



Machine & deep learning for personalized oncology

Loic Verlingue
MD Medical Oncology
PhDc Data Science

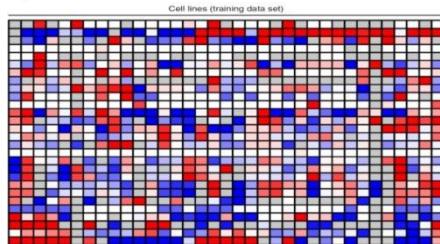
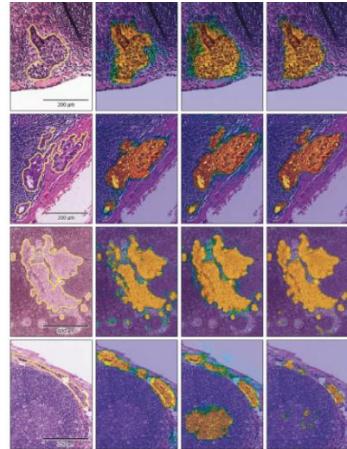
21-22/01/2021



Plan

- A ML / DL workflow
- ML / DL is impacting
 - Diagnosis procedures
 - Drug development
 - Patients' monitoring
- Because of many cool things in deep learning

ML/DL in oncology

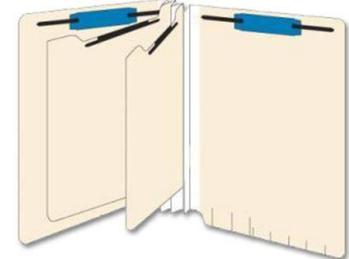


Diagnostics

Molecular biology
& prognosis

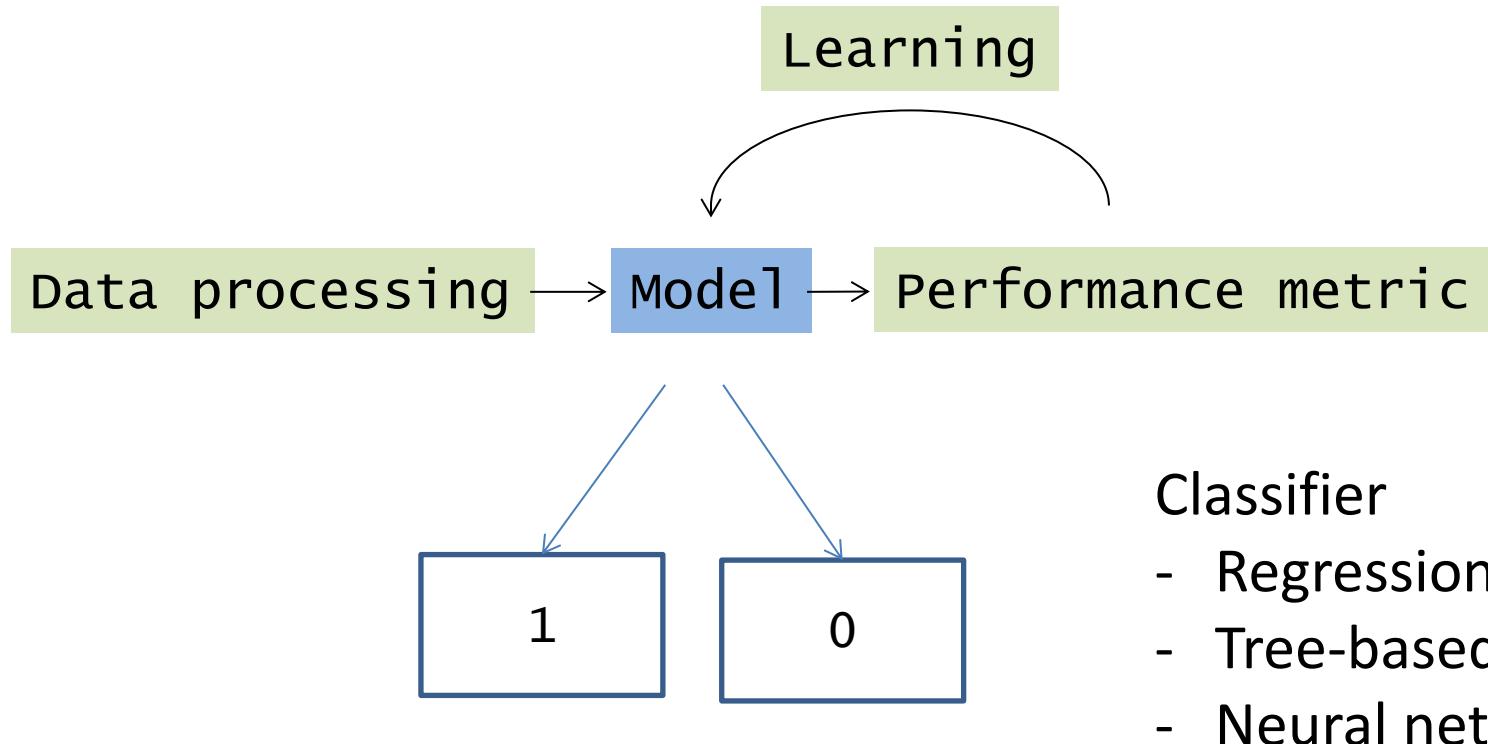


Drug development
& prediction

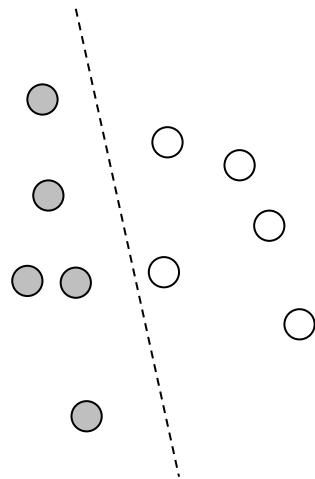


Monitoring patients

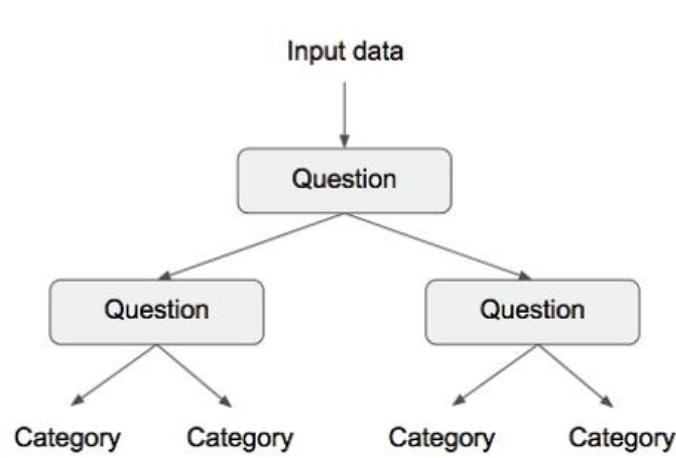
Machine learning workflow



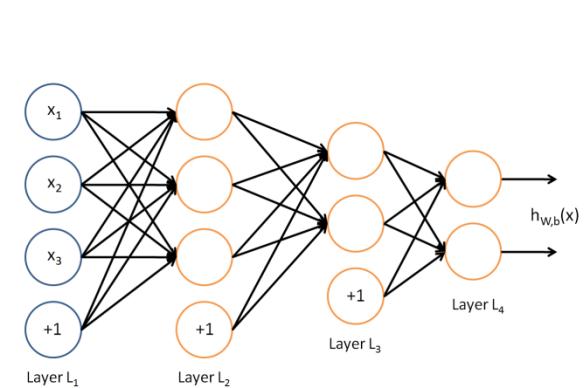
Families of models for classification



Regression /
Kernel approach eg SVM

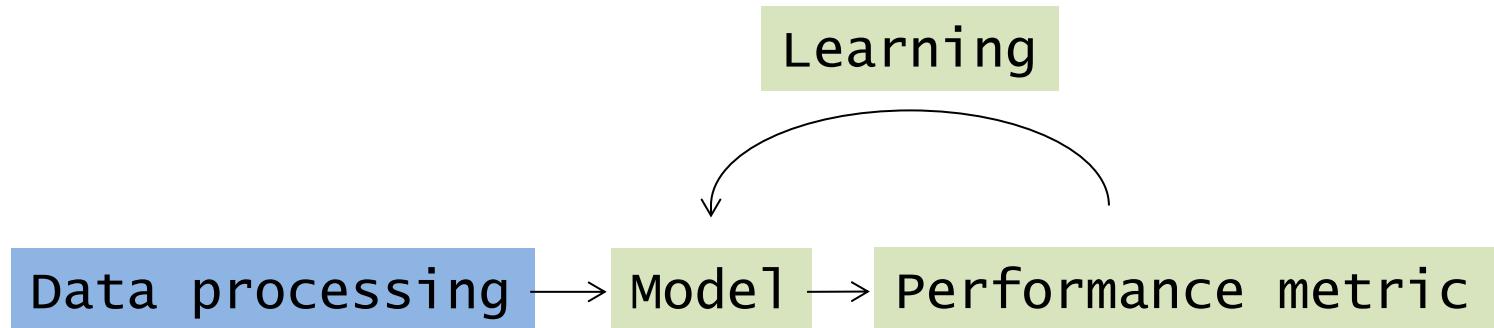


Decision tree /
Ensemble modeling



Multi Layer Perceptron =
Artificial Neural Network

Data Processing

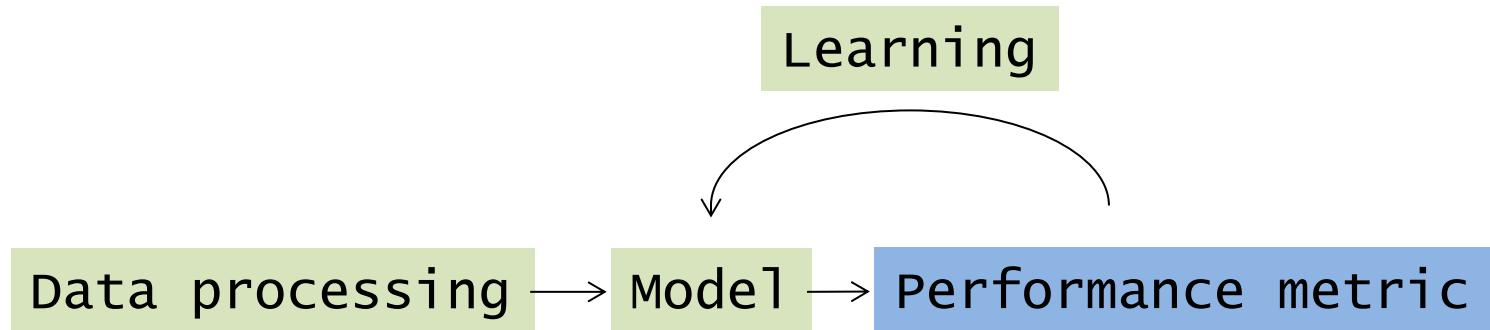


Numerical data -> ok for models
*

Natural language data -> turn to numerical
*

Images -> use pixels values

Evaluation metrics



Comment évaluer les prédictions du modèle?

Métriques pour classification binaire

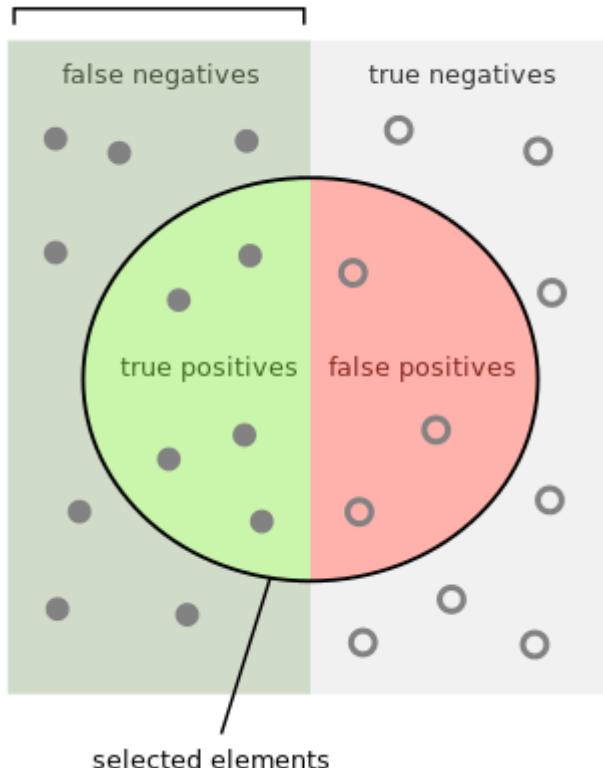
→ Matrice de confusion (1=SF)

		Classe réelle	
		0	1
Classe prédictive	0	TN	FN
	1	FP	TP

→ Que préférez-vous évaluer?

Métriques pour classification binaire

Positives Negatives



How many relevant items are selected?
e.g. How many sick people are correctly identified as having the condition.

$$\text{Sensitivity} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$

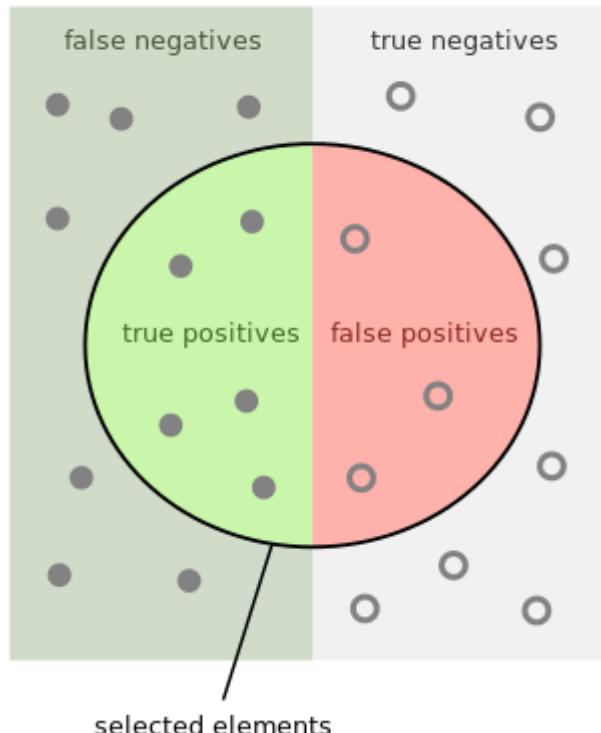
How many negative selected elements are truly negative?
e.g. How many healthy people are identified as not having the condition.

$$\text{Specificity} = \frac{\text{true negatives}}{\text{true negatives} + \text{false positives}}$$

		Classe réelle	
		0	1
Classe prédictive	0	TN	FN
	1	FP	TP

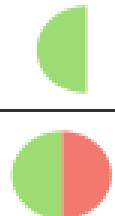
Métriques pour classification binaire

Positives Negatives



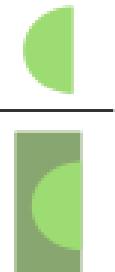
How many selected items are relevant?

Precision =



How many relevant items are selected?

