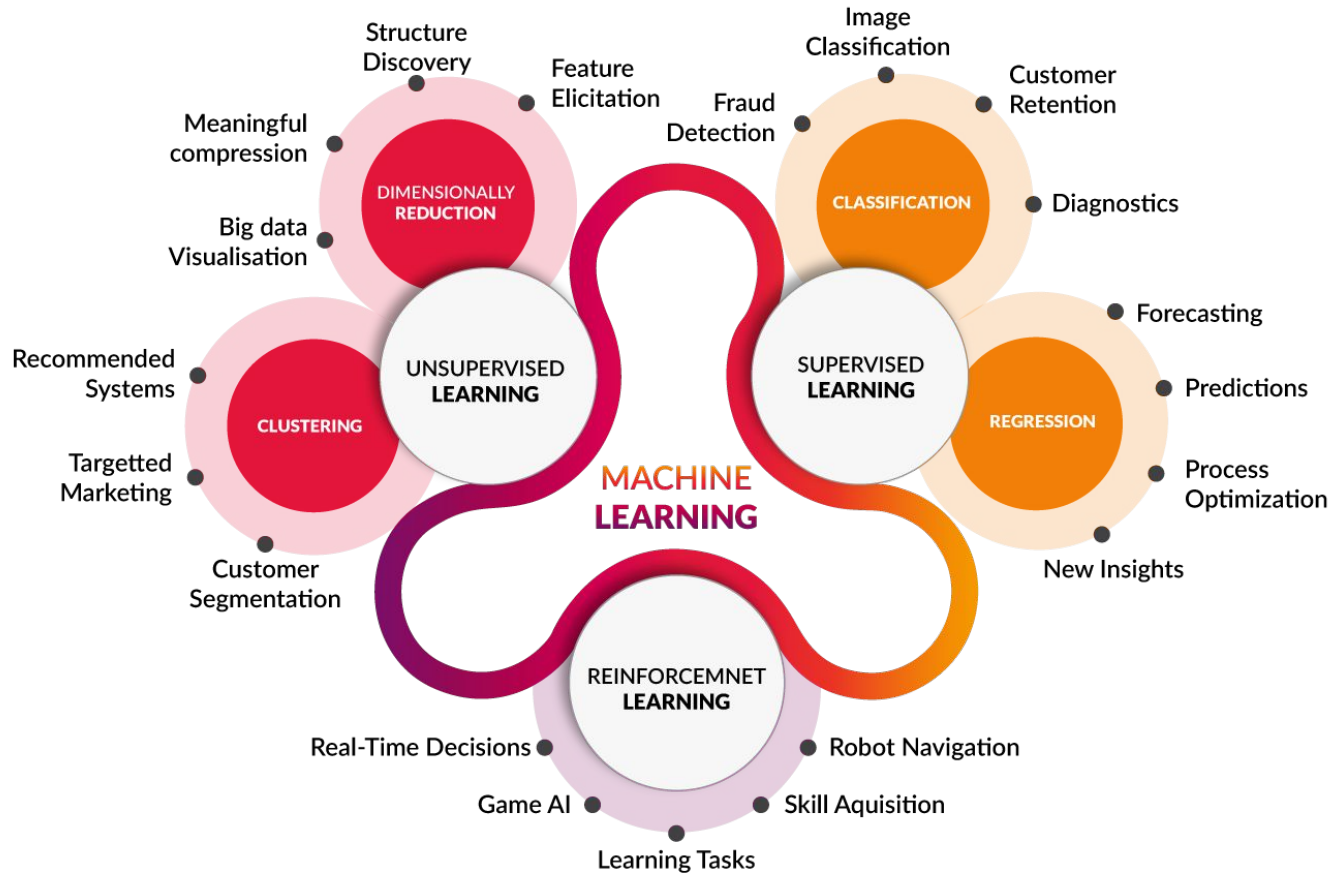


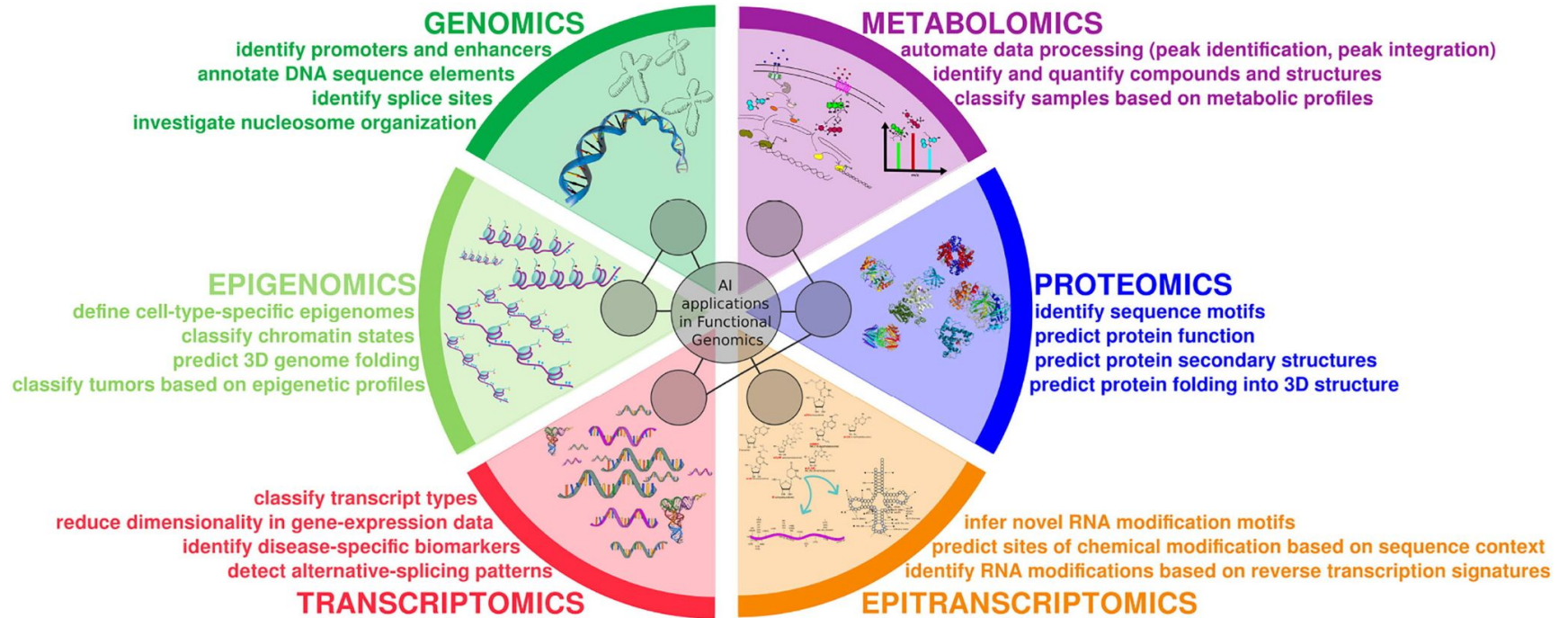
PEPI-IA

Pôle thématique sur l'IA appliquée à la génomique

Amandine Velt et Fabrice Legeai



AI applications in functional genomics



ML in Agricultural and Environment sciences

SEQUENCING

- DNA sequence
- DNA modifications (e.g. methylation)
- Transcription factor binding (e.g. ChIP-seq)
- RNA expression levels (transcriptomics)



GATTACA
TCCTGAT
GAGTAGA
CATGTCT

MICROSCOPY

- Light, fluorescence, electron, flow



iScience

Review

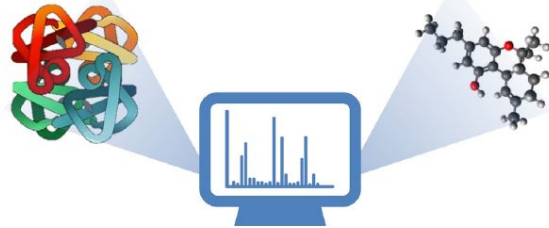
Machine learning in plant science and plant breeding

Aalt Dirk Jan van Dijk,^{1,2,*} Gert Kootstra,³ Willem Kruijer,² and Dick de Ridder¹

SUMMARY

Technological developments have revolutionized measurements on plant genotypes and phenotypes, leading to routine production of large, complex data sets. This has led to increased efforts to extract meaning from these measurements and to integrate various data sets. Concurrently, machine learning has rapidly evolved and is now widely applied in science in general and in plant genotyping and phenotyping in particular. Here, we review the application of machine learning in the context of plant science and plant breeding. We focus on analyses at different phenotype levels, from biochemical to yield, and in connecting genotypes to these. In this way, we illustrate how machine learning offers a suite of methods that enable researchers to find meaningful patterns in relevant plant data.

