

# BIOINFODIAG

Réseau Français de Bioinformatique pour le Diagnostic

## API HOUR

---

# Le deep learning en génétique moléculaire: 3 outils qui changent la donne

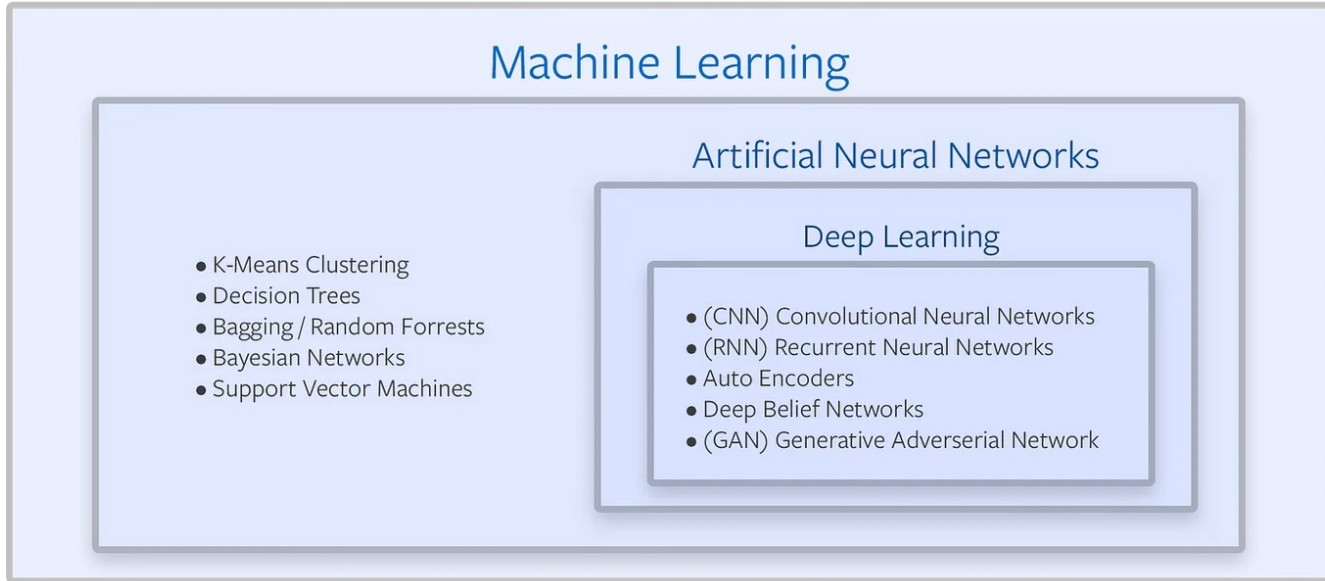
David Baux, CHU de Montpellier

31 mai 2023

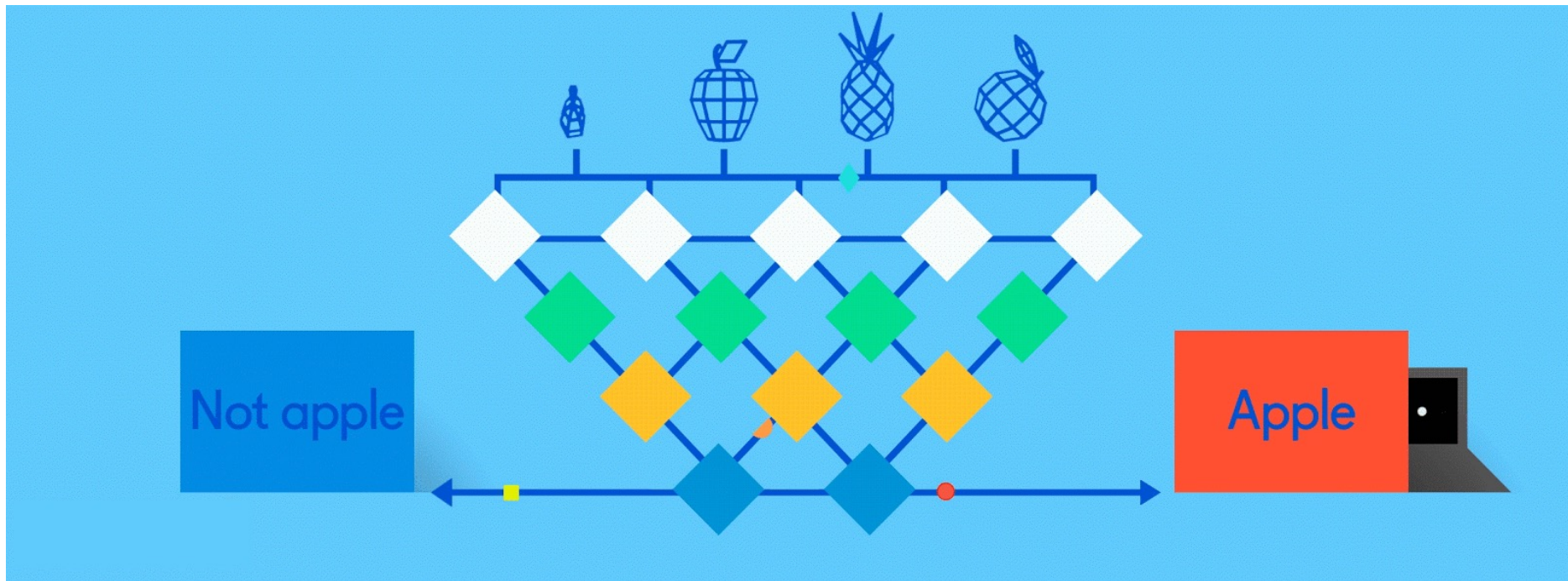