Recall =

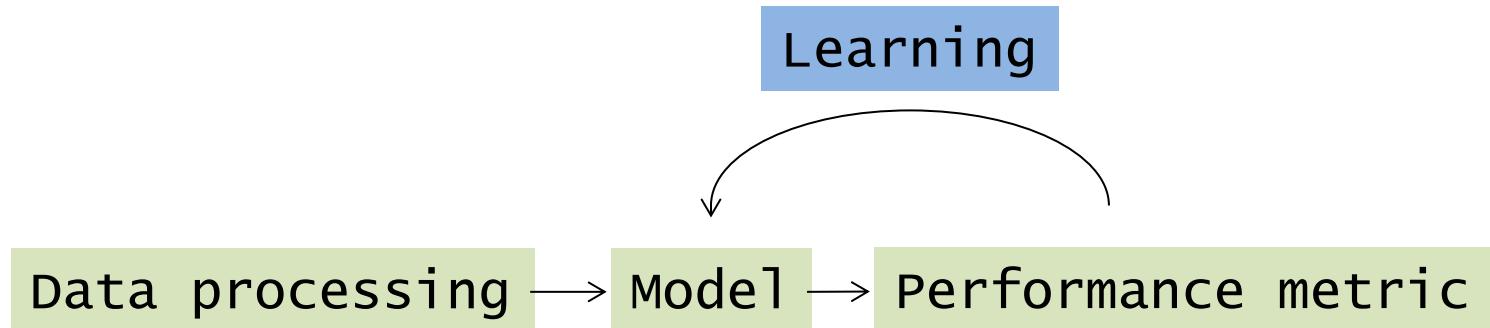


		Classe réelle	
		0	1
Classe prédictive	0	TN	FN
	1	FP	TP

Métriques pour classification binaire

	Total population	True condition		Prevalence = $\frac{\sum \text{Condition positive}}{\sum \text{Total population}}$	Accuracy (ACC) = $\frac{\sum \text{True positive} + \sum \text{True negative}}{\sum \text{Total population}}$
Predicted condition	Predicted condition positive	Condition positive True positive , Power	Condition negative False positive , Type I error	Positive predictive value (PPV), Precision = $\frac{\sum \text{True positive}}{\sum \text{Predicted condition positive}}$	False discovery rate (FDR) = $\frac{\sum \text{False positive}}{\sum \text{Predicted condition positive}}$
	Predicted condition negative	False negative , Type II error	True negative	False omission rate (FOR) = $\frac{\sum \text{False negative}}{\sum \text{Predicted condition negative}}$	Negative predictive value (NPV) = $\frac{\sum \text{True negative}}{\sum \text{Predicted condition negative}}$
	True positive rate (TPR), Recall, Sensitivity, probability of detection = $\frac{\sum \text{True positive}}{\sum \text{Condition positive}}$	False positive rate (FPR), Fall-out, probability of false alarm = $\frac{\sum \text{False positive}}{\sum \text{Condition negative}}$	Positive likelihood ratio (LR+) = $\frac{\text{TPR}}{\text{FPR}}$	Diagnostic odds ratio (DOR) = $\frac{\text{LR+}}{\text{LR-}}$	F_1 score = $\frac{2}{\frac{1}{\text{Recall}} + \frac{1}{\text{Precision}}}$
	False negative rate (FNR), Miss rate = $\frac{\sum \text{False negative}}{\sum \text{Condition positive}}$	Specificity (SPC), Selectivity, True negative rate (TNR) = $\frac{\sum \text{True negative}}{\sum \text{Condition negative}}$	Negative likelihood ratio (LR-) = $\frac{\text{FNR}}{\text{TNR}}$		

Quelle méthode d'apprentissage ?

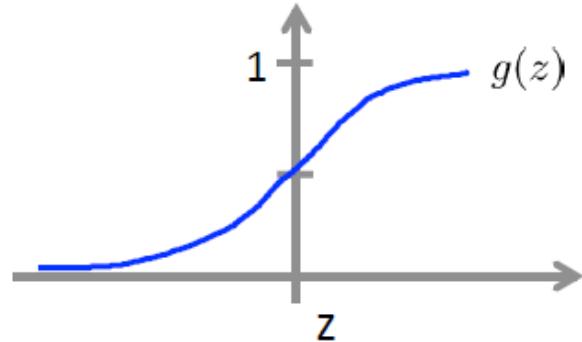


Probability to belong to a class

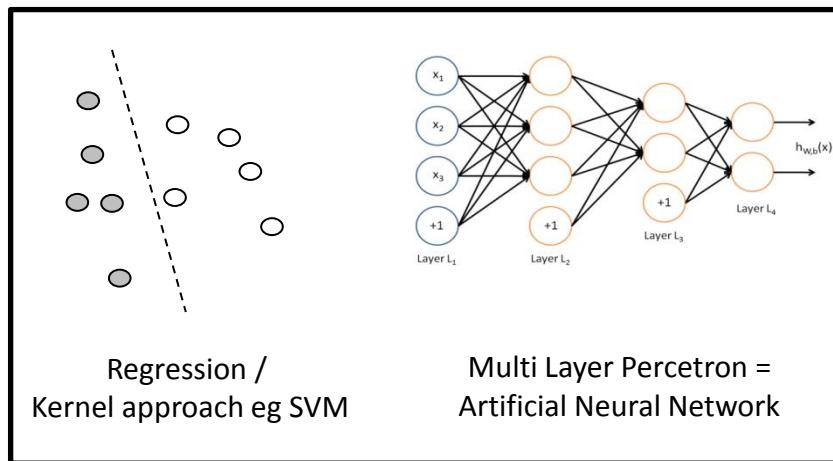
Logistic regression with a sigmoid function

$$h_{\theta}(x) = g(\theta^T x)$$

$$g(z) = \frac{1}{1+e^{-z}}$$



Ok for:



Calculate the cost / error

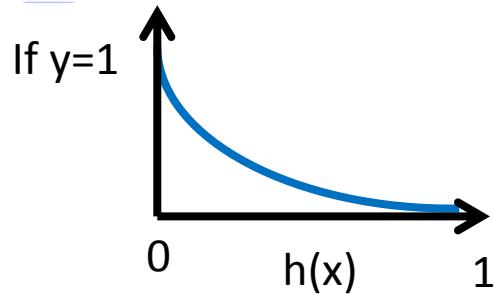
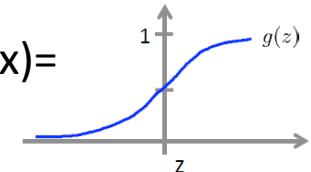
Logistic regression cost function

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

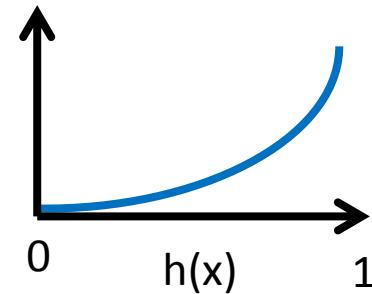
$$= -\frac{1}{m} \left[\sum_{i=1}^m y^{(i)} \log h_\theta(x^{(i)}) + (1 - y^{(i)}) \log (1 - h_\theta(x^{(i)})) \right]$$

prediction

truth



If $y=0$



Parameters' updates with gradient descent

Compute the cost

$$J(\theta) = \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_\theta(x^{(i)}), y^{(i)})$$

Gradient descent:

Repeat {

Parameters' updates

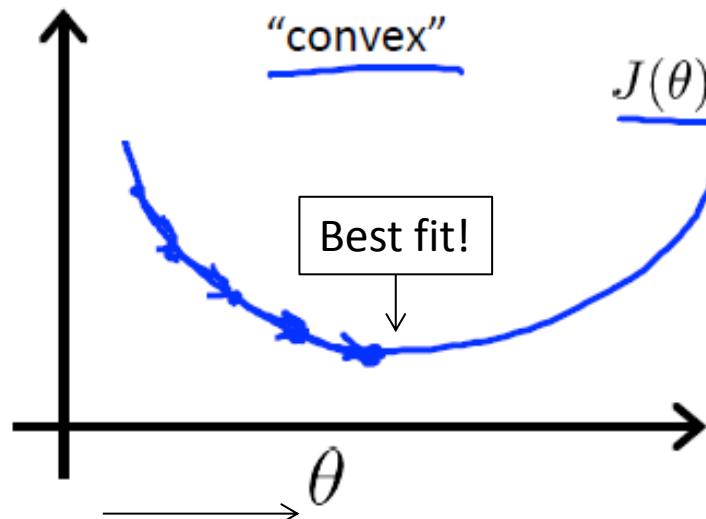
$$\theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta)$$

Derivative of the cost

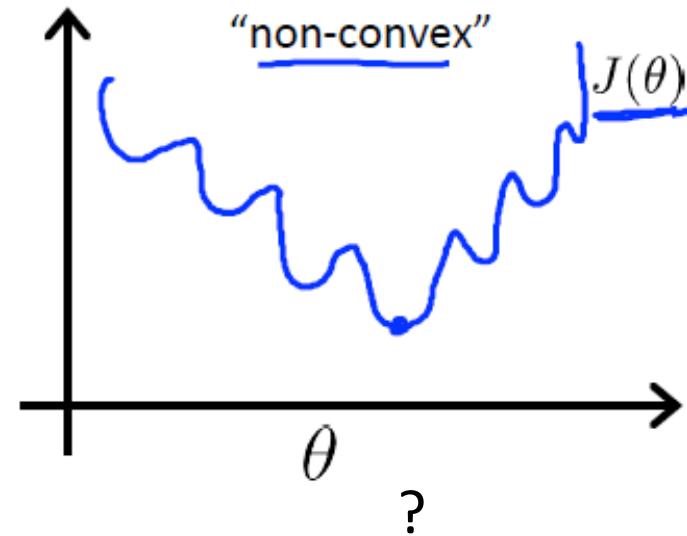
}

Learning rate

Optimisation of the cost



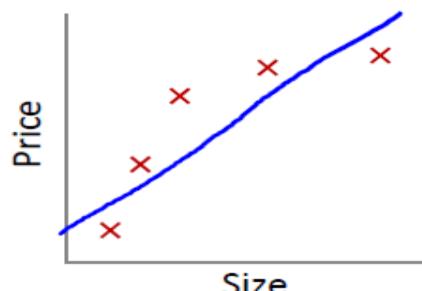
Updates with
gradient descent



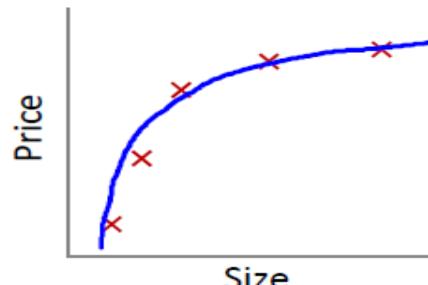
Figures from Andrew Ng Coursera

Etes vous sûr que les
futurs exemples seront
bien prédits?

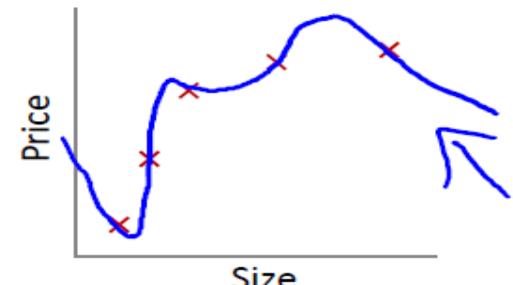
Généralisation du modèle



$\rightarrow \theta_0 + \theta_1 x$
"Underfit" "High bias"



$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2$
"Just right"

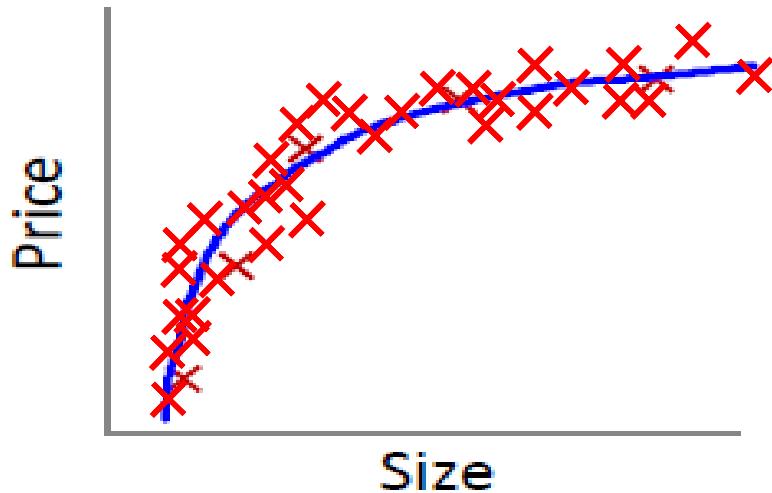


$\rightarrow \theta_0 + \theta_1 x + \theta_2 x^2 + \theta_3 x^3 + \theta_4 x^4$
"Overfit" "High variance"

Overfitting: If we have too many features, the learned hypothesis may fit the training set very well ($J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_\theta(x^{(i)}) - y^{(i)})^2 \approx 0$), but fail to generalize to new examples (predict prices on new examples).

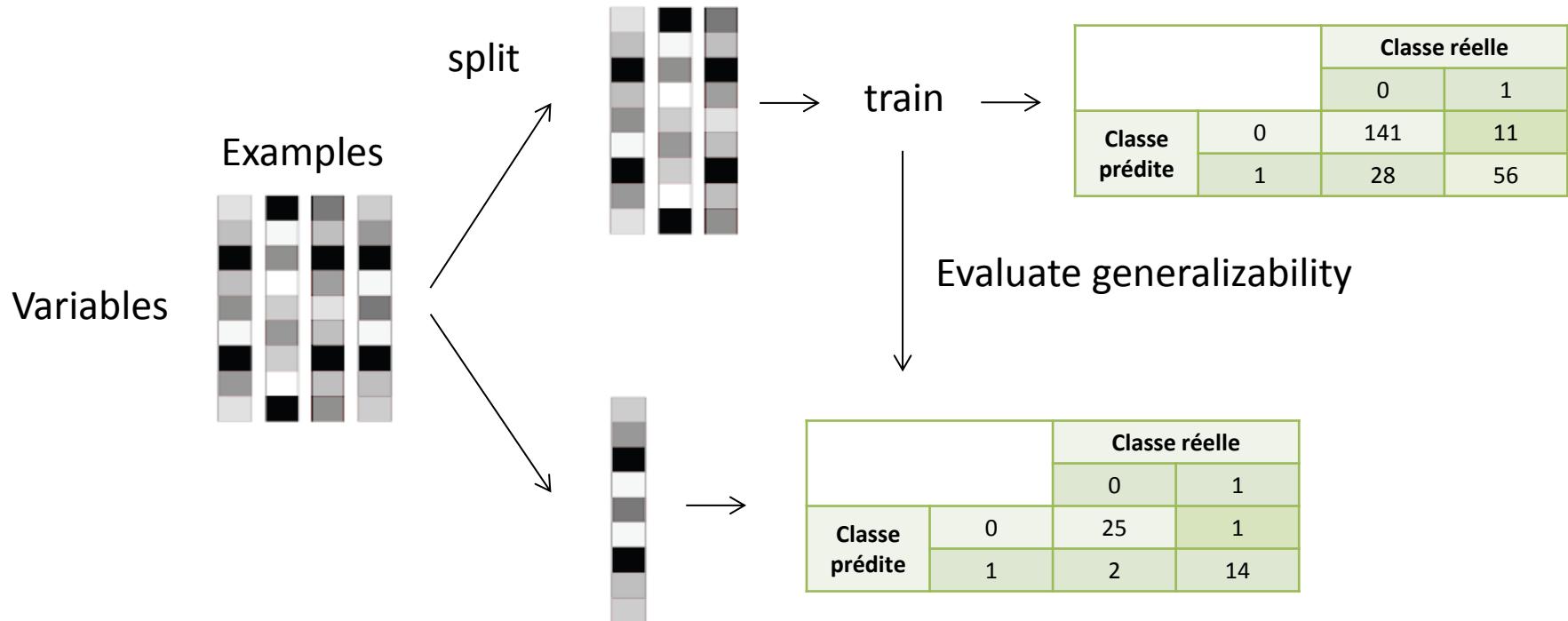
The problem of overfitting

One solution: increase observations

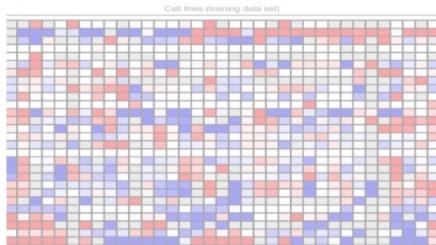
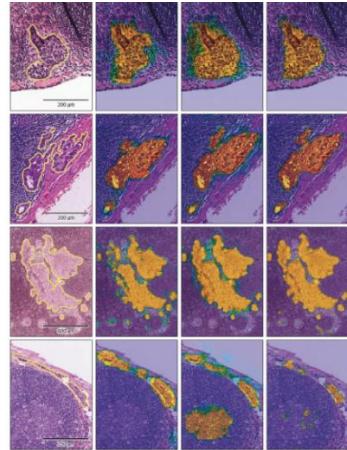