CellPress
OPEN ACCESS



MASS SPECTROMETRY

- Protein / metabolite quantification
- Protein modifications
- Protein(-ligand) interactions

MACHINE LEARNING

GENOMIC REGIONS

- Pseudogenes
- Crossovers
- Selective sweeps

GENE FUNCTION/REGULATION

- Promoter activity
- Regulatory networks
- Candidates for traits

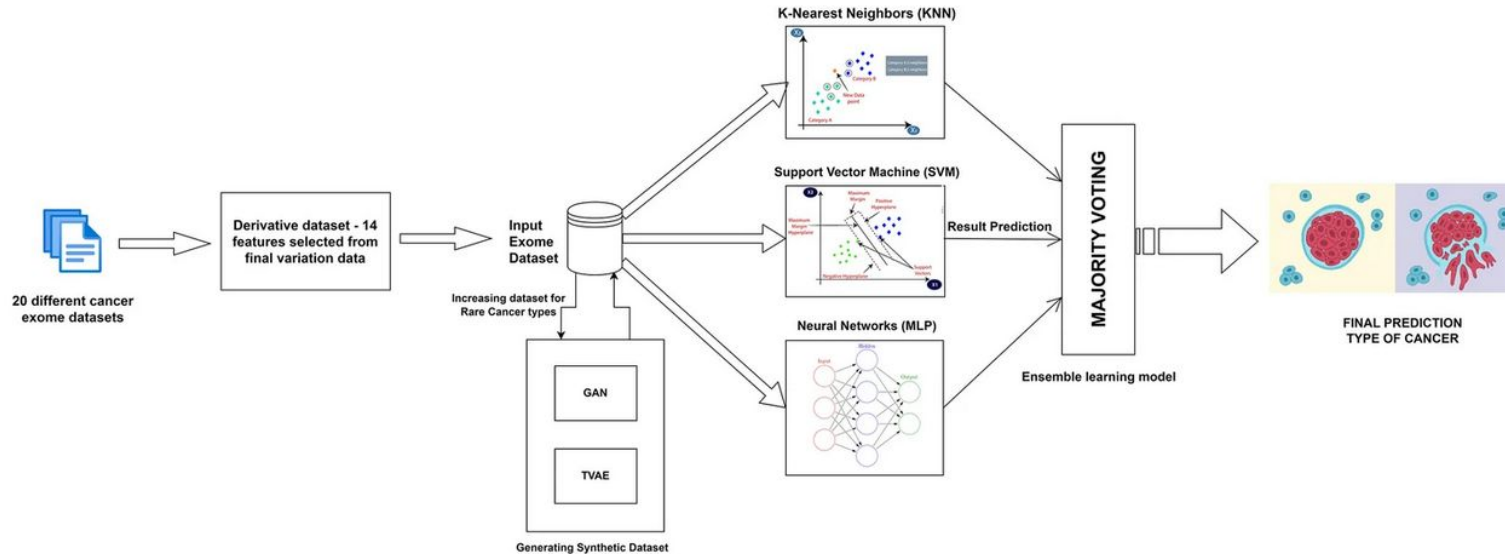
OTHER

- Single-cell RNA sequencing
- Metabolic pathways
- Multi-omics

Applications - GWAS

Implementation of ensemble machine learning algorithms on exome datasets for predicting early diagnosis of cancers

[Abdu Rehaman Pasha Syed](#), [Rahul Anbalagan](#), [Anagha S. Setlur](#), [Chandrashekar Karunakaran](#), [Jyoti Shetty](#), [Jitendra Kumar](#) & [Vidya Niranjana](#)