**Disclaimer: je ne suis pas du tout  
expert en IA**

# Artificial Intelligence



from Facebook Developer Circles Resources



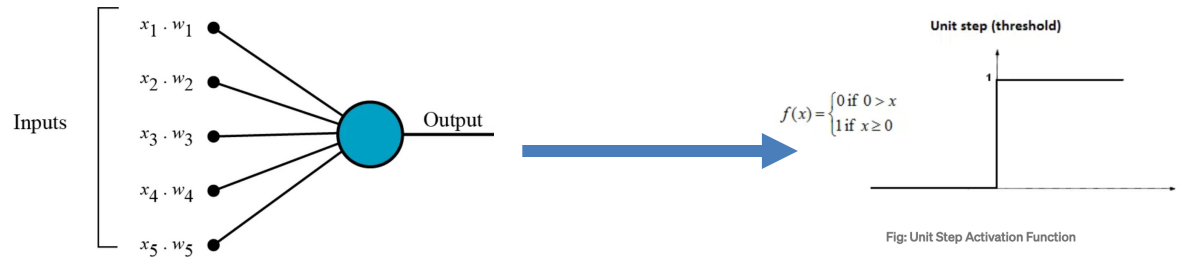
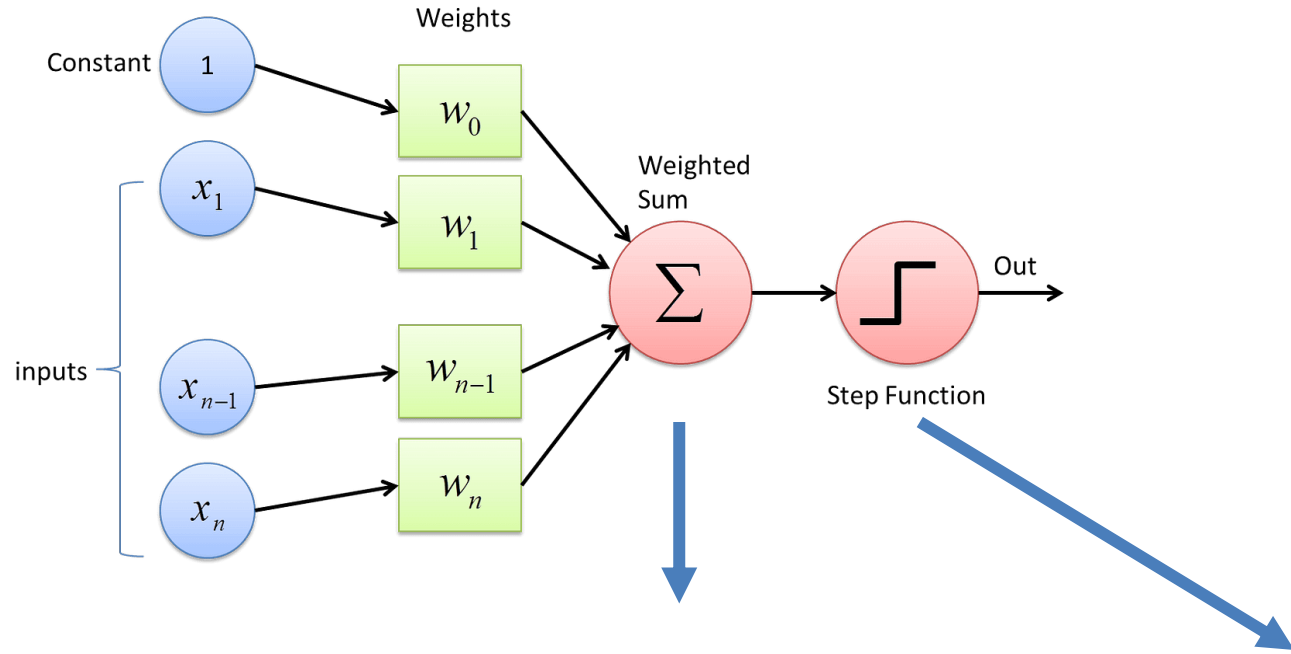
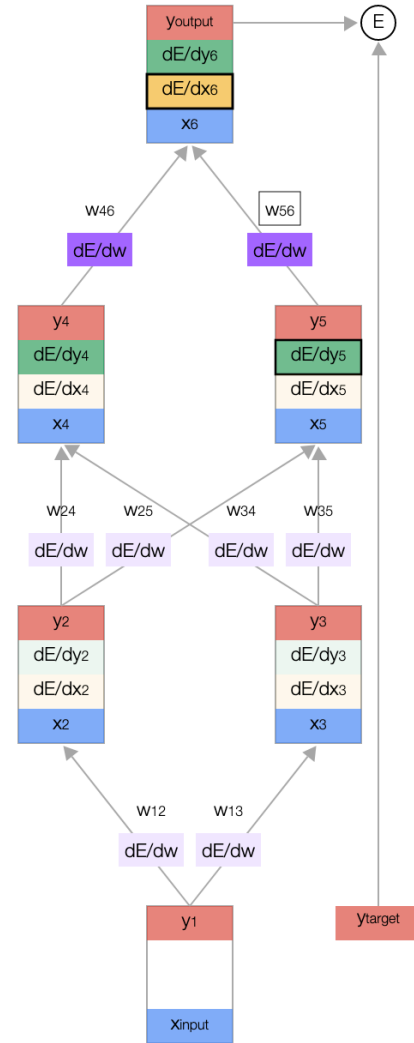
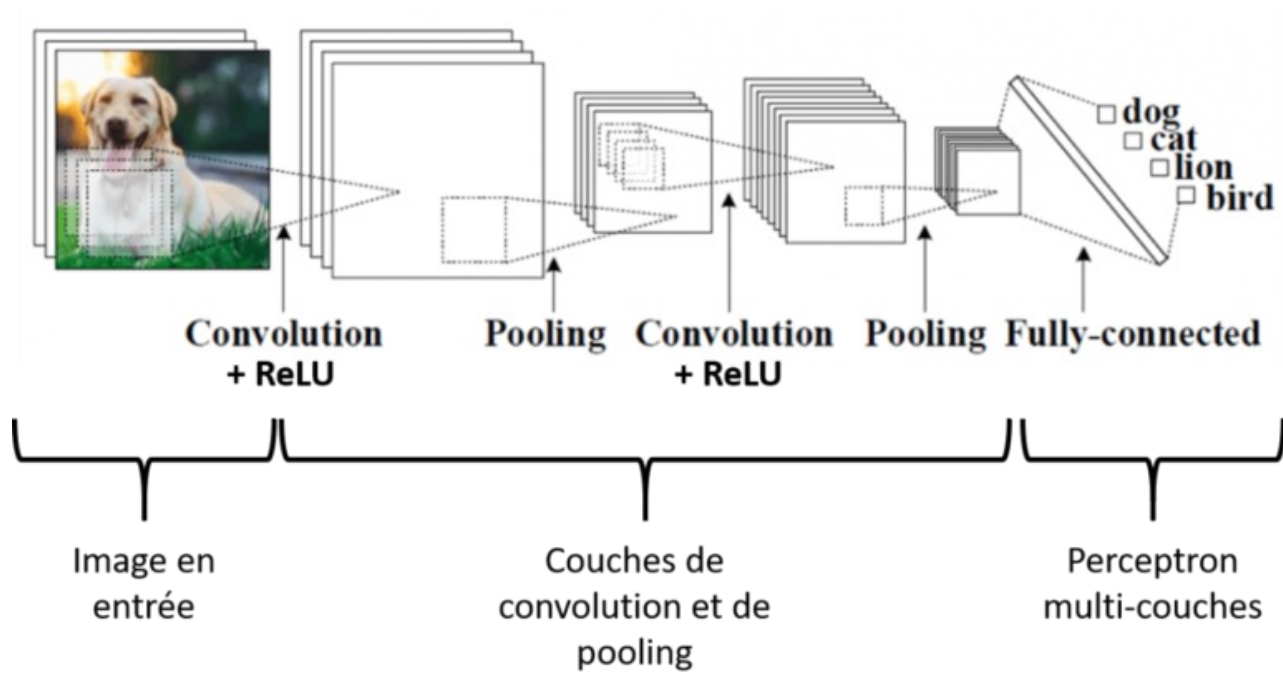