Other solutions: penalize the learning constrain the model

ML generalization testing workflow



ML/DL in oncology



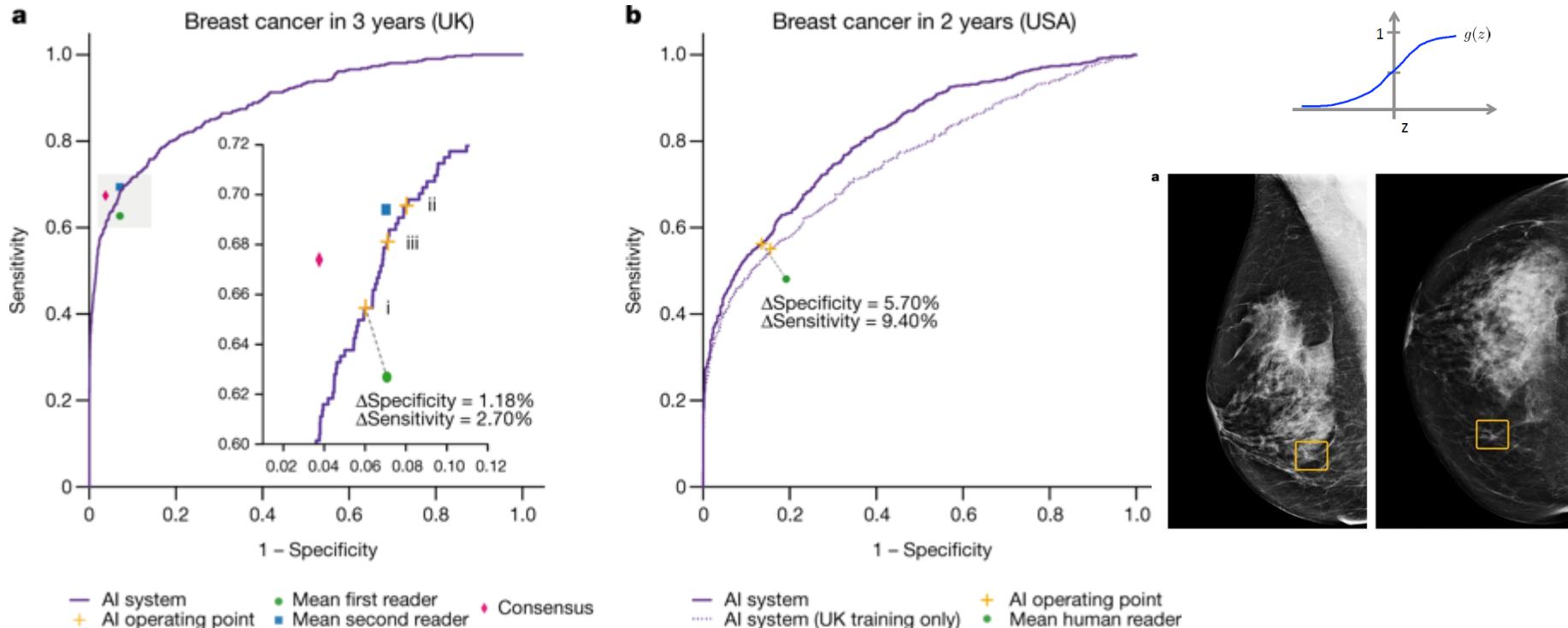
Diagnostics

Molecular biology
& pronostic

Drug development
& prediction

Monitoring patients

Breast cancer screening

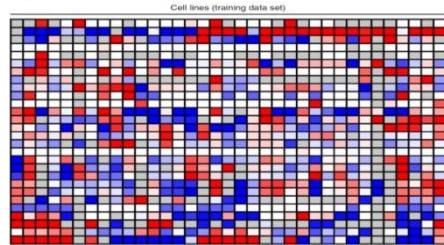
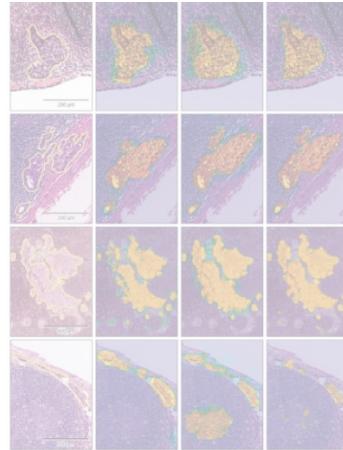


International evaluation of an AI system for breast cancer screening, Scott Mayer McKinney, Marcin Sieniek, [...]Shravya Shetty, Nature volume 577, pages89–94(2020)

Many others

Main Objective	Tumor type	Image Modality	Number of data points (training set + validation/test set)	Endpoints and results	Reference
Diagnosis	Lung cancer	Low dose CT-scan	102 + 70	Prediction of malignancy in lung nodules. Accuracy: 70-77%	Liu et al., 2016 ⁵
Diagnosis	Lung cancer	CT-scan	312 + 288	Malignant nodules accuracy: 80% Benign nodules accuracy: 79%	Hawkins et al., 2016 ⁶
Diagnosis	Lung ADK	CT-scan	54 + 86	Correlation with invasive versus non-invasive adenocarcinomas: Spearman R = 0.87-0.89 (p<0.0001)	Maldonado et al., 2013 ⁷
Diagnosis	NSCLC	CT-scan	81 + 48	Distinguish ADK versus SCC: AUC=0.893-0.905	Zhu et al., 2018 ⁸
Diagnosis	NSCLC	CT-scan	160 + 235	Discrimination between invasive and non-invasive ADK: Accuracy: 94.0-90.8%	Fan et al., 2018 ⁹
Diagnosis	Prostate Cancer	MRI	54 + 0	Radiomics model for cancer diagnosis: AUC=0.955	Wang et al., 2017 ¹⁰
Diagnosis	Prostate Cancer	MRI	147 + 0	Benign versus malignant lesions: Accuracy 93%	Fehr et al., 2015 ¹¹
Molecular characterization	Breast Cancer	MRI	91 + 0	Radiomics features used to predict molecular subtypes, AUC: 0.65-0.89	Li et al., 2016 ¹²
Molecular characterization	Breast cancer	MRI	48 + 0	Significant correlation between radiomics features and luminal B subtype (p<0.0015)	Mazurowski et al., 2014 ¹³
Molecular characterization	Breast cancer	MRI	461 + 461	Prediction of molecular subtype: Luminal A, AUC=0.697 TNBC, AUC=0.682 ER status, AUC=0.649 PR status, AUC=0.622	Saha et al., 2018 ¹⁴
Molecular characterization	NSCLC	PET-CT	26 + 0	Association with gene expression pathways with an accuracy ranging from 59% - 83%	Gevaert et al., 2012 ¹⁵
Molecular characterization	Lung ADK	CT-scan	385 + 0	Prediction of EGFRm status with combination of radiomics signature and clinical variables: AUC=0.778	Liu et al., 2016 ¹⁶
Molecular characterization	GBM	MRI	23 + 0	Correlation of radiomics features with multiple molecular pathways	Jamshidi et al., 2014 ¹⁷
Molecular characterization	GBM	MRI	22 + 110	Significant associations with tumor hypoxia (p<0.012); proliferation gene-expression signature (p<0.0017); EGFR protein overexpression (p<0.002)	Diehn et al., 2008 ¹⁸
Molecular characterization	HNSSC	PET-CT + CT-scan	42 + 79	Correlation with hypoxia: AUC=0.755-0.833 Hypoxia measured by ¹⁸ F-FMISO PET-CT.	Crispin-Ortuzar et al., 2018 ¹⁹
Molecular characterization	CRC	CT-scan	61 + 56	Radiomics signature predicts: KRAS mutation, AUC=0.829-0.869 NRAS mutation, AUC=0.757-0.686 BRAF mutation, AUC=0.833-0.857	Yang et al., 2018 ²⁰
Molecular characterization	Rectal cancer	MRI	114 + 0	pCR diagnosis after neo-adjuvant treatment : AUC=0.93	Horvat et al., 2018 ²¹
Prognosis	NSCLC & HNSCC	CT-scan	474 + 545	Radiomics signature correlates with stage and OS: NSCLC CI=0.65 HNSCC CI=0.69	Aerts et al., 2014 ²²
Prognosis	Lung ADK	CT-scan	98 + 84	Risk of metastases relapse CI=0.61	Coroller et al., 2015 ²³
Prognosis	NSCLC	CT-scan	22 + 0	Recurrence prediction based on post-treatment imaging: Accuracy 73-77%	Mattonen et al., 2015 ²⁴
Prognosis	NSCLC	PET-CT	70 + 31	Risk of distant metastasis: CI=0.71 When combining with histologic type: CI=0.80	Wu et al., 2016 ²⁵
Prognosis	GBM	MRI	79 + 40	Prognosis OS (p<0.001)	Kickingereder et al., 2016 ²⁶
Prognosis	GBM	MRI	121 + 60	Prognosis prediction errors decreased 36% for PFS and 37% for OS compared to clinical features alone	Kickingereder et al., 2018 ²⁷
Prognosis	GBM	MRI	126 + 165	Significant stratification for OS (p<0.001) and PFS (p<0.000021)	Grossman et al., 2017 ²⁸
Prognosis	Colorectal cancer	CT-scan	326 + 200	Lymph node metastasis prediction: CI=0.79	Huang et al., 2016 ²⁹
Prognosis	Nasopharyngeal carcinoma	MRI	88 + 30	Radiomics signature associated with PFS: CI=0.761	Zhang et al., 2018 ³⁰
Prognosis	Breast Cancer	MRI	84 + 0	Prediction of risk recurrence concordance as assessed by PAM50 or Mammoprint: AUC=0.55-0.88	Li et al., 2016 ³¹
Prognosis	Breast Cancer	MRI	92 + 54	Sentinel lymph node metastasis based on radiomics features of pre-operative MRI: AUC=0.770-0.863	Dong et al., 2018 ³²
Prognosis	Breast cancer	MRI	194 + 100	Radiomics features correlate with DFS: p=0.002 in the training set P=0.036 in the validation set	Park et al., 2018 ³³
Prognosis	High-grade osteosarcoma	CT-scan	102 + 48	OS prognosis: AUC=0.73-0.86	Wu et al., 2018 ³⁴
Treatment effect	NSCLC	PET-CT	47 + 0	Correlation of radiomics features with treatment response (log(0.01))	Cook et al., 2015 ³⁵

ML/DL in oncology



Diagnostics

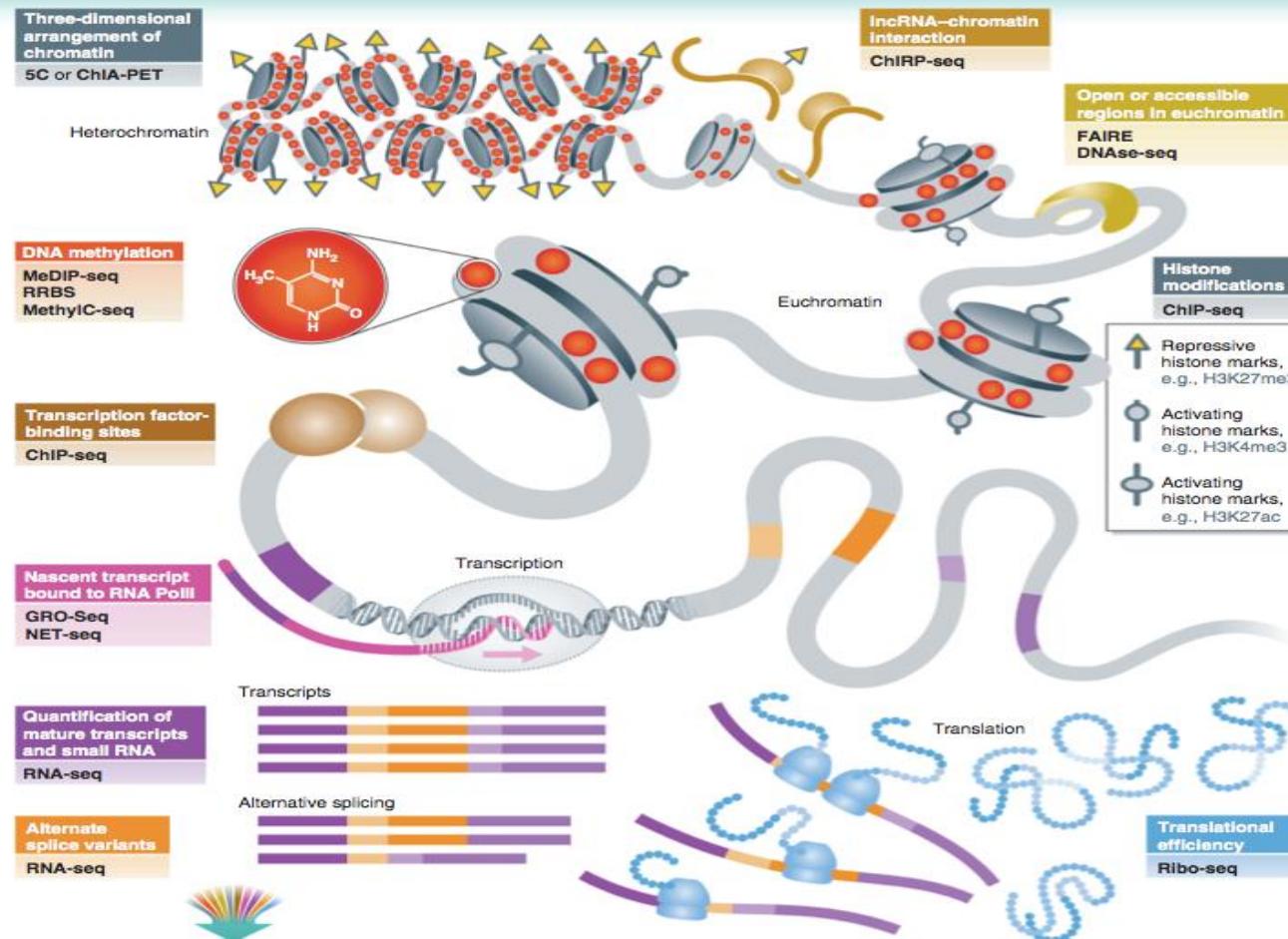
Molecular biology
& pronostic

Drug development
& prediction

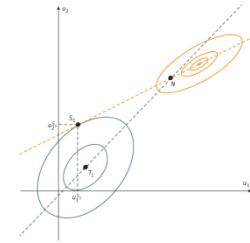
Monitoring patients



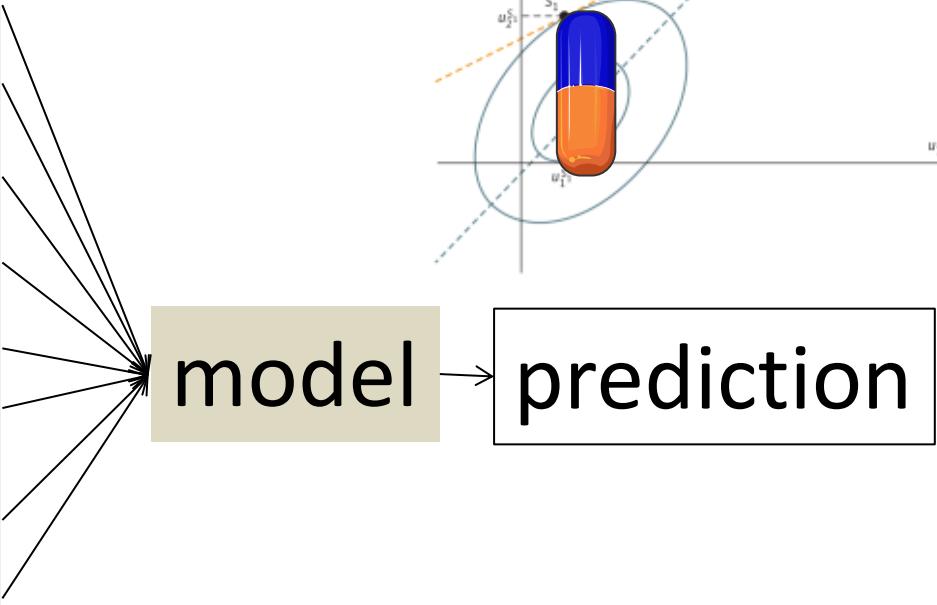
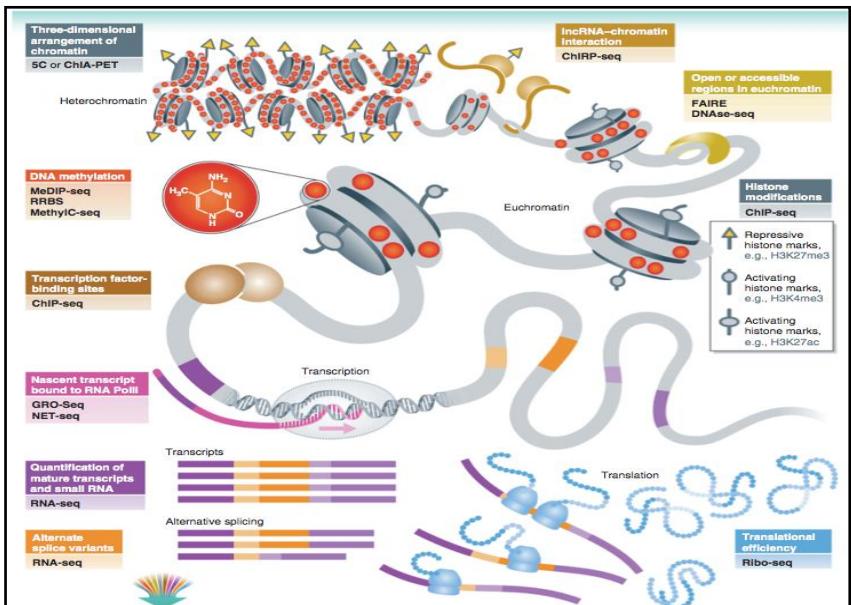
Complexity