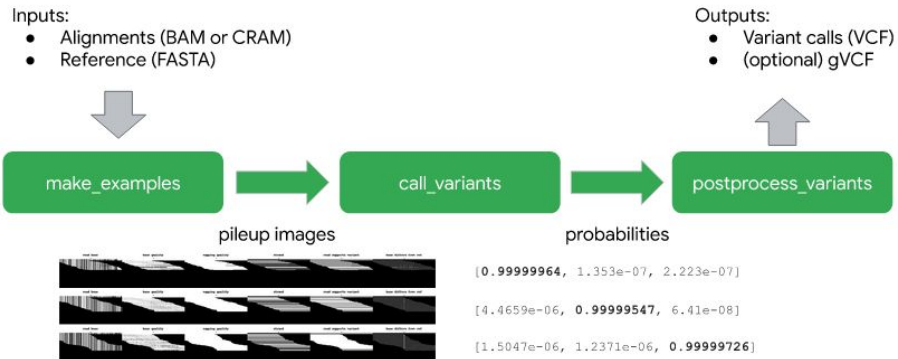
Applications - Variation prediction

nature
biotechnology

LETTERS

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

Ryan Poplin^{1,2}, Pi-Chuan Chang², David Alexander², Scott Schwartz², Thomas Colthurst², Alexander Ku², Dan Newburger¹, Jojo Dijamco¹, Nam Nguyen¹, Pegah T Afshar¹, Sam S Gross¹, Lizzie Dorfman^{1,2}, Cory Y McLean^{1,2} & Mark A DePristo^{1,2}



nature methods

Explore content ▾ About the journal ▾ Publish with us ▾ Subscribe

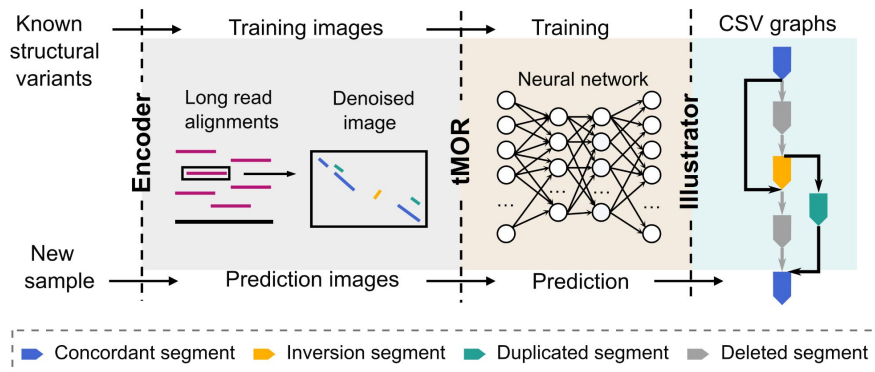
[nature](#) > [nature methods](#) > [brief communications](#) > article

Brief Communication | [Published: 16 September 2022](#)

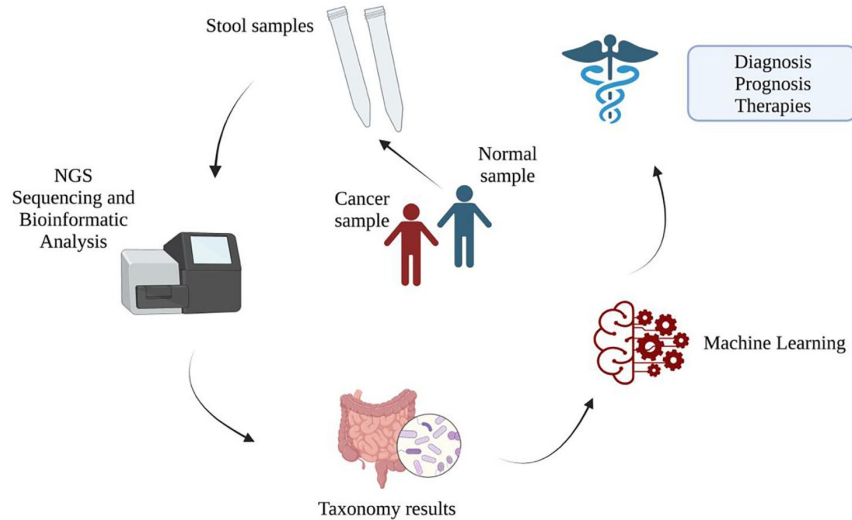
SVision: a deep learning approach to resolve complex structural variants