Fig: Multiplying Inputs with weights for 5 inputs

Fig: Unit Step Activation Function

# Entrainement et rétropropagation de l'erreur: Mise à jour des pondérations





# A universal SNP and small-indel variant caller using deep neural networks

Ryan Poplin<sup>1,2</sup>, Pi-Chuan Chang<sup>2</sup>, David Alexander<sup>2</sup>, Scott Schwartz<sup>2</sup>, Thomas Colthurst<sup>2</sup>, Alexander Ku<sup>2</sup>, Dan Newburger<sup>1</sup>, Jojo Dijamco<sup>1</sup>, Nam Nguyen<sup>1</sup>, Pegah T Afshar<sup>1</sup>, Sam S Gross<sup>1</sup>, Lizzie Dorfman<sup>1,2</sup>, Cory Y McLean<sup>1,2</sup> & Mark A DePristo<sup>1,2</sup>

Inputs:

- Alignments (BAM or CRAM)
- Reference (FASTA)



make\_examples



call\_variants



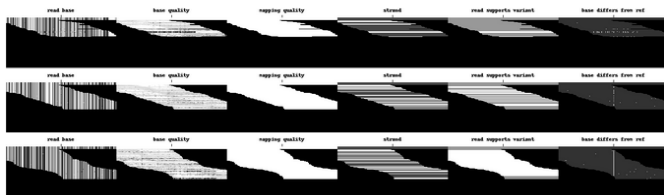
postprocess\_variants

Outputs:

- Variant calls (VCF)
- (optional) gVCF



pileup images



probabilities

[0.99999964, 1.353e-07, 2.223e-07]

[4.4659e-06, 0.99999547, 6.41e-08]

[1.5047e-06, 1.2371e-06, 0.99999726]



Channels are shown in greyscale below in the following order:

1. Read base: different intensities represent A, C, G, and T.
2. Base quality: set by the sequencing machine. White is higher quality.
3. Mapping quality: set by the aligner. White is higher quality.
4. Strand of alignment: Black is forward; white is reverse.
5. Read supports variant: White means the read supports the given alternate allele, grey means it does not.
6. Base differs from ref: White means the base is different from the reference, dark grey means the base matches the reference.



## Deepvariant 1.5 VS NA24385

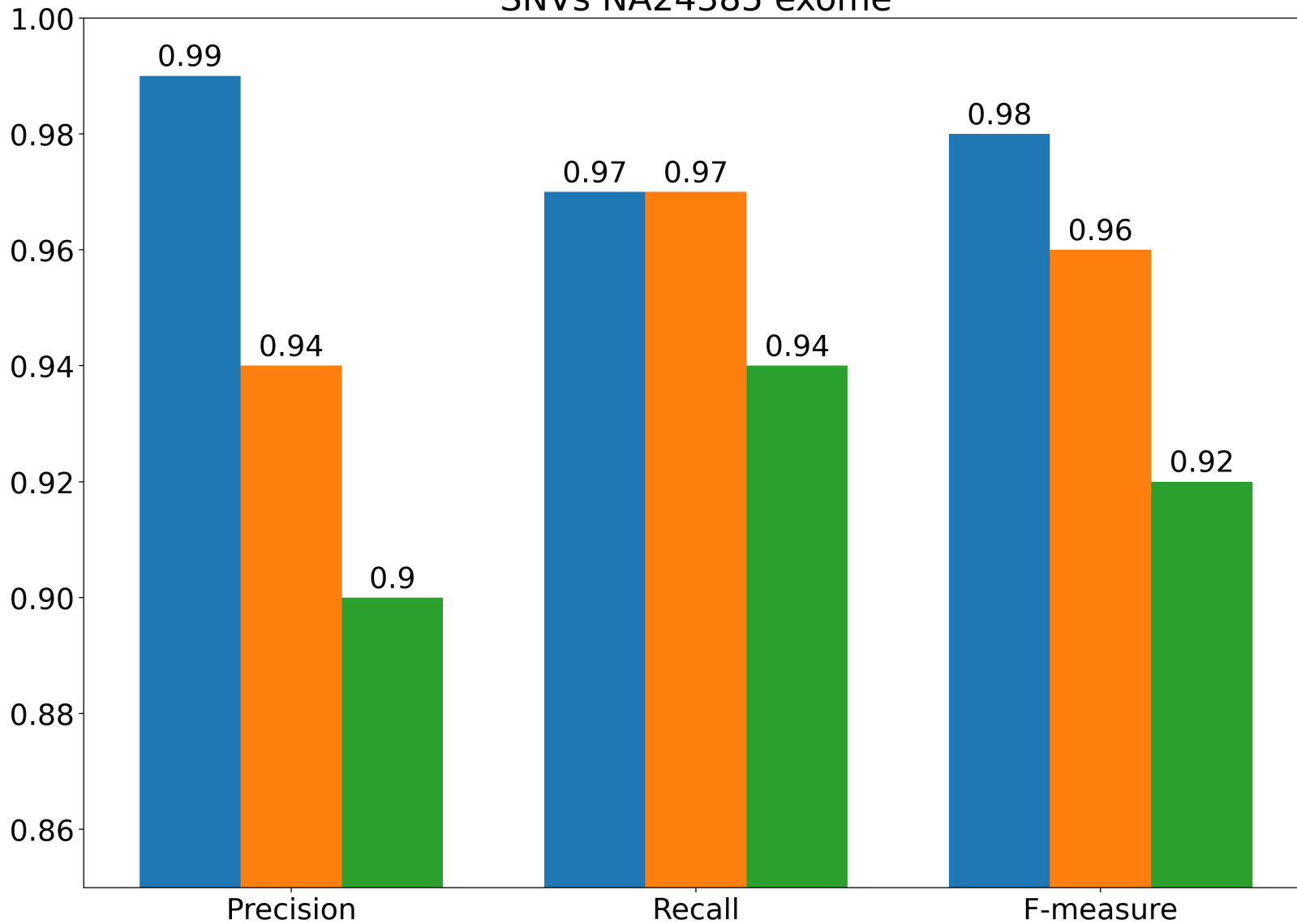
Type	Filter	TRUTH.TOTAL	TRUTH.TP	TRUTH.FN	QUERY.FP	METRIC.Recall	METRIC.Precision	METRIC.F1_Score
INDEL	ALL	1671	1405	266	76	0.840814	0.948959	0.891619
INDEL	PASS	1671	1405	266	76	0.840814	0.948959	0.891619
SNP	ALL	27362	26606	756	249	0.972370	0.990729	0.981464
SNP	PASS	27362	26606	756	249	0.972370	0.990729	0.981464