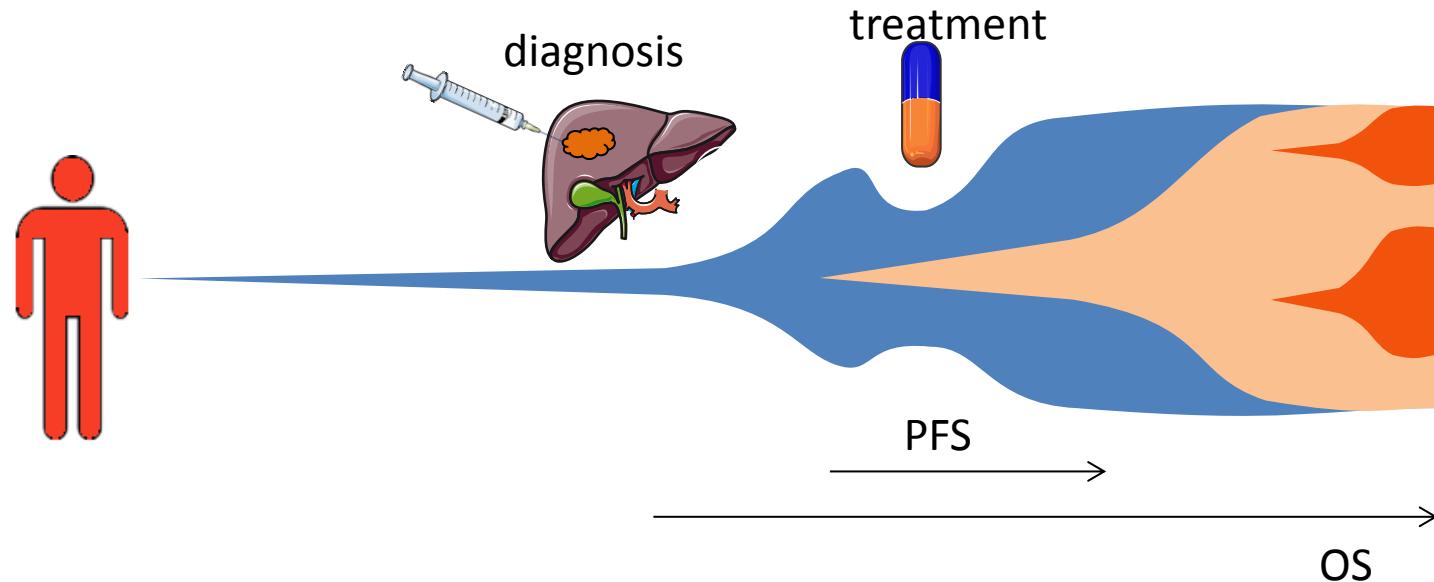
model → prediction



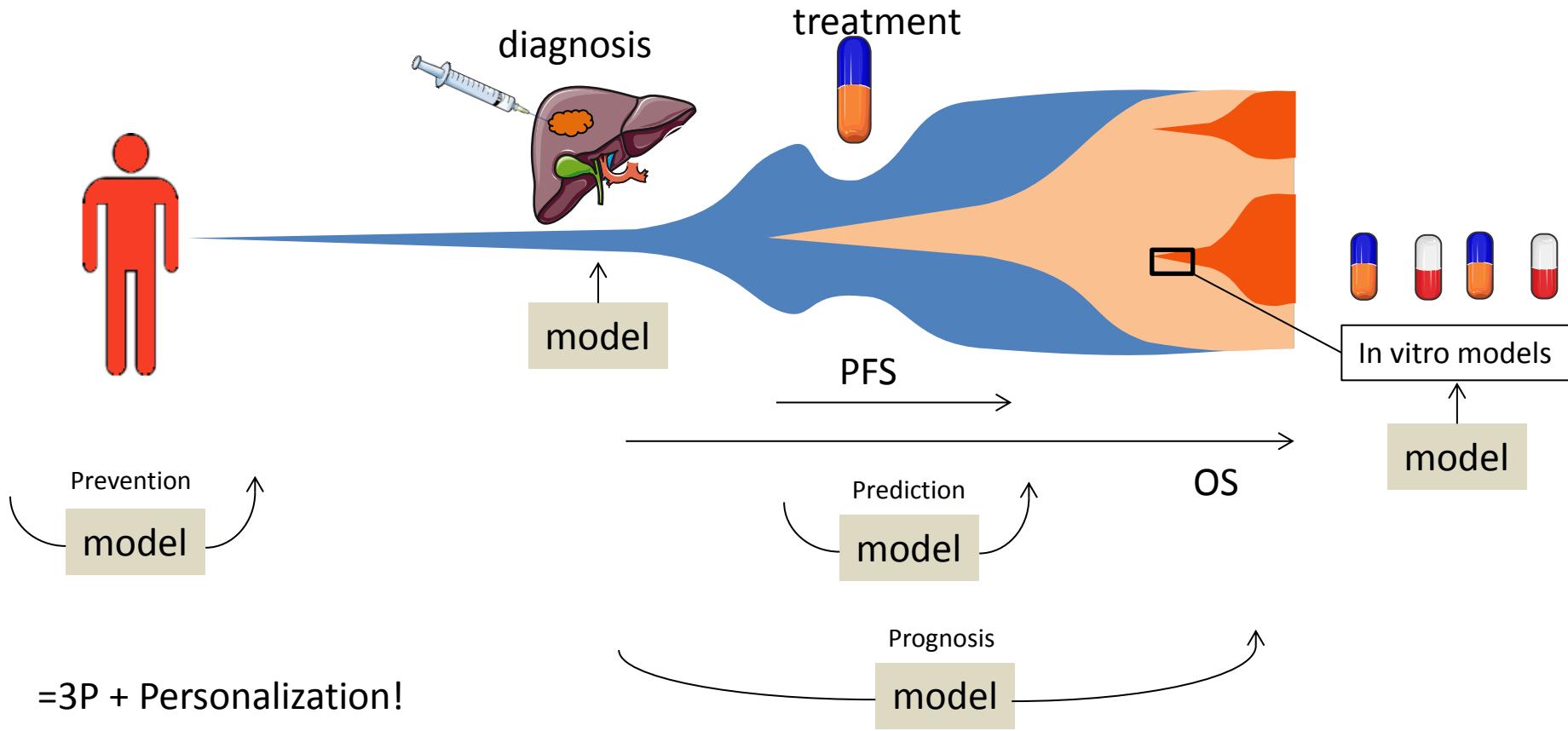
Models



The patient & his tumor's biology



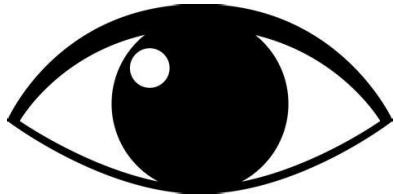
What do we expect from models?



Models

The models you have

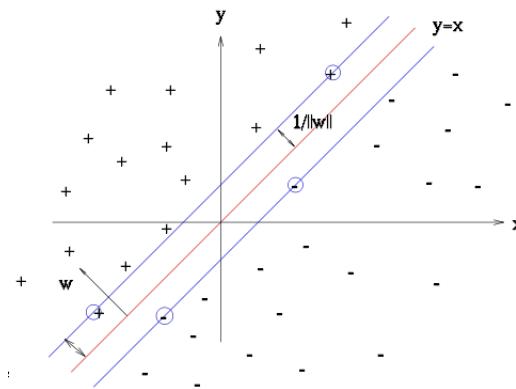
Human



Use your knowledge
and reasoning

Examples: Christophe Massard

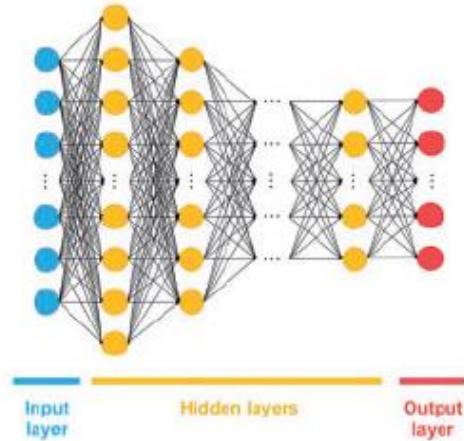
Statistics



Linear models easy to
explain

Cox model
logReg eg: RMH

Deep learning



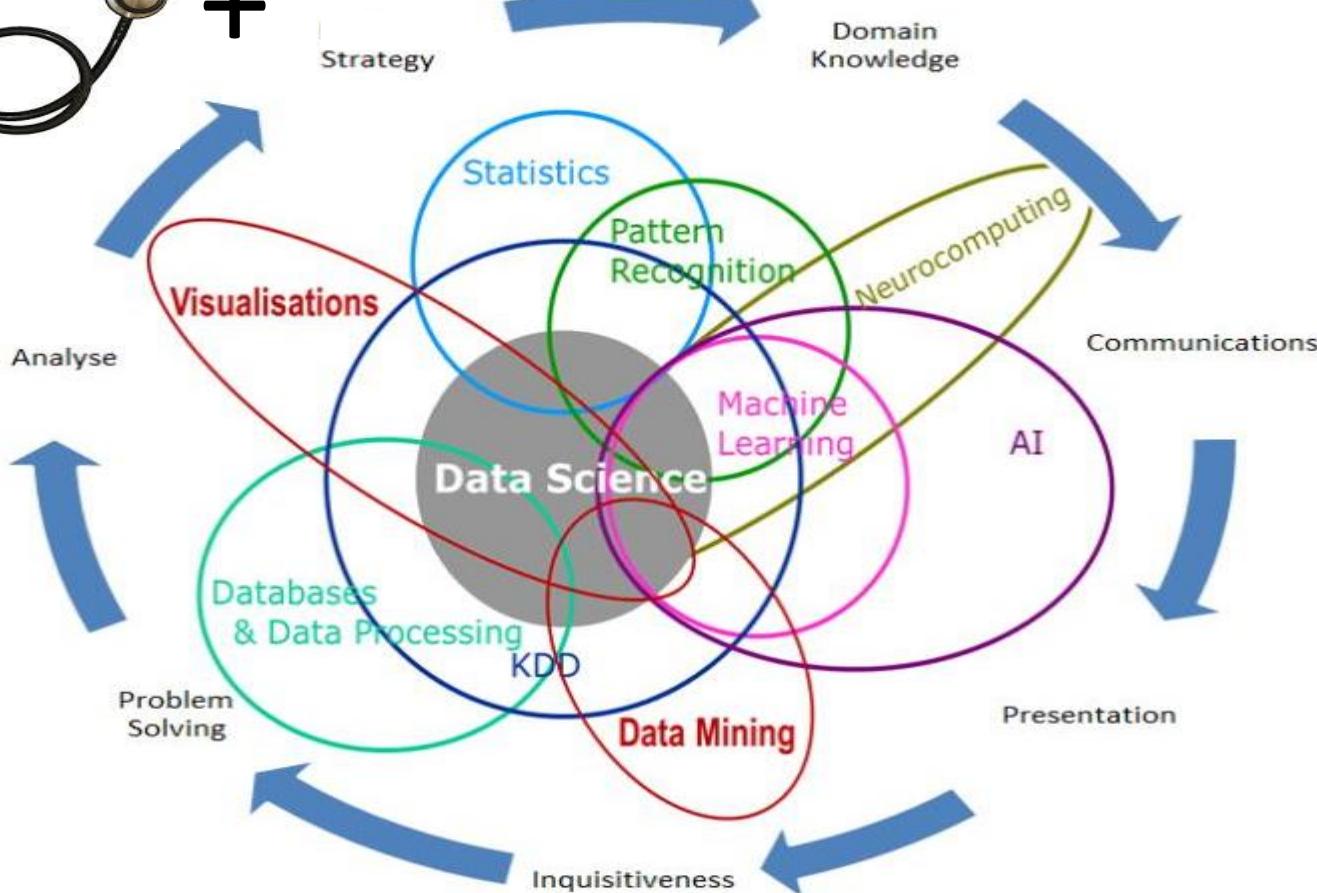
Non-linear models
data-expensive

Deep neural network

Cancer/patient selection

The tools you need

By Brendan Tierney, 2012



What is cool about deep learning?

You can play with formalisms & architectures

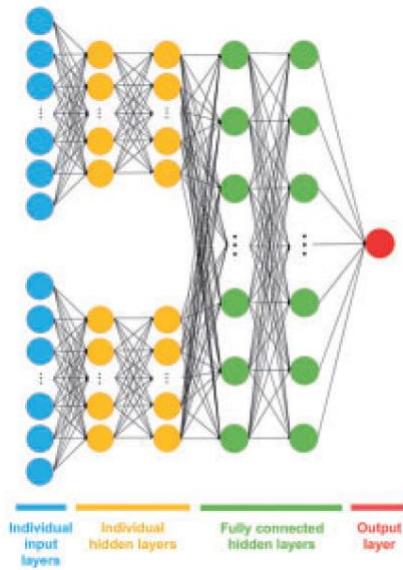
For what ?

To better fit your data and your questions!

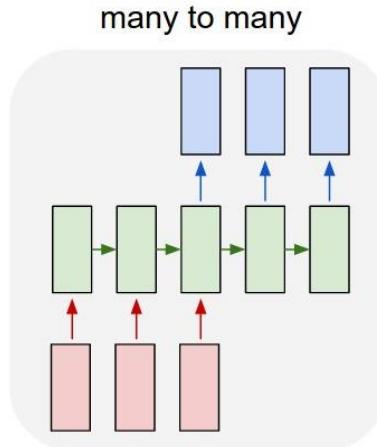
So be innovative!

Find your architecture

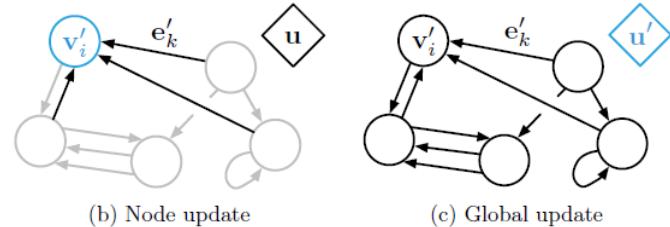
Drugs and proteins
to predict bioactivity
→ **Multitask NN**



Patients' folder
to predict patients' outcome
→ **Recurrent neural nets**

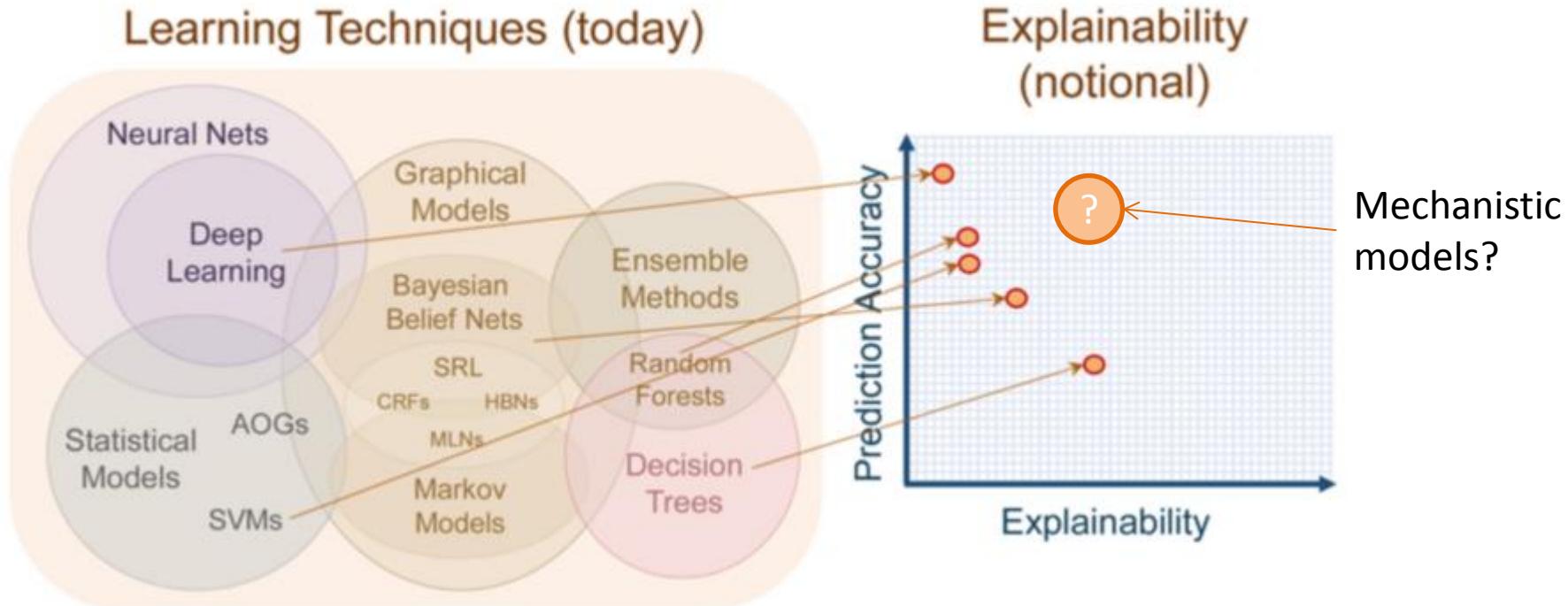


Molecular data
to predict biological behaviors
→ **Graphical networks**



Data science team at

Intepretability of ML/DL



Data

Where are the data?