[Jiadong Lin](#), [Songbo Wang](#), [Peter A. Audano](#), [Deyu Meng](#), [Jacob I. Flores](#), [Walter Kusters](#), [Xiaofei Yang](#), [Peng Jia](#), [Tobias Marshall](#), [Christine R. Beck](#) & [Kai Ye](#) ✉

Nature Methods 19, 1230–1233 (2022) | [Cite this article](#)



Metagenomics



A Review and Tutorial of Machine Learning Methods for Microbiome Host Trait Prediction

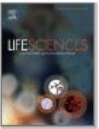
Yi-Hui Zhou^{1*} and Paul Gallins²

¹ Department of Biological Sciences, North Carolina State University, Raleigh, NC, United States, ² Bioinformatics Research Center, North Carolina State University, Raleigh, NC, United States



Life Sciences

Volume 311, Part A, 15 December 2022, 121118



Review article

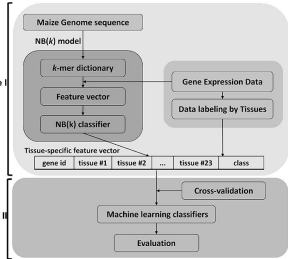
The influence of machine learning technologies in gut microbiome research and cancer studies - A review

Expression data



Predicting Tissue-Specific mRNA and Protein Abundance in Maize: A Machine Learning Approach

Kyoung Tak Cho¹, Taner Z. Sen² and Carson M. Andorf^{3*}



Hanczar et al. BMC Bioinformatics (2020) 21:501
https://doi.org/10.1186/s12859-020-03836-4

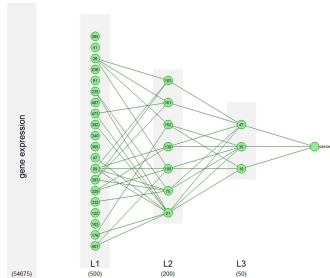
BMC Bioinformatics

RESEARCH ARTICLE

Open Access

Biological interpretation of deep neural network for phenotype prediction based on gene expression

Blaise Hanczar^{1*}, Farida Zehraoui¹, Tina Issa and Mathieu Arles

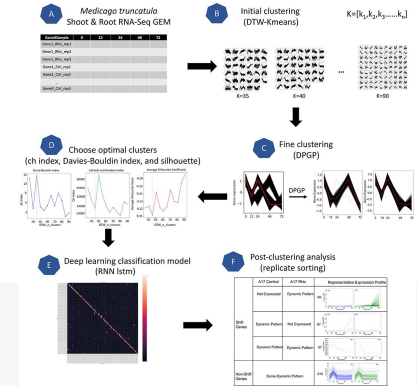
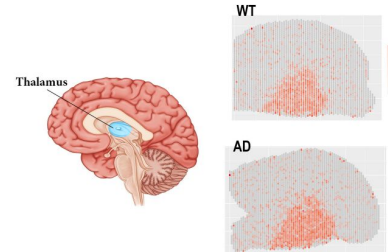


HOME | ABOUT | ARCHIVE | SUBMIT | SUBSCRIBE | ADVERTISE | AUTHOR INFO | CON

Institution: INRAE Inst Natl Recherche pour l'Agriculture, l'Alimentation et l'Environnement

A model-based constrained deep learning approach for clustering spatial-resolved single-cell data

Xiang Lin¹, Le Gao¹, Nathan Whitener², Asraa Ahmed³ and Zhi Wei^{1,4}



ORIGINAL RESEARCH article

Front. Plant Sci., 07 April 2022
Sec. Plant Bioinformatics
https://doi.org/10.3389/fpls.2022.861039

Time Series Transcriptome Analysis in *Medicago truncatula* Shoot and Root Tissue During Early Nodulation

Yueyao Gao¹, Bradley Selee², Elise L. Schnabel¹, William L. Poehlman^{1,3}, Suchitra A. Chavan¹,
Julia A. Frugoli¹ and Frank Alex Feltus^{1,4,5*}

Annotation

BMC Bioinformatics

[Home](#) [About](#) [Articles](#) [Submission Guidelines](#) [Join The Board](#)

Research | [Open Access](#) | [Published: 06 October 2022](#)