~186000 régions

38,8Mb

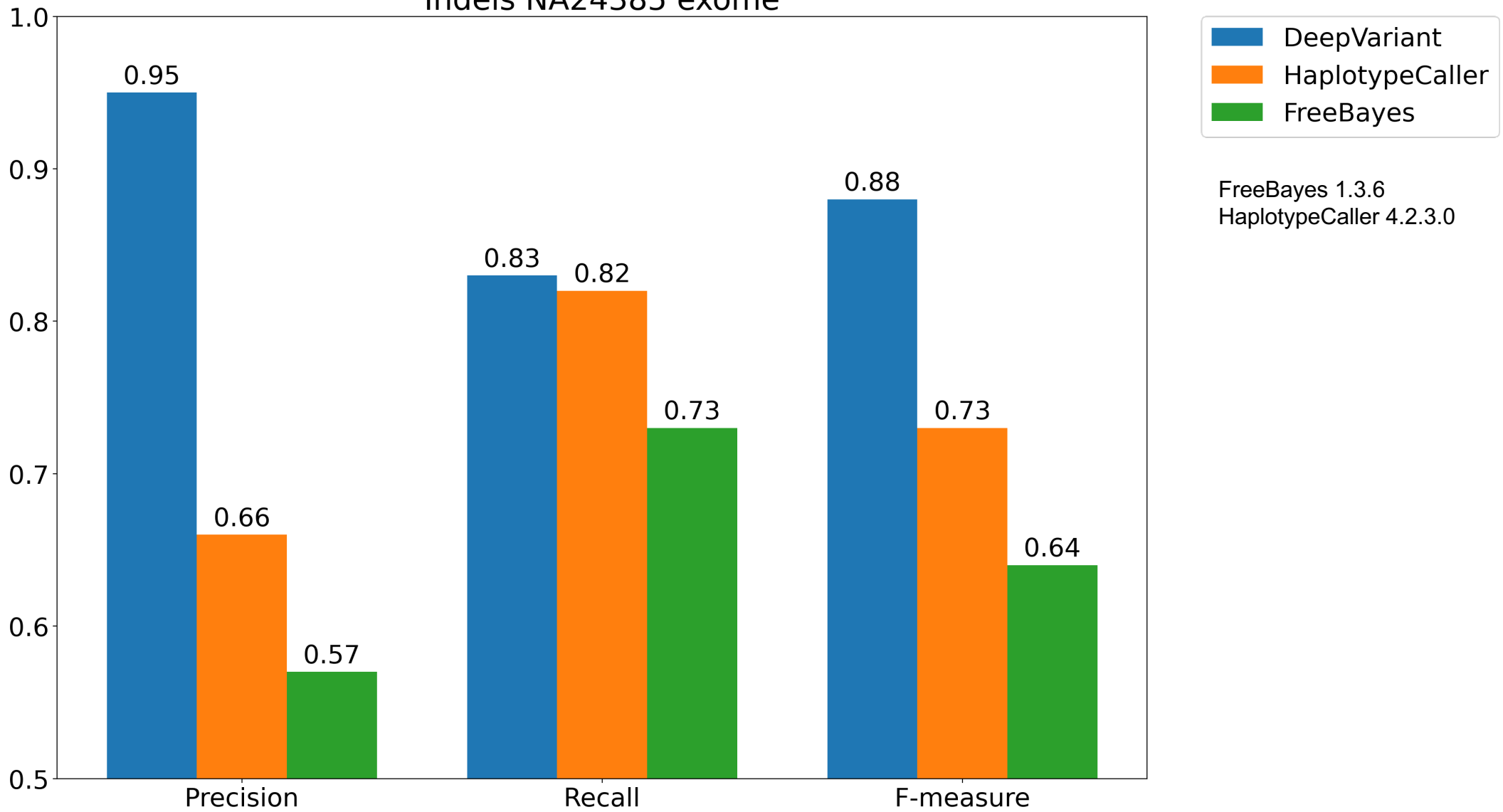
WES 60X

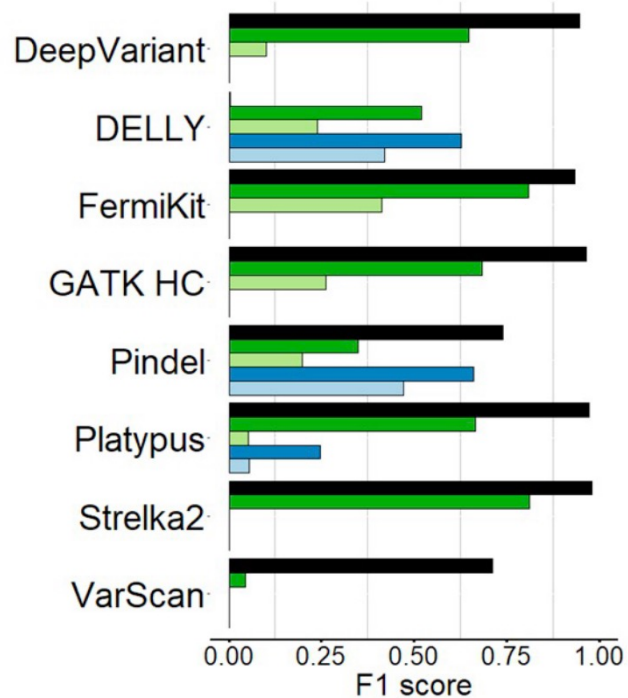
# SNVs NA24385 exome



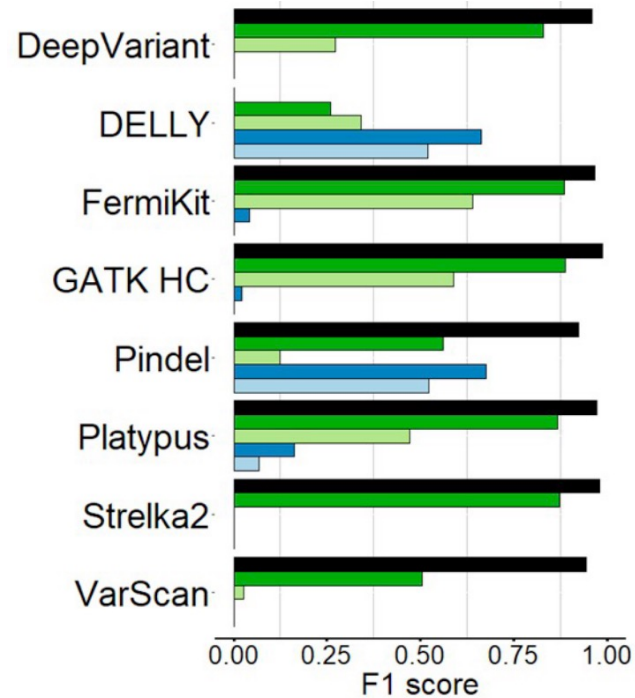
FreeBayes 1.3.6  
HaplotypeCaller 4.2.3.0

# Indels NA24385 exome





30X; 100bp



30X; 250bp

Deletions size ■ 1 - 20 bp ■ 20 - 50 bp ■ 50 - 200 bp ■ 200 - 500 bp ■ ≥ 500 bp

Semi-simulated dataset: 43,066 insertions and 45,223 deletions

## Article

# Highly accurate protein structure prediction with AlphaFold

<https://doi.org/10.1038/s41586-021-03819-2>

Received: 11 May 2021



Accepted: 12 July 2021

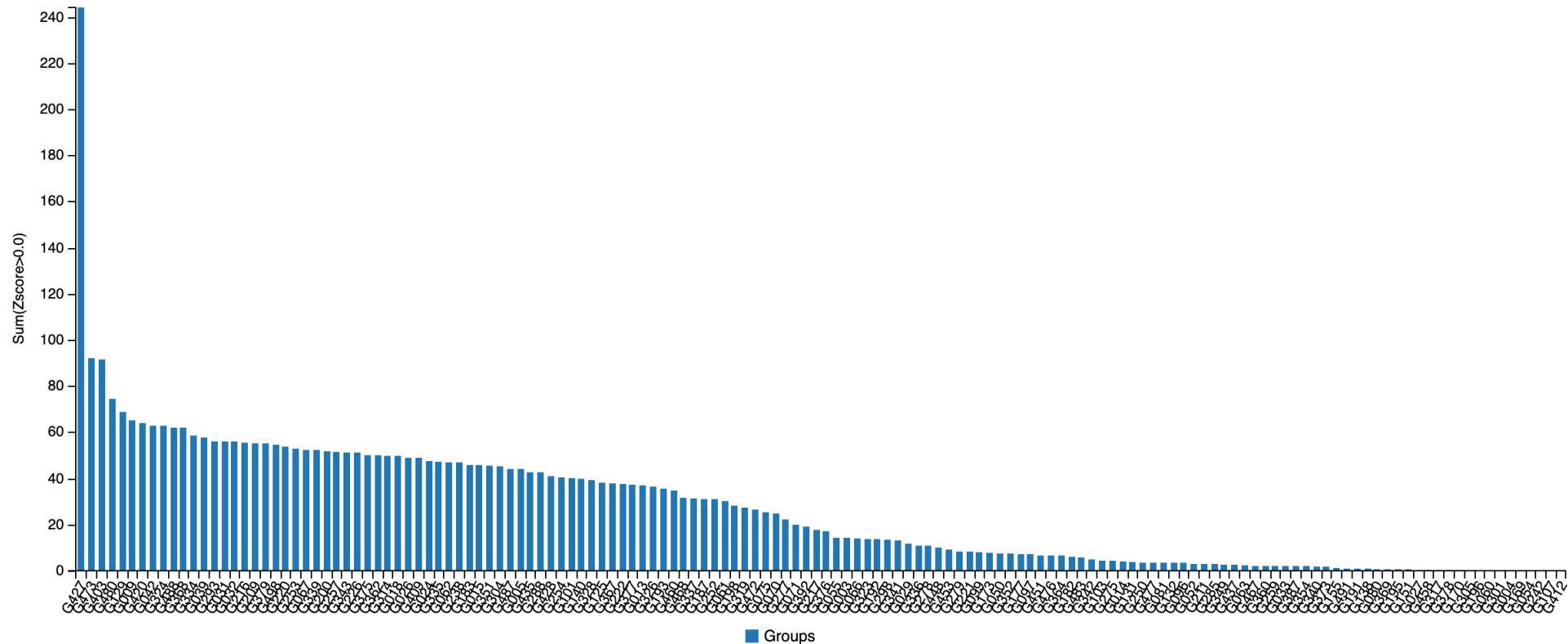
Published online: 15 July 2021

Open access



Check for updates

John Jumper<sup>1,4</sup>, Richard Evans<sup>1,4</sup>, Alexander Pritzel<sup>1,4</sup>, Tim Green<sup>1,4</sup>, Michael Figurnov<sup>1,4</sup>, Olaf