A screenshot of the European Medicines Agency (EMA) website. The header includes the EMA logo and the text "EUROPEAN MEDICINES AGENCY SCIENCE MEDICINES HEALTH". The navigation menu at the top includes links for Home, Find medicine, Human regulatory, Veterinary regulatory, Committees, News & events, and Partners. A sidebar on the left shows categories like Overview, Research and development, Marketing authorisation (which is currently selected), and Advanced therapies. The main content area shows a breadcrumb trail: Home > Human regulatory > Marketing authorisation > Clinical data publication. Below this, there is a section titled "Clinical data publication" with a sub-section about publishing clinical data from October 2016.



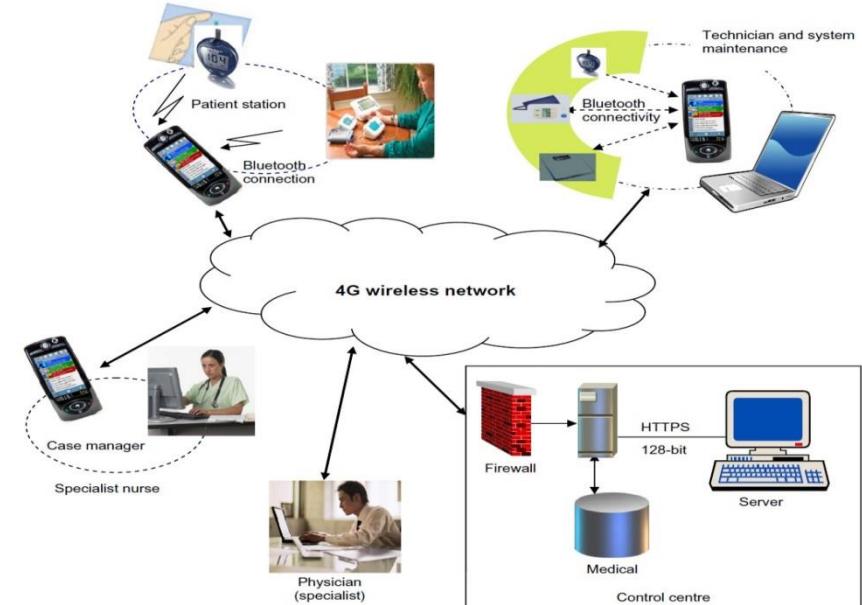
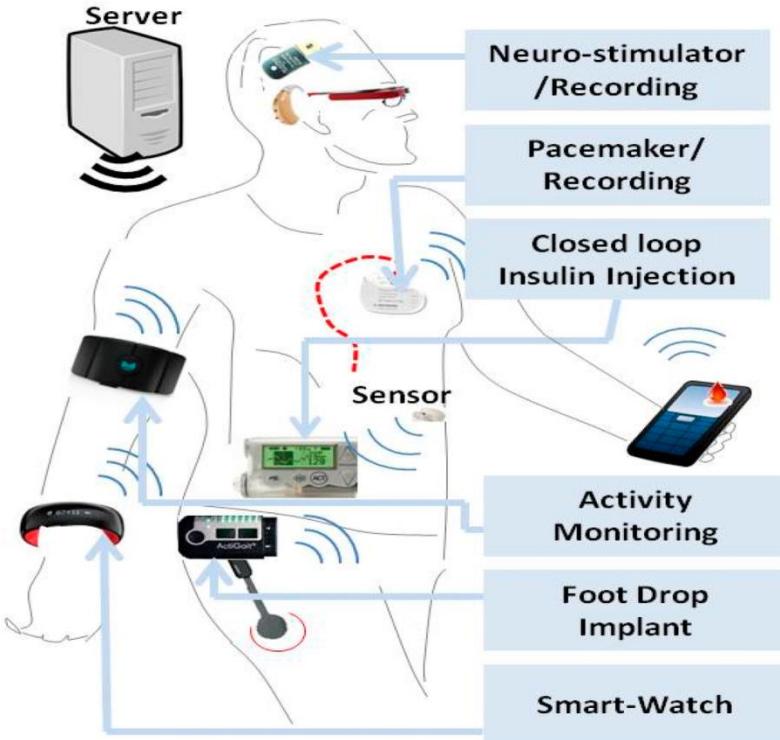
A screenshot of the ClinicalStudy DataRequest.com website. At the top, there is a search bar with the placeholder "View" and a "View" button. Below the search bar, there is a message about viewing studies from multiple sponsors. A grid of logos for various pharmaceutical companies is displayed, including Pfizer, Boehringer Ingelheim, GSK, Eli Lilly, Novartis, Roche, Sanofi, Takeda, UCB, and ViiV Healthcare. At the bottom of the page, there is a link to "The NEW ENGLAND JOURNAL of MEDICINE".

A screenshot of an article from The New England Journal of Medicine. The title of the article is "Data Sharing from Clinical Trials — A Research Funder's Perspective". The authors listed are Robert Kiley, Tony Peatfield, Jennifer Hansen, and Fiona Reddington. The article is categorized as a "SPECIAL ARTICLE".

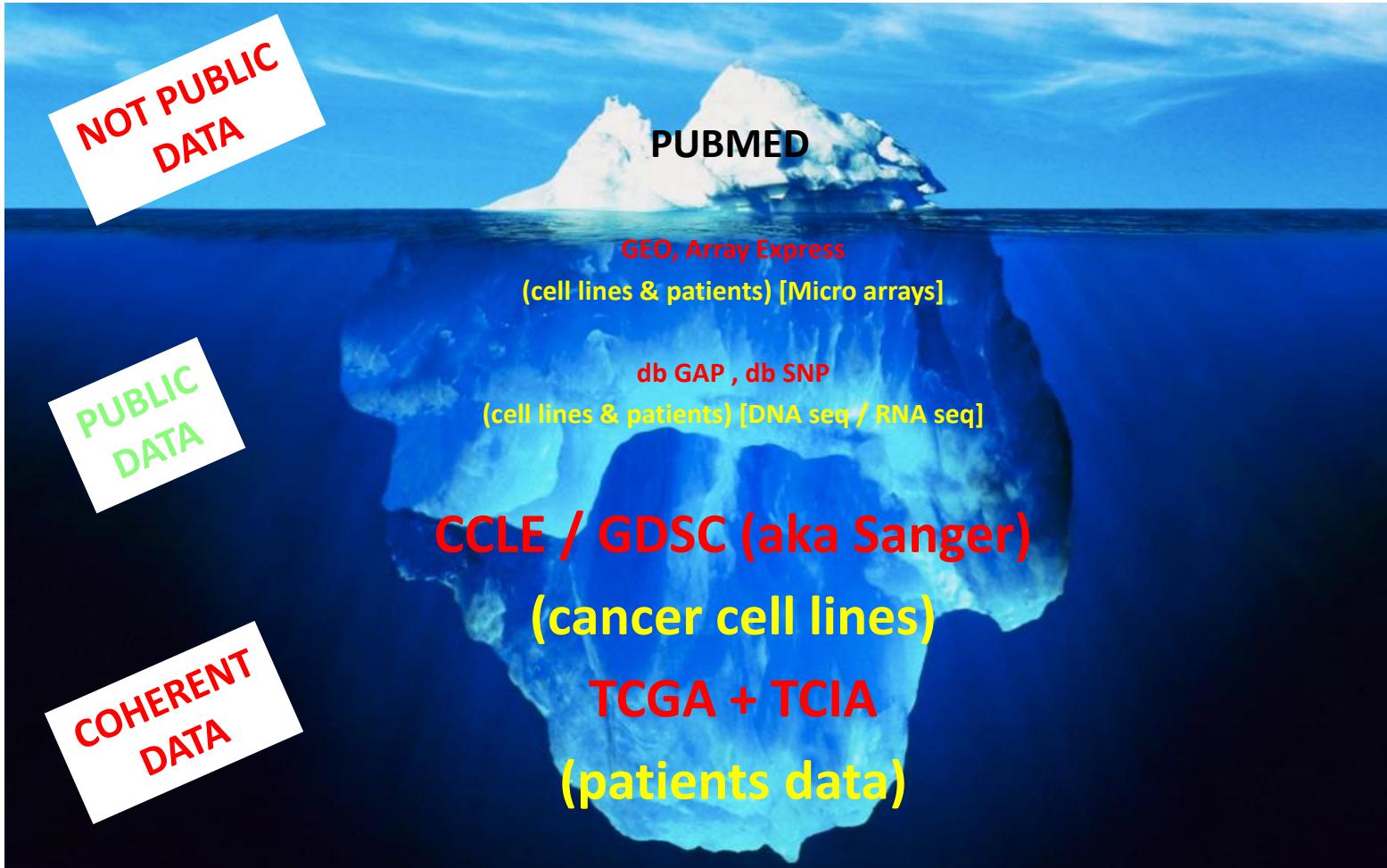
Clinical Trial Participants' Views
of the Risks and Benefits of Data Sharing

Michelle M. Mello, J.D., Ph.D., Van Lieou, B.S.,
and Steven N. Goodman, M.D., Ph.D.

New flow of clinical data



Where are the data?



Search is made easy

Traditional databases



A screenshot of the Human Protein Atlas website. At the top, it says "THE HUMAN PROTEIN ATLAS" with a small logo of three overlapping colored shapes (red, green, blue). Below that are "MENU", "HELP", and "NEWS". A search bar contains "PDX1" and "e.g. RBM3; insulin, CD36". Below the search bar are three thumbnail images: "TISSUE ATLAS", "CELL ATLAS", and "PATHOLOGY ATLAS".

Specialized search engine

A composite screenshot showing three specialized search engines. The top part shows the "Harmonizome" homepage with its logo and a search bar. The middle part shows the "OmicsDI" homepage with a search bar and a large word cloud of biological terms like "sequencing", "transcriptome", "methylation", etc. The bottom part shows a network visualization titled "OmicsDI" where nodes represent different tissues and organs, such as "Lung", "Blood", "Leaf", "Brain", "Liver", and "Others", connected by lines.

At Gustave Roussy

Non sécurisé | 31.10.2.38/drwh_dev/index.php

- Applications
- List of autoimmune...
- Cigogne - Platefor...
- Messagerie - Loic.V...
- OmicsDI: Home
- Harmonizome
- SAPHealth
- IntranetCurie
- cluster | Gustave Ro...
- Dr wh

- Onglet génétique
 - Onglet carte
 - Onglet clustering

Les documents autorisés :
 tout

Guide utilisateur

Télécharger le guide utilisateur en cliquant ici

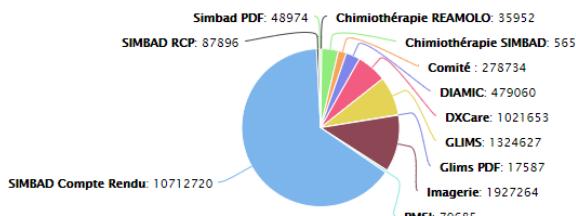
A citer pour une publication

"We recruited patients from this study using the data warehouse at XXX Hospital, Dr Warehouse (ref). It allows to search for patients from structured data (biology) and free text (hospital reports). It contains xxxx patients and x millions health reports."

Garcelon N, Neuraz A, Salomon R, Faour H, Benoit V, Delapalme A, Munnich A, Burgun A, Rance B. A clinician friendly data warehouse oriented toward narrative reports: Dr. Warehouse. J Biomed Inform. 2018 Apr;80:52-63. doi: 10.1016/j.jbi.2018.02.019. Epub 2018 Mar 1. PubMed PMID: 29501921 https://doi.org/10.1016/j.jbi.2018.02.019 pubmed.

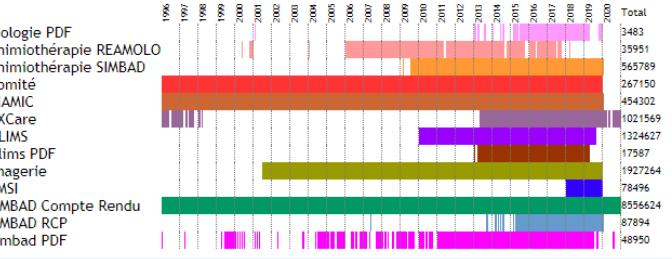
Derniers chargements

22/01/2020 DXCare	2915 documents chargés
22/01/2020 SIMBAD Compte Rendu	123 documents chargés
16/01/2020 Chimiothérapie SIMBAD	1176 documents chargés
08/01/2020 PMSI	5787 documents chargés
08/01/2020 COMITE	7128 documents chargés



Simbad PDF: 48974
 SIMBAD RCP: 87896
 SIMBAD Compte Rendu: 10712720
 Chimiothérapie REAMOLO: 35952
 Chimiothérapie SIMBAD: 565789
 Comité: 278734
 DIAMIC: 479060
 DXCare: 1021653
 GLIMS: 1324627
 Glims PDF: 17587
 Imagerie: 1927264
 PMSI: 79685

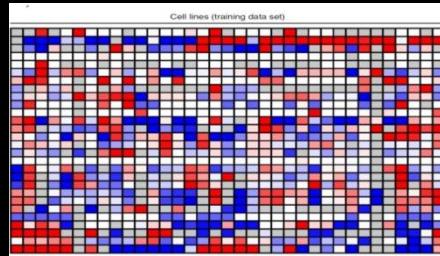
Les sources présentes



Source	1996	1997	1998	1999	2000	2001	2002	2003	2004	2005	2006	2007	2008	2009	2010	2011	2012	2013	2014	2015	2016	2017	2018	2019	2020	Total
Biologie PDF																									3483	
Chimiothérapie REAMOLO																									35951	
Chimiothérapie SIMBAD																									565789	
Comité																									267150	
DIAMIC																									454302	
DXCare																									1021569	
GLIMS																									1324627	
Glims PDF																									17587	
Imagerie																									1927264	
PMSI																									78496	
SIMBAD Compte Rendu																									8556624	
SIMBAD RCP																									87894	
Simbad PDF																									48950	

IESRM Bio Data nl view ↻

Molecular biology & pronostic



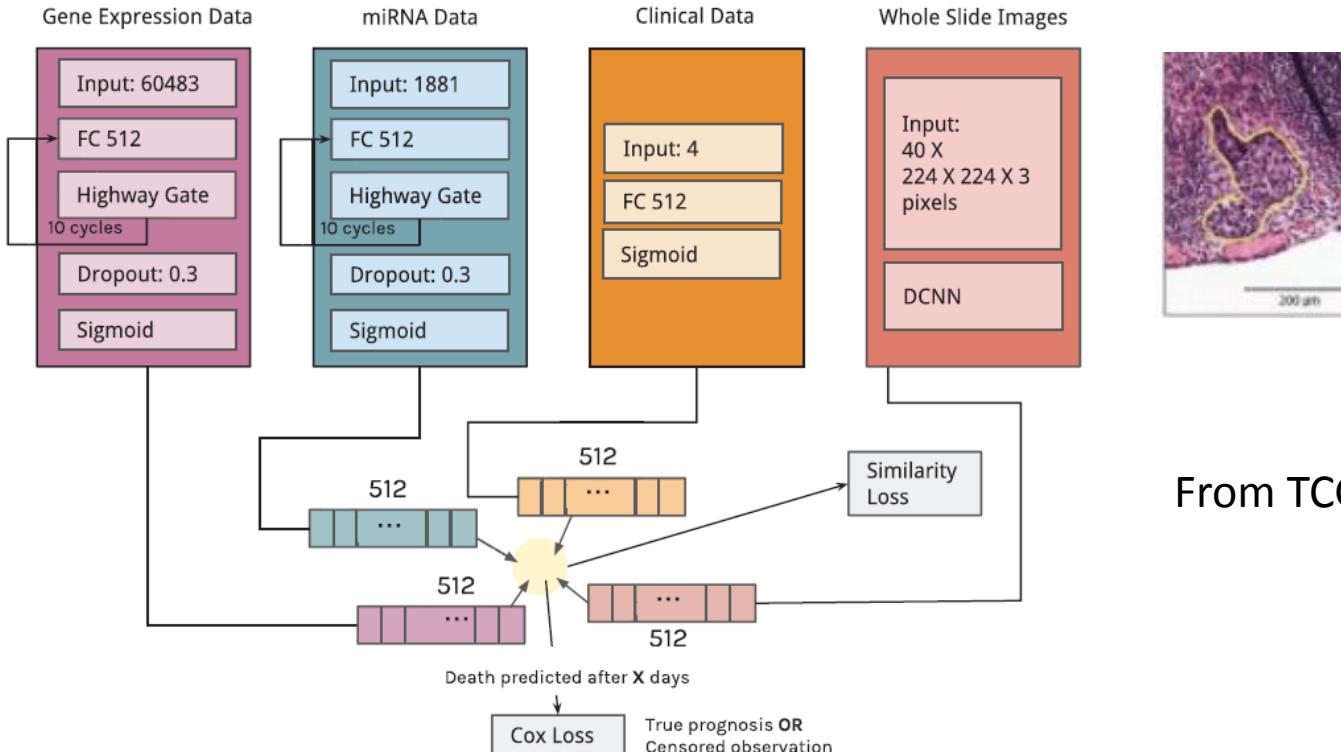
Does it work?

