles
- Meilleures prédictions des réactivités pour le DMS
- Bonne généralisation aux données du test pour le modèle obtenu par optimisation des hyperparamètres
- Bonne capacité de prédiction des valeurs extrêmes



Notre modèle optimisé pourrait être un bon modèle de prédiction des structures d'ARN

Perspectives :

- Modèle simples : modèle combinant couches de convolution et couches LSTM à explorer davantage (*cf SPOTRNA1, SPOTRNA2*)
- Utiliser des modèles basés sur les transformers afin de concevoir un embedding qui sera utilisé pour mieux prédire les valeurs de réactivités
- Utiliser des données redondantes en appliquant un poids en fonction de la redondance des séquences
- Appliquer un poids selon la valeur de signal_to_noise

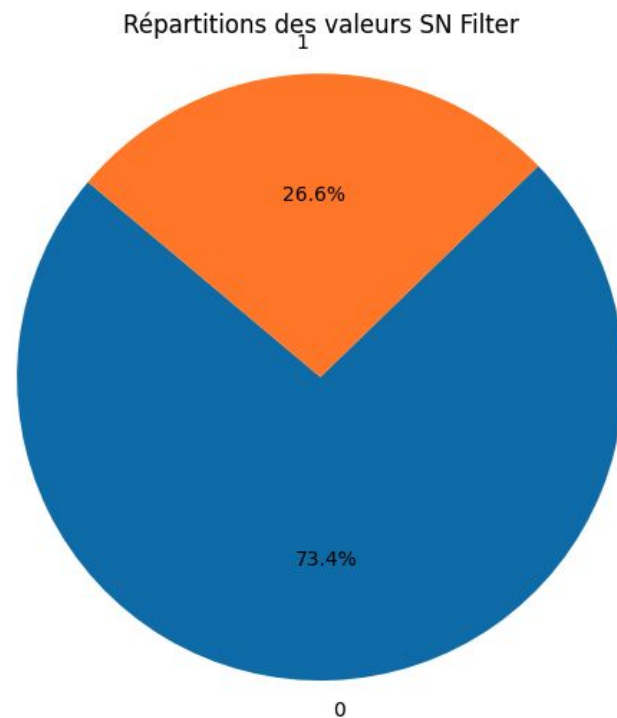