Ronneberger<sup>1,4</sup>, Kathryn Tunyasuvunakool<sup>1,4</sup>, Russ Bates<sup>1,4</sup>, Augustin Židek<sup>1,4</sup>, Anna Potapenko<sup>1,4</sup>, Alex Bridgland<sup>1,4</sup>, Clemens Meyer<sup>1,4</sup>, Simon A. A. Kohl<sup>1,4</sup>, Andrew J. Ballard<sup>1,4</sup>, Andrew Cowie<sup>1,4</sup>, Bernardino Romera-Paredes<sup>1,4</sup>, Stanislav Nikolov<sup>1,4</sup>, Rishub Jain<sup>1,4</sup>, Jonas Adler<sup>1</sup>, Trevor Back<sup>1</sup>, Stig Petersen<sup>1</sup>, David Reiman<sup>1</sup>, Ellen Clancy<sup>1</sup>, Michal Zielinski<sup>1</sup>, Martin Steinegger<sup>2,3</sup>, Michalina Pacholska<sup>1</sup>, Tamas Berghammer<sup>1</sup>, Sebastian Bodenstern<sup>1</sup>, David Silver<sup>1</sup>, Oriol Vinyals<sup>1</sup>, Andrew W. Senior<sup>1</sup>, Koray Kavukcuoglu<sup>1</sup>, Pushmeet Kohli<sup>1</sup> & Demis Hassabis<sup>1,4</sup>