EnsembleSplice: ensemble deep learning model for splice site prediction

[Victor Akpokiro](#), [Trevor Martin](#) & [Oluwatosin Oluwadare](#) 

BMC Bioinformatics **23**, Article number: 413 (2022) | 

716 Accesses | 5 Altmetric | [Metrics](#)

PNAS

Identification of the expressome by machine learning on omics data

Ryan C. Sartor^a, Jaclyn Noshay^b, Nathan M. Springer^b, and Steven P. Briggs^{a,1}

^aDivision of Biology, University of California San Diego, La Jolla, CA 92093; and ^bDepartment of Plant Biology, University of Minnesota, St. Paul, MN 55108

Contributed by Steven P. Briggs, July 11, 2019 (sent for review August 14, 2018; reviewed by James A. Birchler and Virginia Walbot)

Accurate annotation of plant genomes remains complex due to the presence of many pseudogenes arising from whole-genome duplication-generated redundancy or the capture and movement

Most researchers study the predicted genes that are derived from whole-genome annotations. These annotation approaches can be complicated by the presence of sequences with homology



The Plant Genome

OPEN ACCESS 

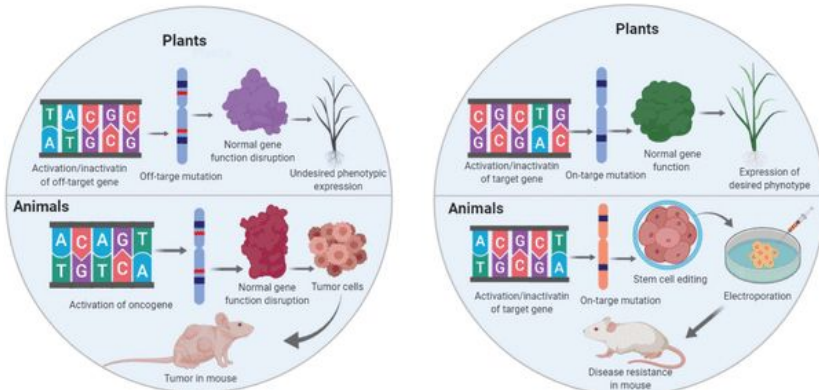
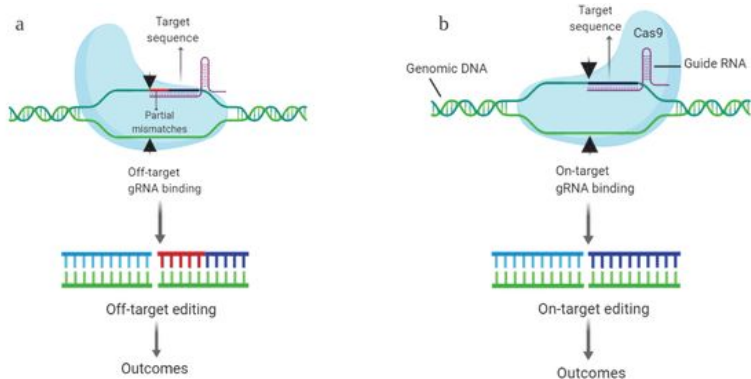
Crop Science
UNIVERSITY OF MINNESOTA

ORIGINAL ARTICLE |  [Open Access](#) |  

ASRpro: A machine-learning computational model for identifying proteins associated with multiple abiotic stress in plants

[Prabina Kumar Meher](#)  [Tanmaya Kumar Sahu](#), [Ajit Gupta](#), [Anuj Kumar](#), [Sachin Rustgi](#) 

Genome editing



<https://www.mdpi.com/2073-4409/9/7/1608>

Prediction of off-target activities for the end-to-end design of CRISPR guide RNAs

Jennifer Listgarten , Michael Weinstein , Benjamin P. Kleinstiver, Alexander A. Sousa, J. Keith Joung, Jake Crawford, Kevin Gao, Luong Hoang, Melih Elibol, John G. Doench  & Nicolo Fusi 

Nature Biomedical Engineering **2**, 38–47 (2018) | [Cite this article](#)