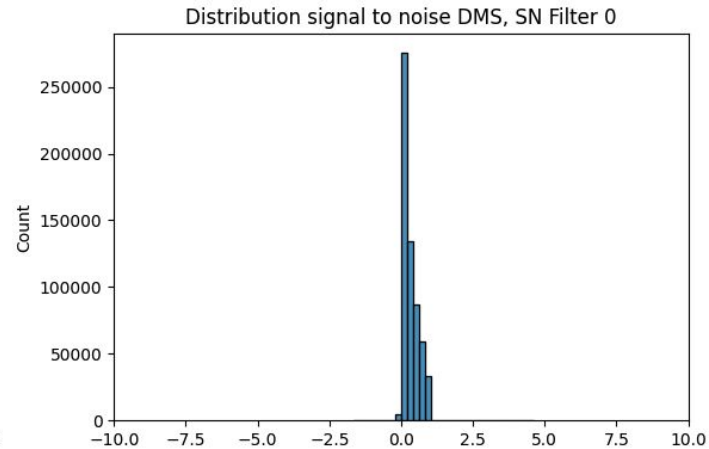
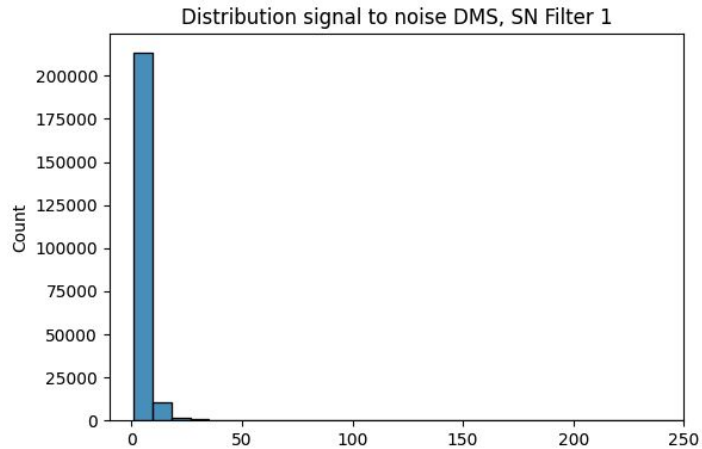
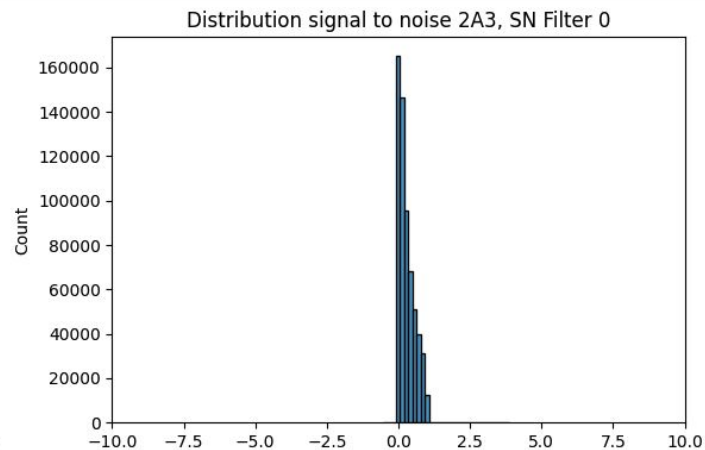
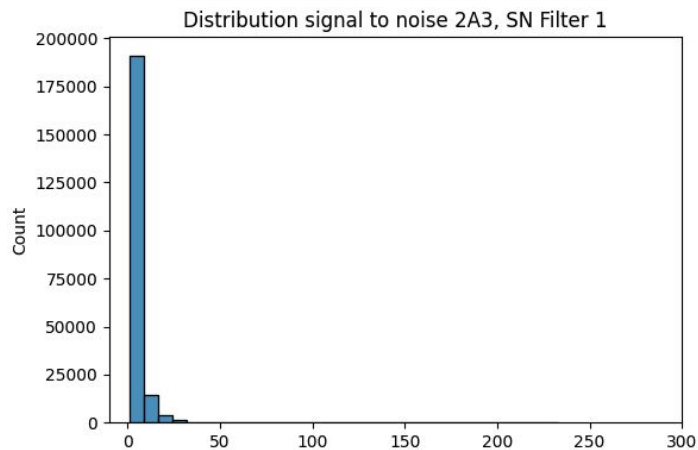
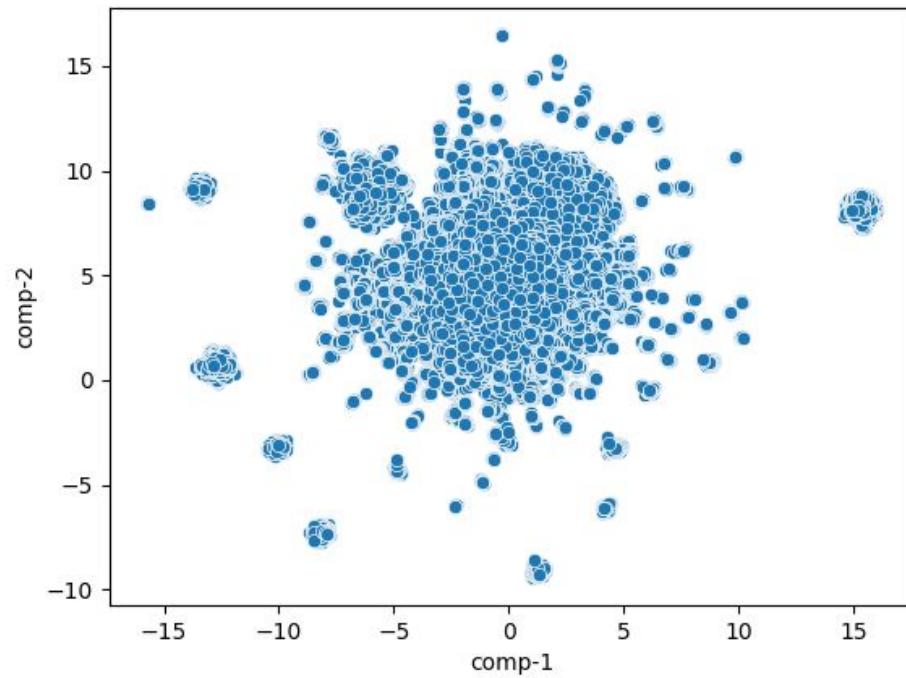
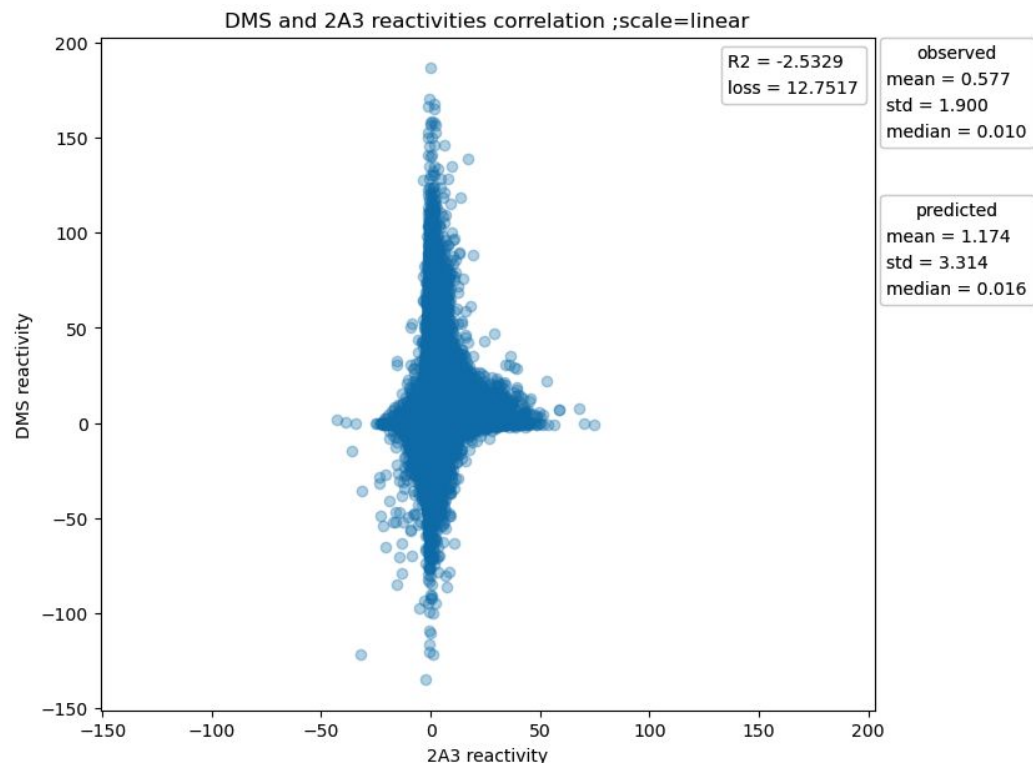
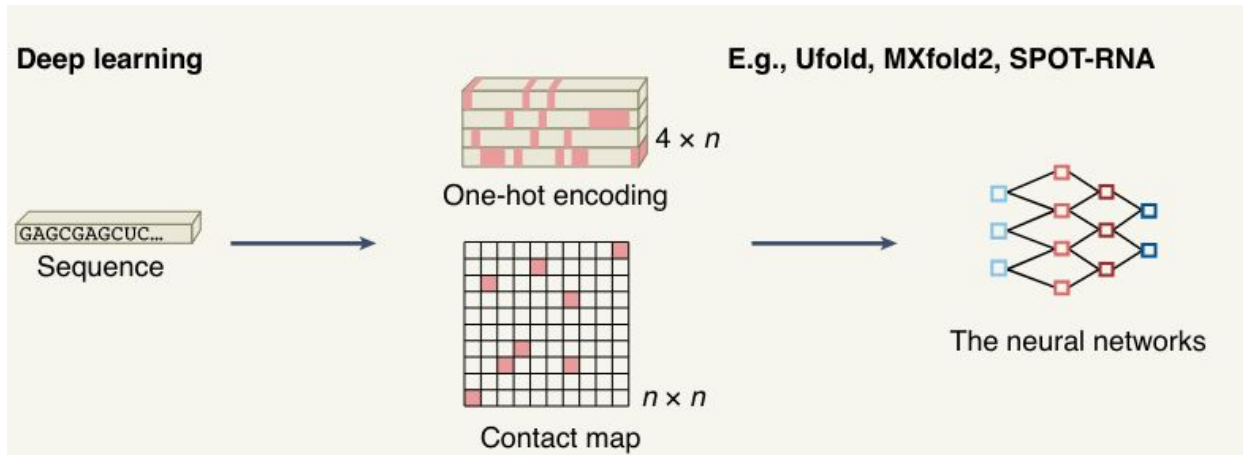


Figure 1 : Répartition des valeurs de SN Filter (0 ou 1) dans le jeu de données.









approches dans la littérature