**14th Community Wide Experiment on the  
Critical Assessment of Techniques for Protein Structure Prediction**

**TS Analysis : Group performance based on combined z-scores**

# AlphaFold Protein Structure Database

Developed by DeepMind and EMBL-EBI

Search for protein, gene, UniProt accession or organism

BETA

Search

Examples: [Free fatty acid receptor 2](#) [At1g58602](#) [Q5VSL9](#) [E. coli](#) Help: [AlphaFold DB search help](#)

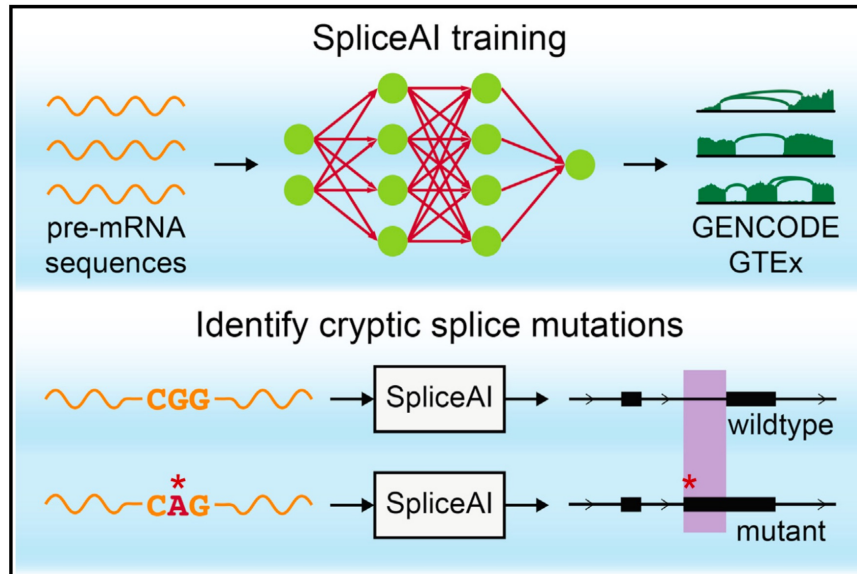
Feedback on structure: [Contact DeepMind](#)

AlphaFold DB provides open access to over 200 million protein structure predictions to accelerate scientific research.



# Predicting Splicing from Primary Sequence with Deep Learning

## Graphical Abstract



## Authors

Kishore Jaganathan,  
Sofia Kyriazopoulou Panagiotopoulou,  
Jeremy F. McRae, ..., Serafim Batzoglou,  
Stephan J. Sanders, Kyle Kai-How Farh

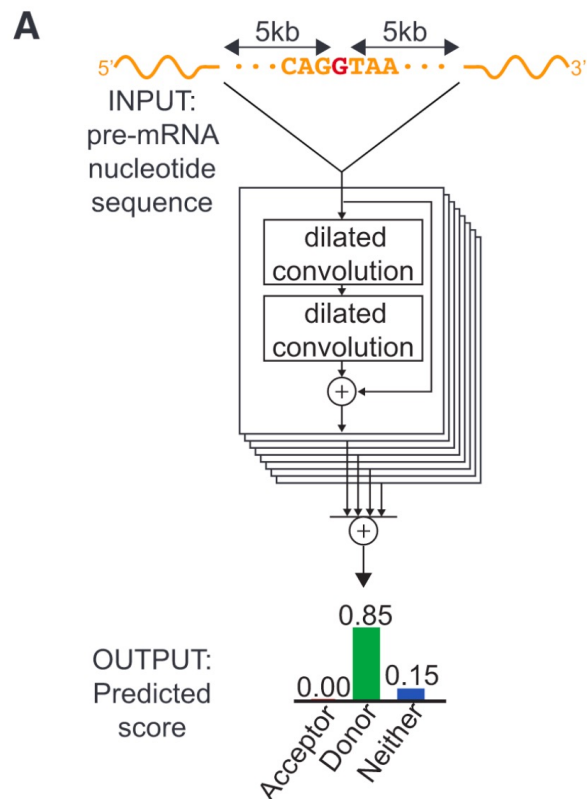
## Correspondence

kfarh@illumina.com

## In Brief

A deep neural network precisely models mRNA splicing from a genomic sequence and accurately predicts noncoding cryptic splice mutations in patients with rare genetic diseases.

# Operation mode - performance



**E**

	Top-k accuracy	PR-AUC
SpliceAI-80nt	0.57	0.60
SpliceAI-400nt	0.90	0.95
SpliceAI-2k	0.93	0.97
SpliceAI-10k	<b>0.95</b>	<b>0.98</b>
GeneSplicer	0.30	0.23
MaxEntScan	0.22	0.15
NNSplice	0.22	0.15

Top-k accuracy is the fraction of correctly predicted splice sites at the threshold where the number of predicted sites is equal to the actual number of sites present.