6319 Accesses | 131 Citations | 128 Altmetric | [Metrics](#)

Abstract

Off-target effects of the CRISPR–Cas9 system can lead to suboptimal gene-editing outcomes and are a bottleneck in its development. Here, we introduce two interdependent machine-learning models for the prediction of off-target effects of CRISPR–Cas9. The approach, which we named Elevation, scores individual guide–target pairs, and also aggregates them into a single, overall summary guide score. We demonstrate that Elevation consistently outperforms competing approaches on both tasks. We also introduce an evaluation method that balances errors between active and inactive guides, thereby encapsulating a range of practical use cases. Because of the large-scale and computational demands of the prediction of off-target activities, we have developed a fast cloud-based service (<https://crispr.ml>) for end-to-end guide-RNA design. The service makes use of pre-computed on-target and off-target activity prediction for every genic region in the human genome.

Infrastructure requirements for AI

01

High computing capacity

CPUs for basic AI, GPUs for deep learning

02

Storage capacity

It's fundamental to have the ability to scale storage as the volume of data grows : need to monitor capacity and plan expansion

03

Security

AI can involve handling sensitive data such as patient records, so it is imperative the infrastructure is secured end-to-end

04

Networking infrastructure

Good, fast and reliable networks are essential to maximise the delivery of results.

genOuest, INRAE CollabIA

<https://imotep.inrae.fr/collabIA>

Formations

L'IA en sciences du vivant



Vous ne savez pas ce qu'est l'IA et vous vous demandez si elle peut vous aider à répondre à votre question biologique ? Vous développez des approches d'IA et êtes à la recherche de cas concrets pour les éprouver ? DIGIT-BIO vous invite à suivre son cycle d'animations l'IA en Sciences du Vivant.

Séance introductive - Introduction au Machine Learning - 31/03/2022

- > **Digit-BIO & IA** - Christèle Robert-Granité, Marie-Laure Martin, Julien Chiquet
- > **Introduction à l'IA et au Machine Learning** - Livio Ralaivola

Machine learning pour la classification supervisée - 07/04/2022

- > **Comment faire des prédictions à partir de vos données biologiques grâce au Machine Learning ?** De la régression logistique aux réseaux de neurones - Blaise Hanczary
- > **Explorer le paysage cellulaire avec le Deep Learning** - Emmanuel Moebel, Fadwa Fatmaoui (binôme biologiste - méthodologiste)

Identification de structure par classification non-supervisée - 23 /05/2022

- > **Comment trouver des structures dans vos données biologiques avec la classification non supervisée ?** Méthodes de distance, modèles de mélange ou réseaux de neurones - Cathy Maugis Rabusseau
- > **Comprendre les interactions et l'organisation des réseaux écologiques grâce à l'analyse de réseaux** - Apprentissage non supervisé de structures de graphes - Sophie Donnet, François Massol (binôme biologiste - méthodologiste)

Identification de structure par réduction de dimension - 27/ 06/ 2022

- > **De l'analyse en composantes principales aux auto-encodeurs variationnels** - Stéphanie Allasoinnière
- > **Quelle est la dimension pertinente de l'espace d'expression des gènes ?** - Olivier Gardillon, Franck Picard (binôme biologiste - méthodologiste)

<https://www6.inrae.fr/digitbio/Animations/L-IA-en-sciences-du-vivant>

<https://training.galaxyproject.org/training-material/topics/statistics/>

The screenshot shows the Galaxy Training website interface. At the top, there's a navigation bar with 'Galaxy Training!' and links for 'Contributors', 'Languages', 'Help', 'Extras', and a search bar. The main heading is 'Statistics and machine learning', with a subtitle 'Statistical Analyses for omics data and machine learning using Galaxy tools'. Below this, it says 'You can view the tutorial materials in different languages by clicking the dropdown icon next to the slides (📄) and tutorial (📖) buttons below.' The 'Requirements' section recommends looking at 'Introduction to Galaxy Analyses'. The 'Material' section features a table of lessons with search and navigation icons.