Pan-cancer survival estimation from RNAseq with Deep Learning

Main Objective	Tumor type, data types	Number of examples	Type of model	Endpoints and results	Reference
Patients' survival	4 cancers TCGA RNAseq	953 samples	Survival random forest and Cox lasso	Cross-validation C-index ~ 0.60 - 0.75	Yuan Y, et al. Assessing the clinical utility of cancer genomic and proteomic data across tumor types. Nat Biotechnol. 2014
Patients' survival	10 TCGA organs related datasets RNA-Seq expression	5031 patient samples (train/test split 80/20%)	MLP with Cox loss	Test C-index ~ 0.62	Ching T, Zhu X, Garmire LX. Cox-nnet: An artificial neural network method for prognosis prediction of high-throughput omics data. PLoS Comput Biol. 2018 Apr
Patients' survival per groups	32 cancer types TCGA RNA-Seq expression	6645 patients	Survival	Validation C-index = 0.59-0.75	Ramazzotti D, et al. Multi-omic tumor data reveal diversity of molecular mechanisms that correlate with survival. Nat Commun. 2018 Oct
Patients' survival	20 cancer types TCGA RNA-Seq expression	6404 patients	CNN	Clinic + mRNA, test C-index = 0.60	<i>Deep learning with multimodal representation for pan-cancer prognosis prediction, Anika Cheerla and Olivier Gevaert, Bioinformatics, 35, 2019</i>

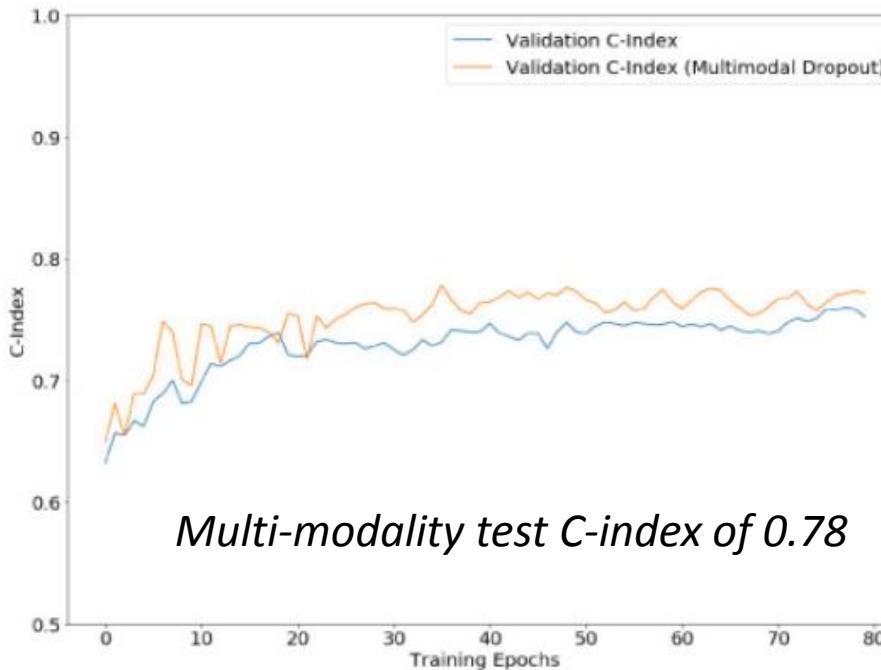
Pronostic estimation the clinic

Best and most recent try

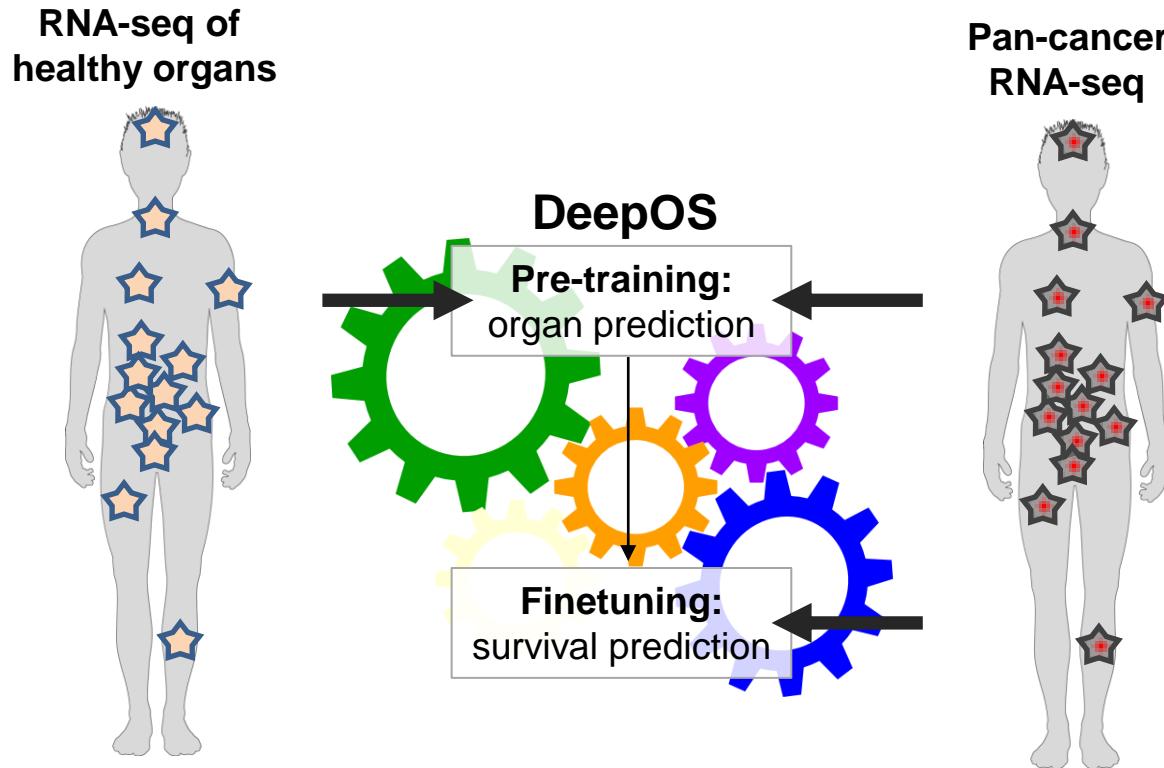


Pronostic estimation the clinic

Best and most recent try

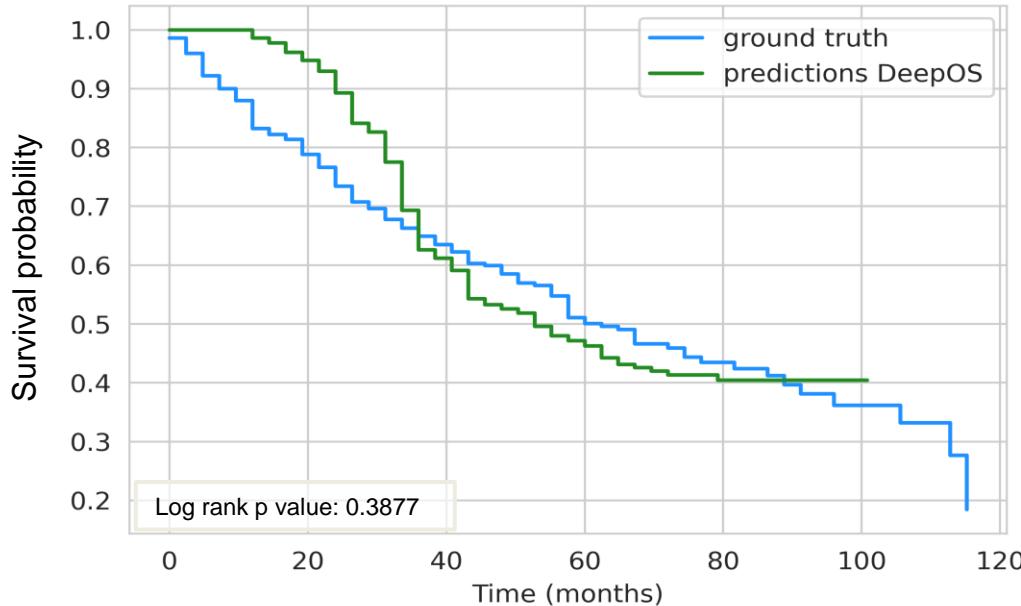


DeepOS



DeepOS

Predicted and true Kaplan-Meier curves on the test cohort



Test set C-index 0.7147

How to improve models' efficacy?

- Need more samples (cell lines, organoids)
- Need more molecular layers ?
- Need better readouts
- Need more functional (perturbation) data
- Adapt your models to the tasks you want

→ Saez-Rodriguez's take

How to improve models' efficacy?

- Focus on your data-analytical approaches
- Use mechanistic prior knowledge
- Adapt your models to the tasks you want

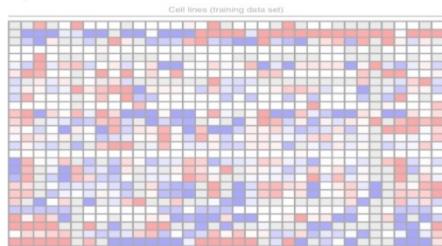
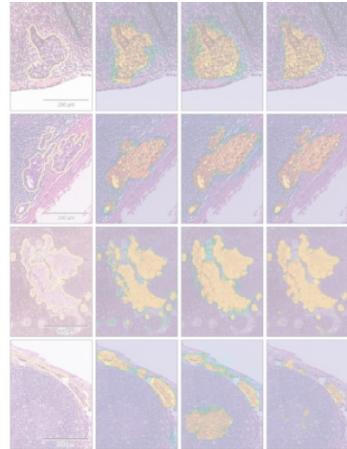
→ Camacho, Costello, Collins takes

→ & my take!

For patients' prognostic estimation

Models are improving...

ML/DL in oncology



Diagnostics

Molecular biology
& pronostic



"Here's my sequence"
The New Yorker

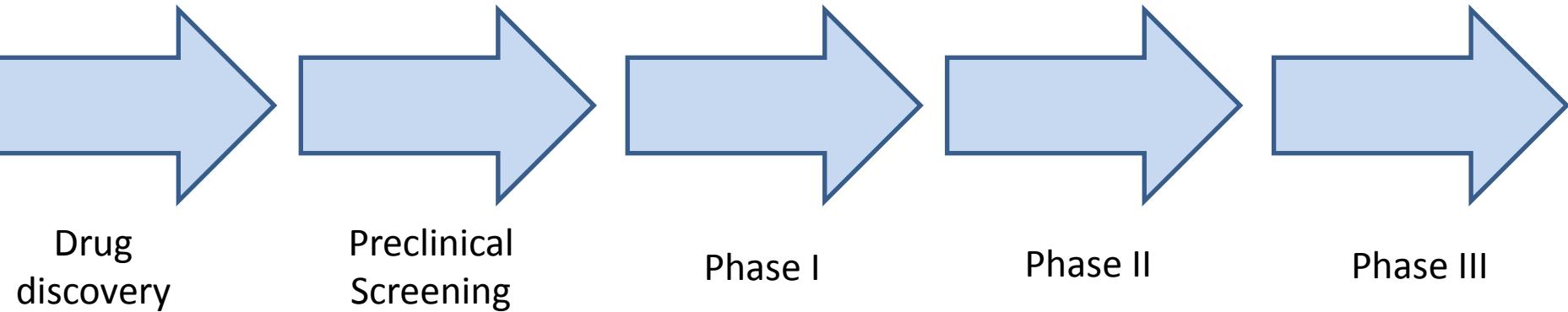
Drug development
& prediction



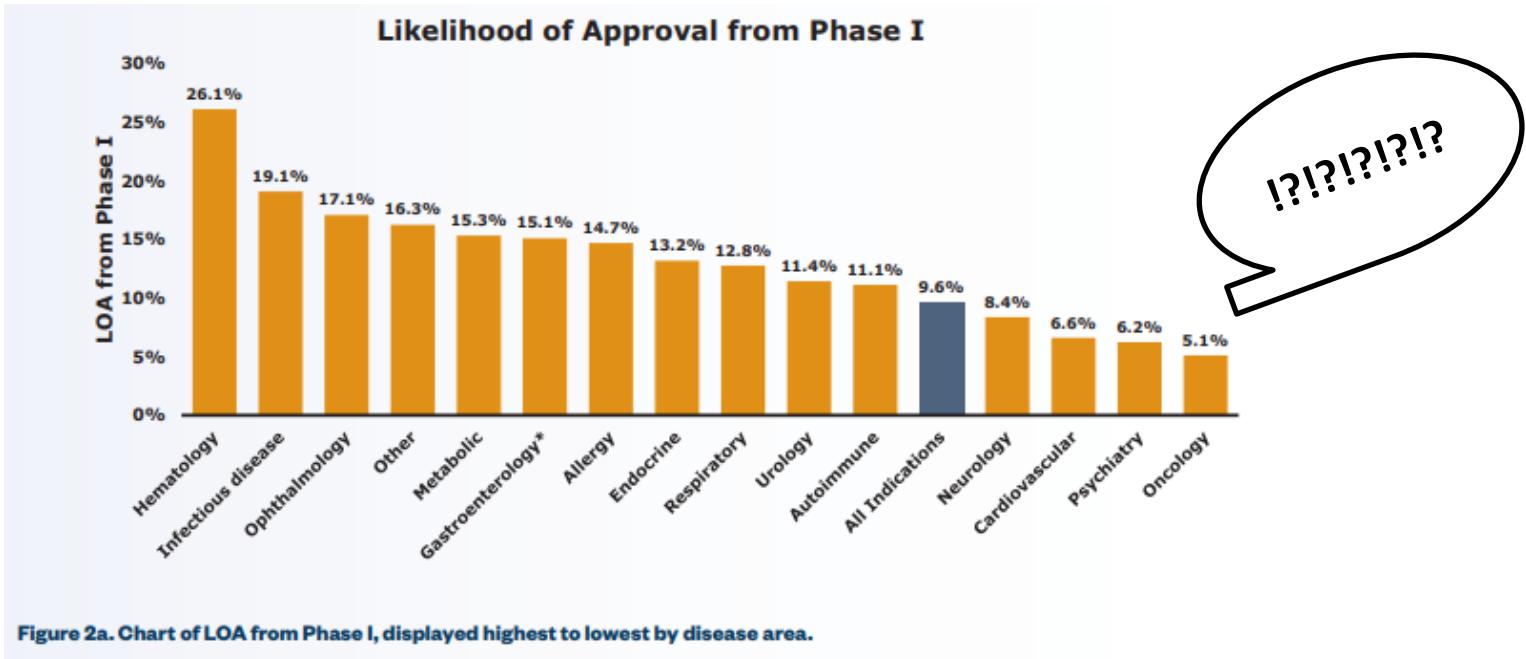
Monitoring patients

PLAN

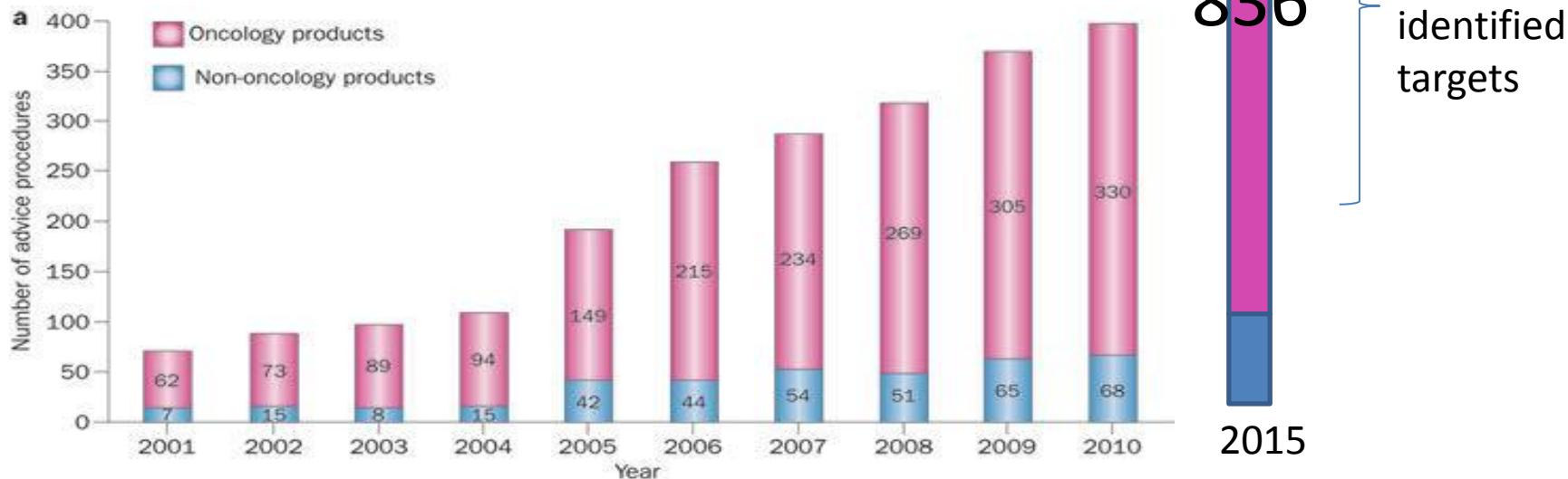
DITEP 
Drug Development Department



Drug development success rate 2006-2015



Drugs' production...