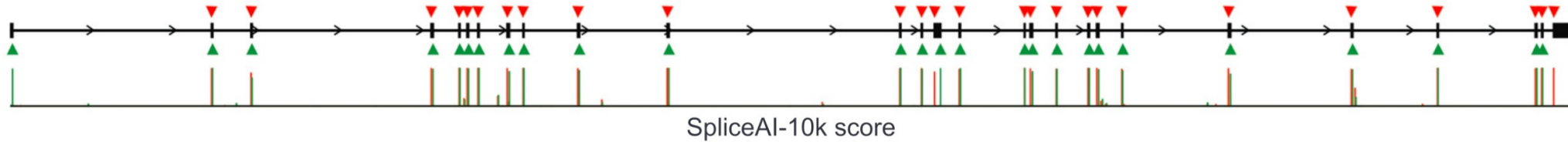
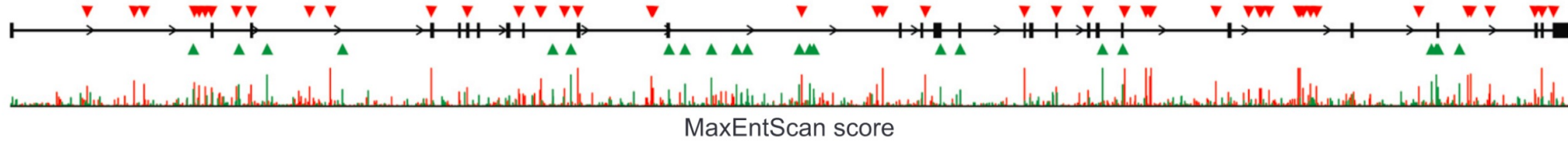
# SpliceAI: splice sites recognition

B

chr7:117,120,017-117,308,719 188,703 nt

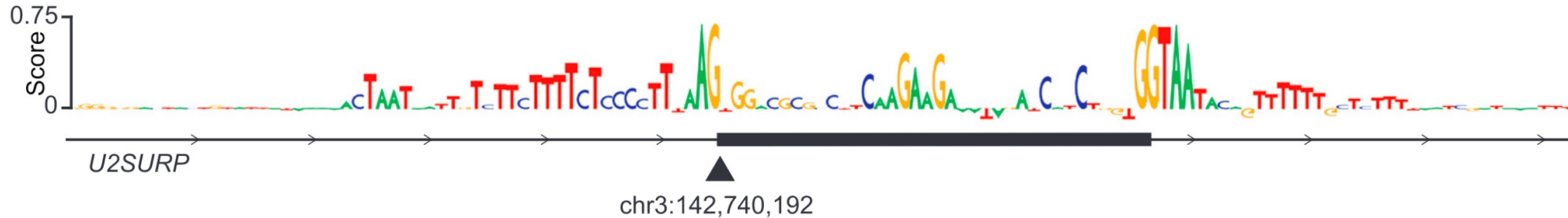
■ Splice acceptor ■ Splice donor

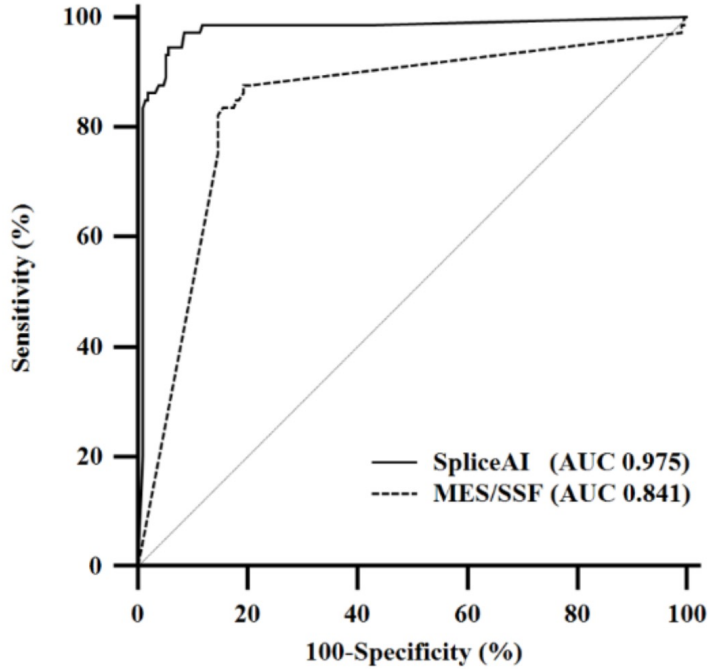
*CFTR*



D

chr3:142,740,137-142,740,263 127 nt



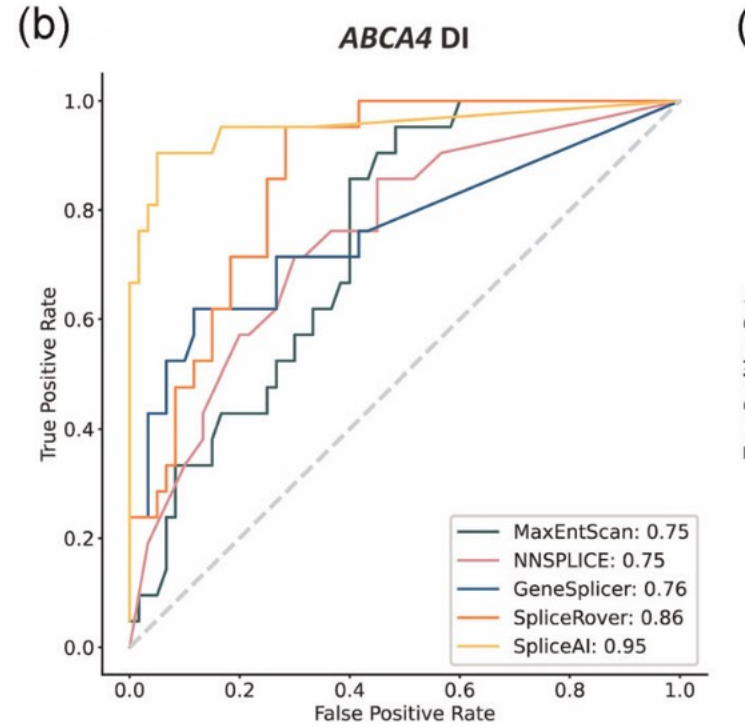
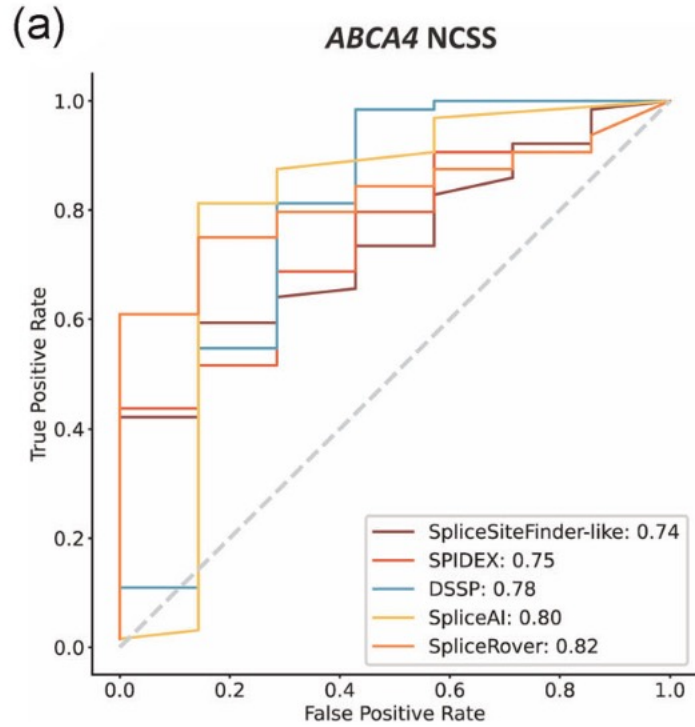