Lesson	Slides	Hands-on	Recordings	Input dataset	Workflows	Galaxy servers
A Docker-based interactive Jupyterlab powered by GPU for artificial intelligence in Galaxy interactive-tools machine-learning deep-learning jupyter-lab image-segmentation protein-3D-structure	📄	📖	📺	📄	📄	🌐
Age prediction using machine learning	📄	📖	📺	📄	📄	🌐
Basics of machine learning	📄	📖	📺	📄	📄	🌐
Classification in Machine Learning	📄	📖	📺	📄	📄	🌐
Clustering in Machine Learning	📄	📖	📺	📄	📄	🌐
Deep Learning (Part 1) - Feedforward neural networks (FNN)	📄	📖	📺	📄	📄	🌐
Deep Learning (Part 2) - Recurrent neural networks (RNN)	📄	📖	📺	📄	📄	🌐

<https://www.fun-mooc.fr/en/courses/machine-learning-python-scikit-learn/>

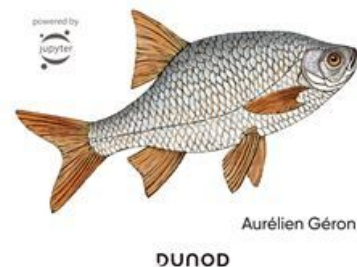
The screenshot shows the FUN MOOC website for the course 'Machine learning in Python with scikit-learn'. It features the FUN MOOC logo, a search bar, and navigation links for 'Home', 'News', 'Courses', 'GRADED', 'Diplômes', and 'Organizations'. The course title is prominently displayed, along with the Inria logo. Below the title, it lists 'Duration: 12 weeks', 'Effort: 36 hours', and 'Pace: ~2h45/week'. A brief description states: 'Build predictive models with scikit-learn and gain a practical understanding of the strengths and limitations of machine learning!'. Social media icons for Facebook, Twitter, LinkedIn, and YouTube are also present.

O'REILLY
Machine Learning
avec
Scikit-Learn

Mise en œuvre et cas concrets



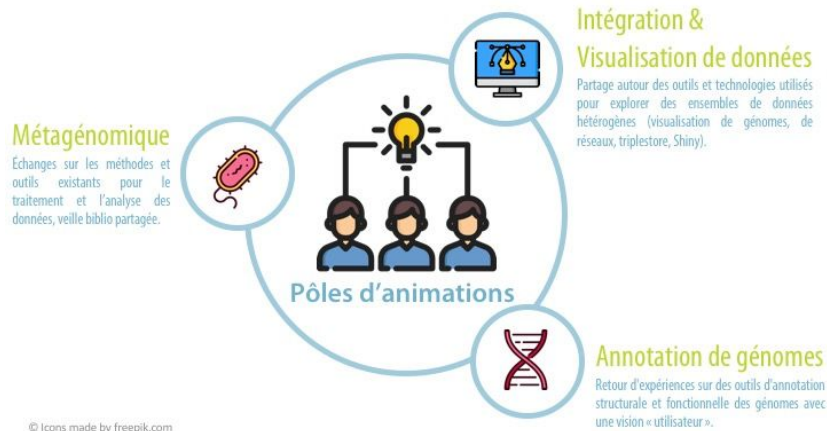
O'REILLY
Deep Learning
avec
Keras et TensorFlow
Mise en œuvre et cas concrets



PEPI-IBIS Pôle IA

Un nouveau pôle du PEPI-IBIS qui prend la suite du pôle Intégration et Visualisation de données et qui complète

- PEPI-Annot
- PEPI-métagénomique
- PEPI-Text-mining



Les activités du pôle

Réunions régulières (1h / 2 mois)

1 exposé

- présentations de résultats, exposé de problématiques, besoins et conseils
- retours d'expérience, tests en cours (outils, méthodes)
- synthèse de publications d'intérêts
- infrastructure
- retour sur des formations
- des discussions

Journée tous les 2 ans sur un site (si financement SAPI)

- présentations générales
- ateliers
- formations

Retour lors des AG du PEPI (tous les 2 ans)

Contact

Liste de diffusion : pepi-bioinfostats-ia@inrae.fr

Site web : <https://pepi-ibis.inrae.fr/ia-genomique>

Orateurs

Raphaël Mourad (CBI, Université Toulouse 3, en délégation à INRAE MIAT-Mathnum) - Deep learning pour la génomique

Camille Kergal & Christophe Hitte (Equipe génétique du chien, Institute Genetics & Development, Rennes) - Méthode d'apprentissage profond pour l'analyse génomique des cancers canins comme modèles des cancers humains