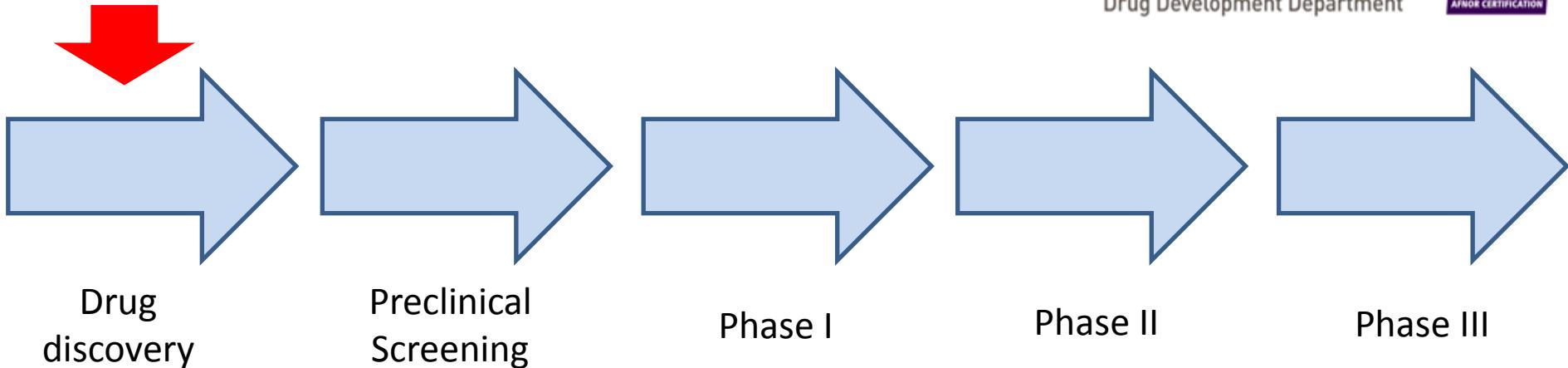


Jonsson B, Bergh J. Hurdles in anticancer drug development from a regulatory perspective. Nat Rev Clin Oncol. 2012

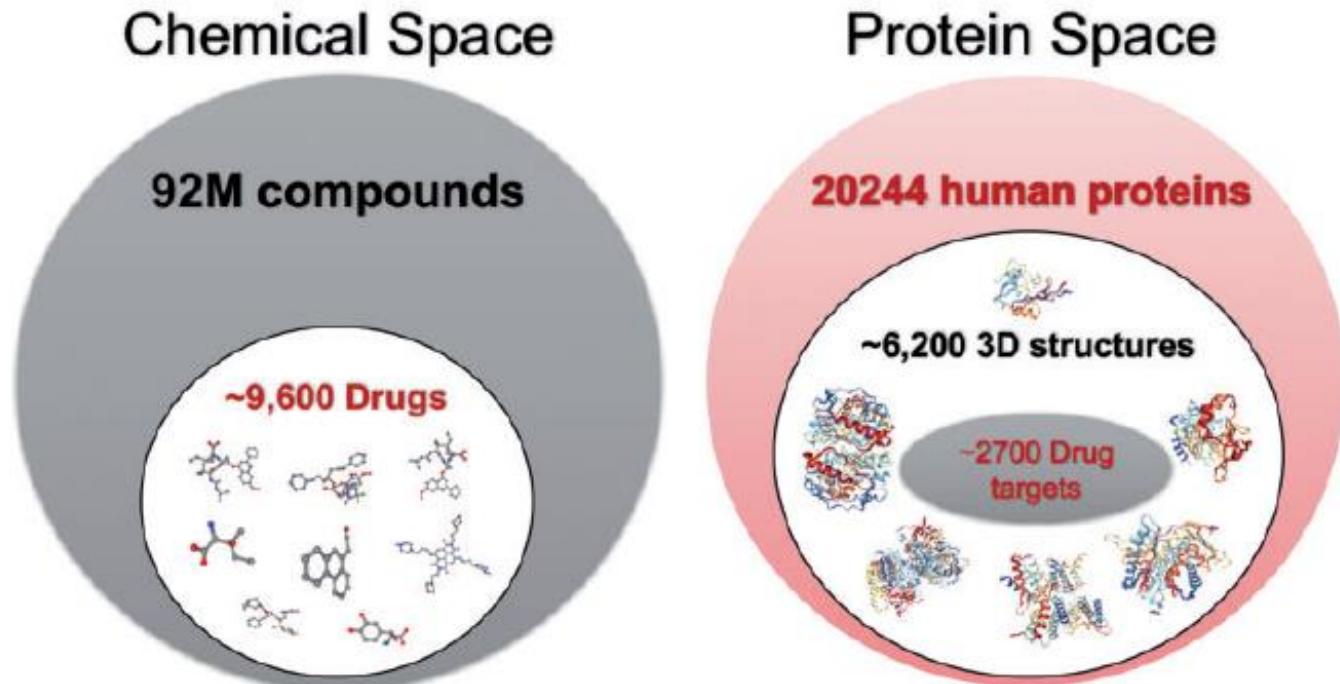
You are here

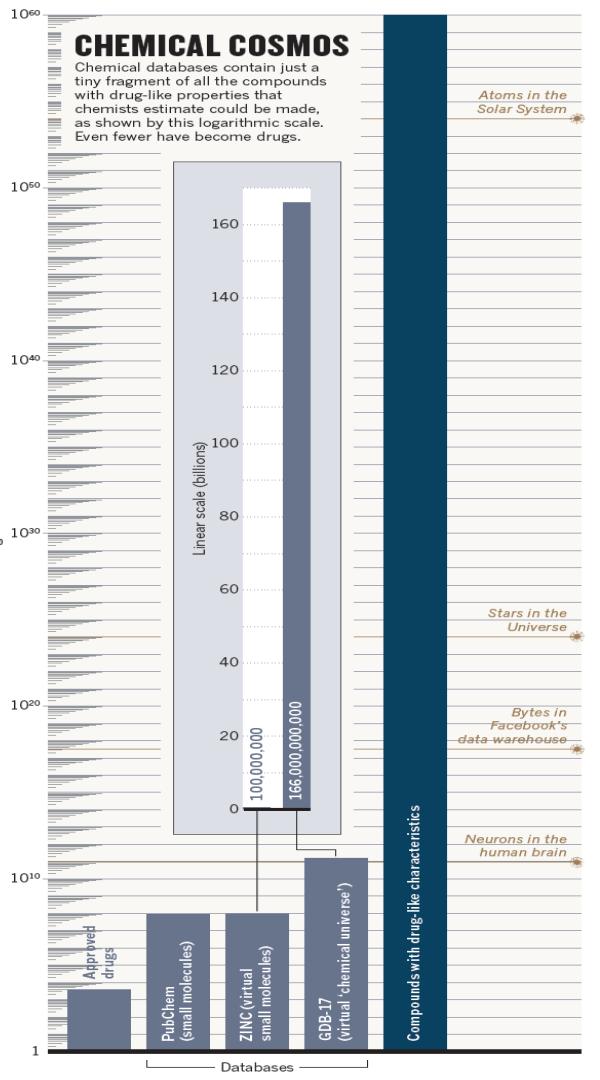
Modèles / ML / IA ? DITEP

Drug Development Department



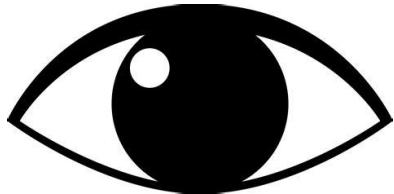
The known space





The models you have

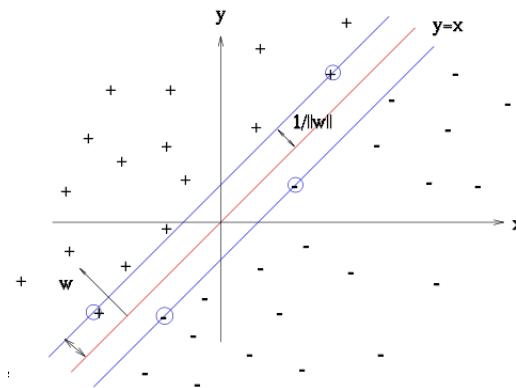
Human



Use your knowledge
and reasoning

Examples: Christophe Massard

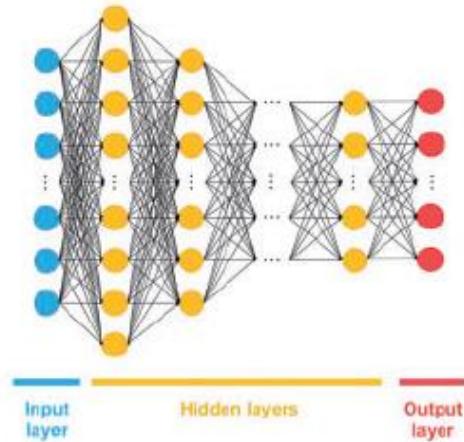
Statistics



Linear models easy to
explain

Cox model
logReg eg: RMH

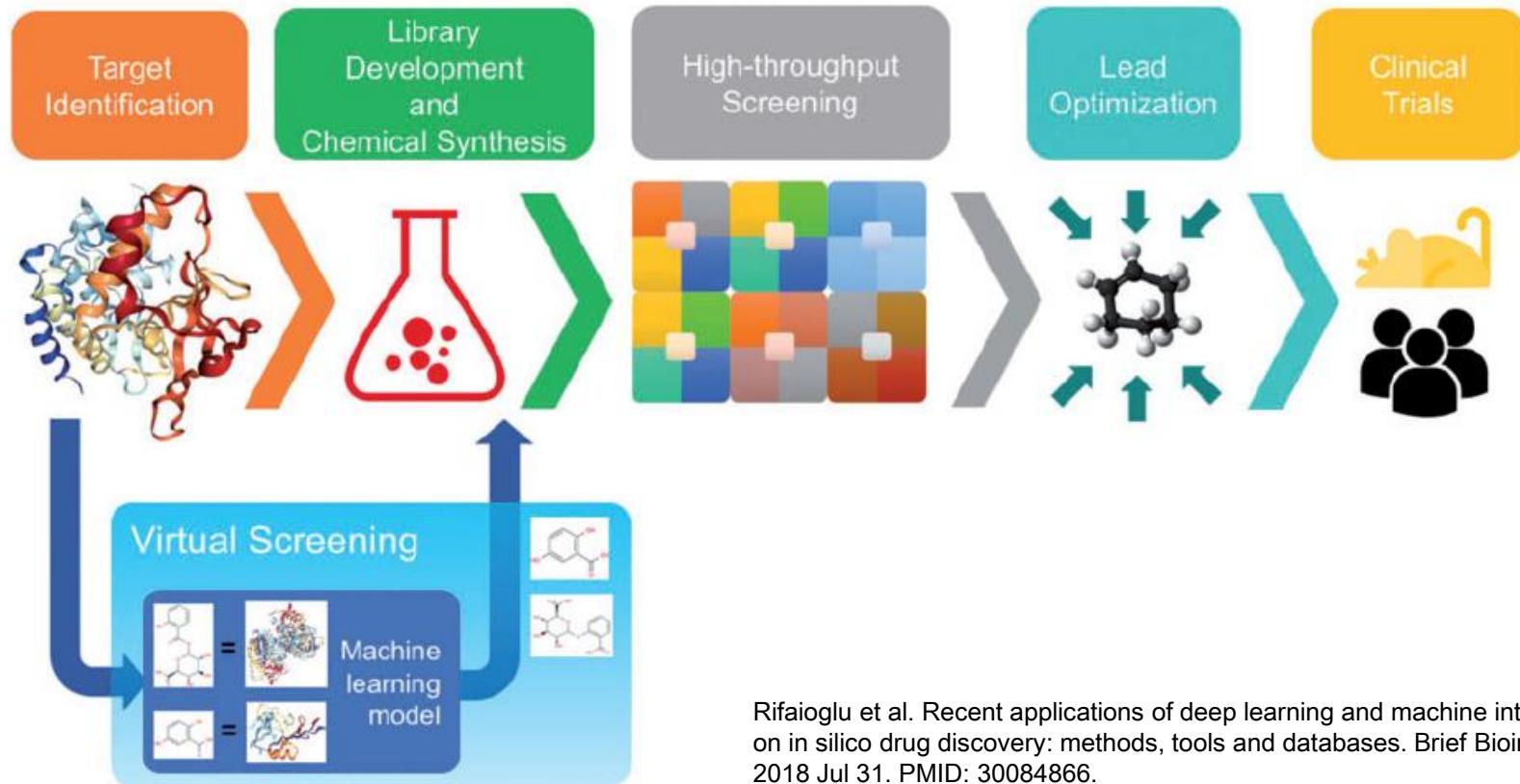
Deep learning



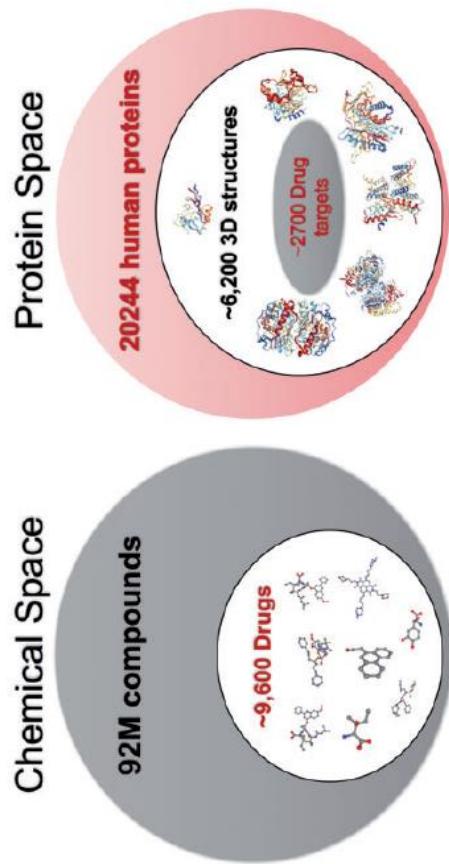
Non-linear models
data-expensive

Deep neural network

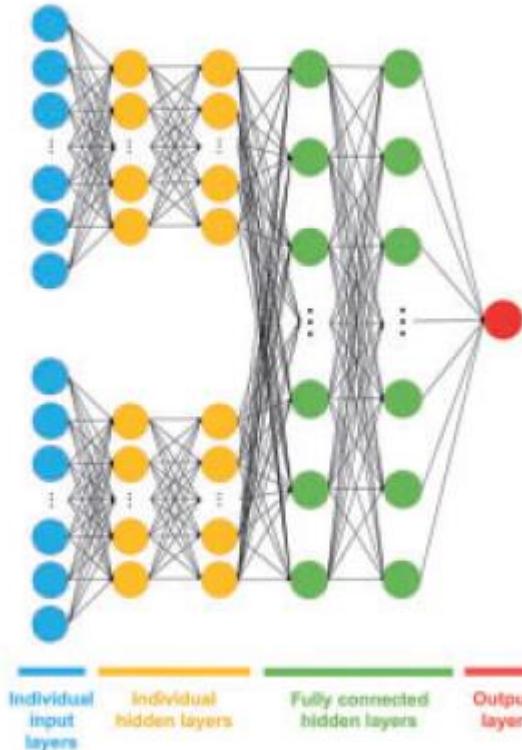
Drug discovery in details



Overview



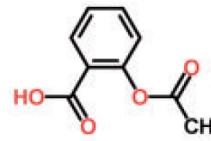
Pairwise Input Neural Network (PINN)