Method	Sensitivity	Specificity
	N/Total N% (95% CI)	N/Total N% (95% CI)
SpliceAI	69/73 94.5% (86.6–98.5%)	200/212 94.3% (90.3–97.0%)

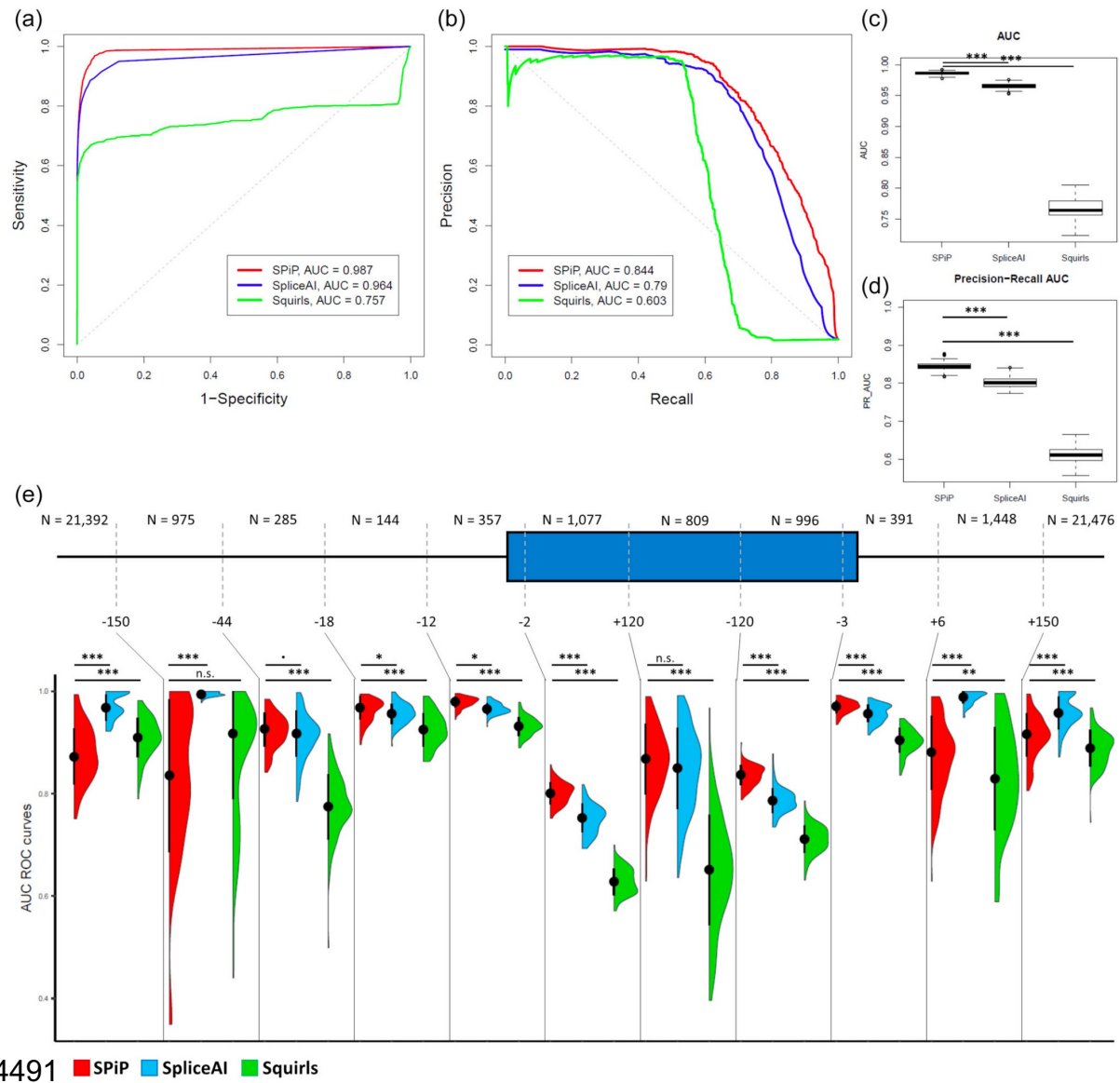
Performance Evaluation of SpliceAI for the Prediction of Splicing of *NF1* Variants (cut-off: 0.22, 285 variants, 73 avec effet, 212 sans effet, séquençage du cDNA)

Ha et al., 2021

<https://doi.org/10.3390/genes12091308>



Tool	<i>ABCA4</i> NCSS	<i>ABCA4</i> DI	<i>MYBPC3</i> NCSS	Suggested threshold
CADD	2.66	0.24	2.09	5' extended: 7.39, 3' intronic: 0.0964, exonic: 0.39
DSSP	0.01	0.13	0.01	0.30
GeneSplicer	0.18	0.05	0.21	-
MaxEntScan	0.26	0.31	0.24	0.10
MMSplice	1.42	-	1.37	2
NNSPLICE	0.13	0.40	0.30	0.05
Spidex	0.86	-	1.72	5
SpliceAI	0.19	0.18	0.11	0.20
SpliceRover	0.18	0.26	0.10	-
SpliceSiteFinder-like	0.01	0.12	0.09	0.05



4616 variants  
(1924 impacting)  
RNA in vitro studies

# Take home message





## People

Olivier Ardouin

Corinne Bareil

David Baux

Thomas Guignard

Simon Cabello

Souphatta Sasorith

Charles Van Goethem

Laboratoire de génétique moléculaire:  
groupe neurosensoriel

- Anne-Françoise Roux
- Luke Mansard
- Christel Vaché
- Valérie faugère
- Julie Bianchi
- Corinne Baudoin