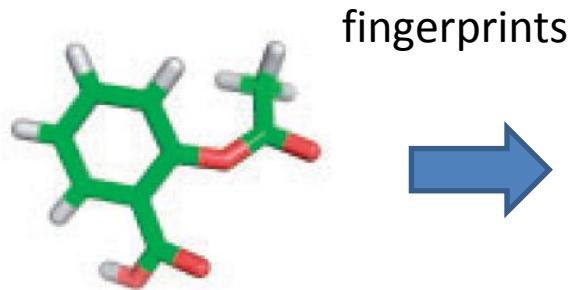
Bioactivity

How to encode drugs

Table 1. Chemical formula, 2D/3D graphical representation, SMILES and InChI notations of aspirin

Category	Representation
Compound name	Aspirin
Chemical formula	C ₉ H ₈ O ₄
3D/2D structure	 
SMILES	CC(=O)OC1=CC=CC=C1C(=O)O
InChI	InChI = 1S/C9H8O4/c1-6(10)13-8-5-3-2-4-7 (8)9(11)12/h2-5H, 1H3, (H, 11, 12)

How to encode drugs or proteins?



fingerprints

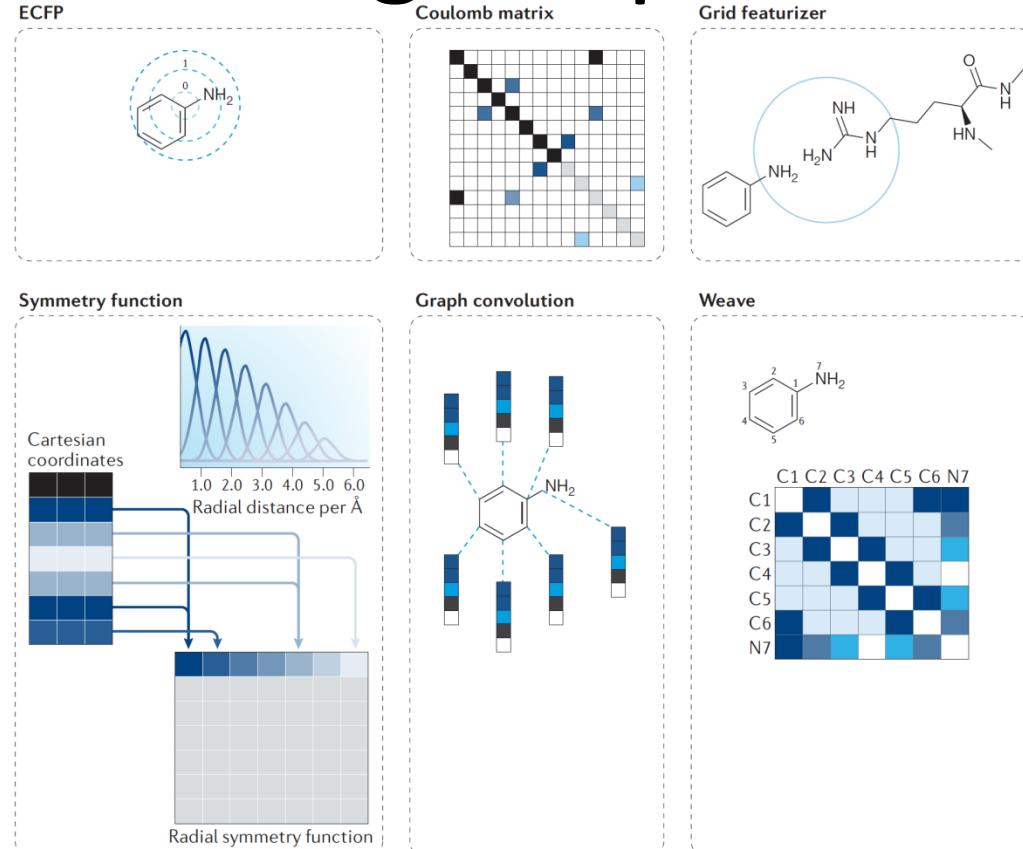
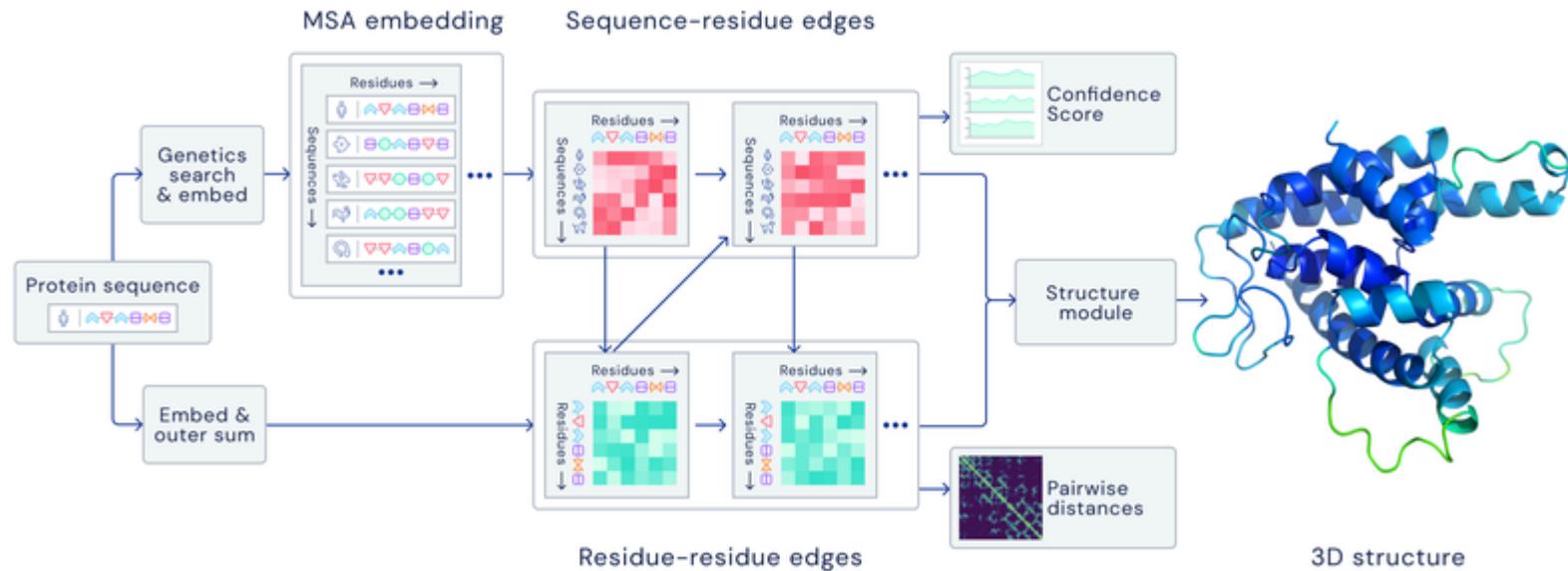


Fig. 3 | The challenges of compound structure representation in machine learning models. The appropriate

Protein folding prediction

Prediction of ‘spatial graph’: eg backbone atom distance matrices and torsion angles



AlphaFold’s architecture in schematic form. Image Credit: DeepMind

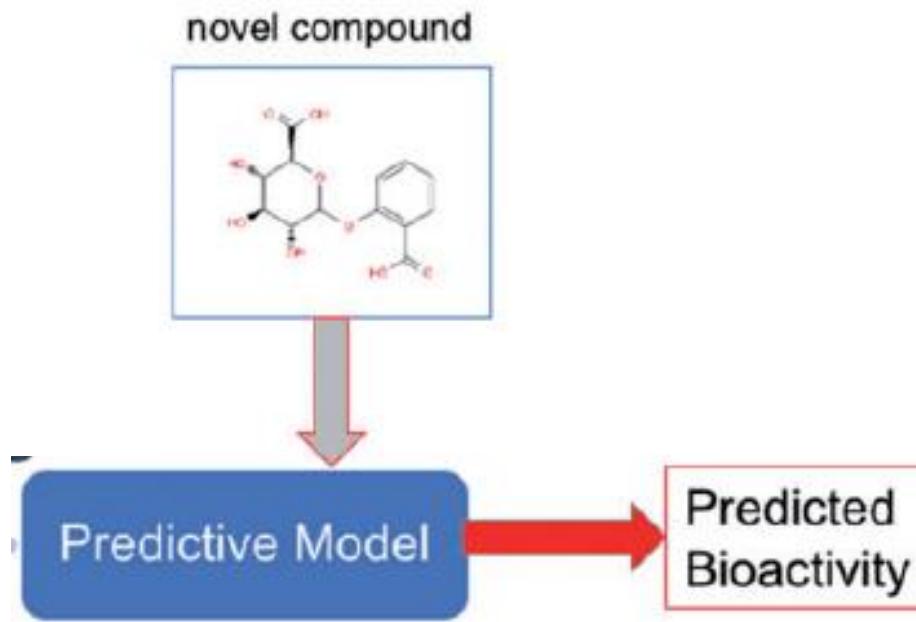
other: arXiv:1911.05531 [q-bio.BM]

Gold standard databases

- **DUD-E**
 - 22 886 ligands and their affinities against 102 targets retrieved from the ChEMBL database
- **TOX-21**
 - 12 000 compounds screened for their toxic effects (in bioassays)
- **MoleculeNet**
 - ~ 700 000 compounds retrieved from publicly available databases

Compound and bioactivity databases	Statistics ^a			Website	Version
	Compounds	Targets	Interactions		
PubChem [1]	93 977 773 (C) 235 653 (S)	10 341 (P)	233 799 255 (I) 1 252 820 (E)	https://pubchem.ncbi.nlm.nih.gov	03.12.2017
ChEMBL [2]	1 735 442 (C)	11 538 (P)	14 675 320 (I) 1 302 147 (E)	https://www.ebi.ac.uk/chembl	v23
DrugBank [5]	9591 (D)	4270 (P)	16 748 (I)	http://www.drugbank.ca	v5.0

Predictions



Deep learning wins

2 studies on multitask NN for bioactivity measures
~ hidden layer sizes: [1000, 2000, 3000]

- 743 336 compounds & 5069 targets
- 2.1 million bioactivity measurements
- Compounds were represented by 43 x 340 dimensions ECFP12 fingerprints

Method	AUC	p-value
Deep network	0.830	
SVM	0.816	1.0e-07
BKD	0.803	1.9e-67
Logistic Regression	0.796	6.0e-53
k-NN	0.775	2.5e-142
Pipeline Pilot Bayesian Classifier	0.755	5.4e-116
Parzen-Rosenblatt	0.730	1.8e-153
SEA	0.699	1.8e-173

J&J + univ Austria 2014

Unterthiner T, Mayr A, Klambauer G, et al. Deep learning as an opportunity in virtual screening. Deep learn Represent Learn Work NIPS 2014;2014:1–9.

- 1.6 million compounds & 259 targets
- 37.8 million experimental compound–protein interactions
- Compounds were represented by ECFP4 fingerprints

Model	PCBA (n = 128)	MUV (n = 17)	Tox21 (n = 12)	Sign Test CI
Logistic Regression (LR)	.801	.752	.738	[.04, .13]
Random Forest (RF)	.800	.774	.790	[.06, .16]
Single-Task Neural Net (STNN)	.795	.732	.714	[.04, .12]
Pyramidal (2000, 100) STNN (PSTNN)	.809	.745	.740	[.06, .16]
Max {LR, RF, STNN, PSTNN}	.824	.781	.790	[.12, .24]
1-Hidden (1200) Layer Multitask Neural Net (MTNN)	.842	.797	.785	[.08, .18]
Pyramidal (2000, 100) Multitask Neural Net (PMTNN)	.873	.841	.818	

Google + Standford 2015

Ramsundar B, Kearnes S, Edu K, et al. Massively Multitask Networks for Drug Discovery. arXiv 2015; arXiv: 1502.02072.

You are here



Drug
discovery

Preclinical
Screening

Phase I

Phase II

Phase III

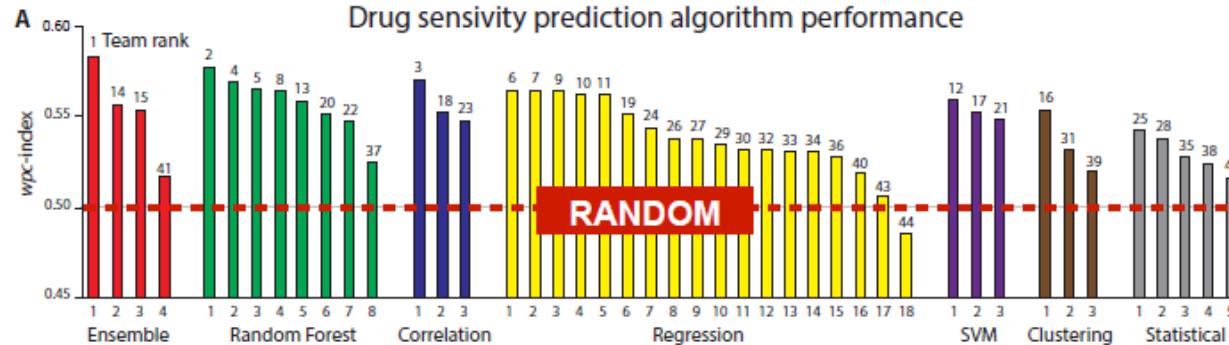
DITEP

Drug Development Department



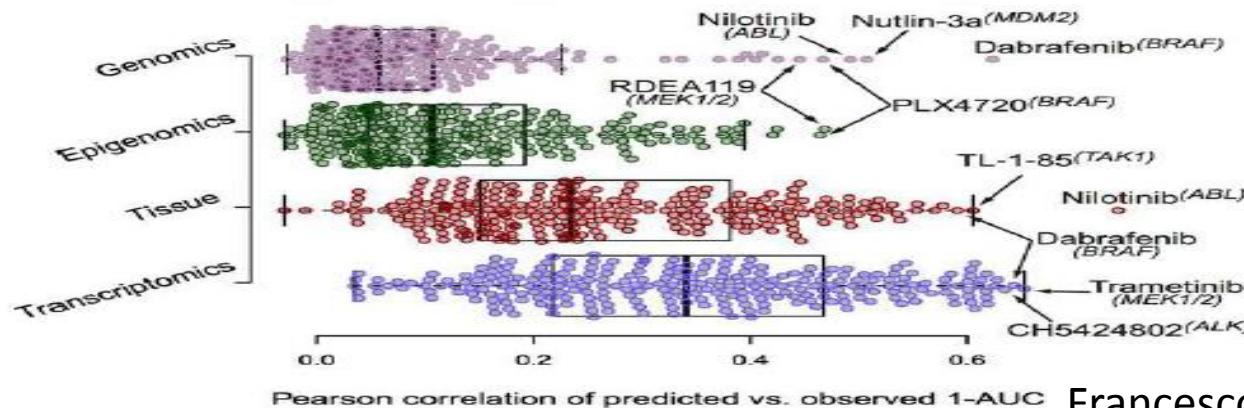
Treatment prediction in cell lines

Saez-Rodriguez « predictability is very low »



DREAM CHALLENGES 

Costello JC, et al. A community effort to assess and improve drug sensitivity prediction algorithms. Nat Biotechnol. 2014

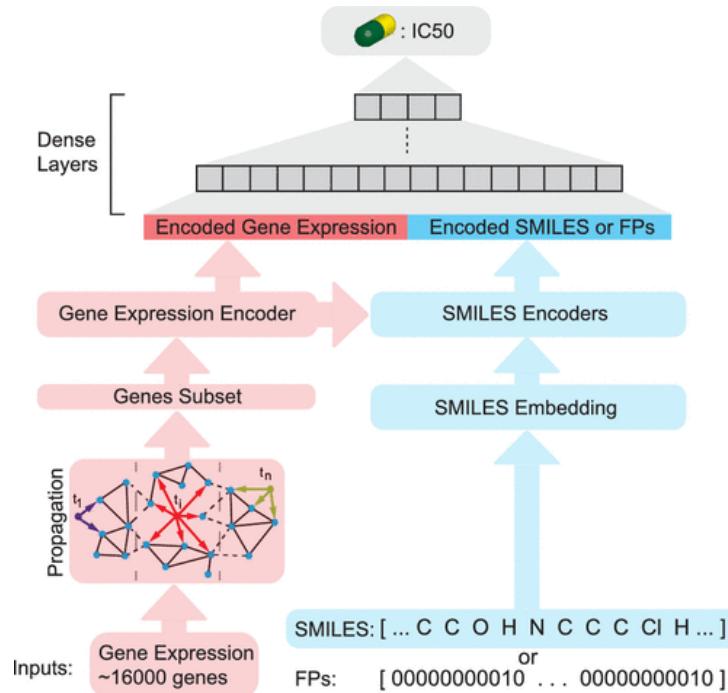


Francesco Iorio et al, Cell, 2016

Treatment prediction in cell lines

Best and most recent try

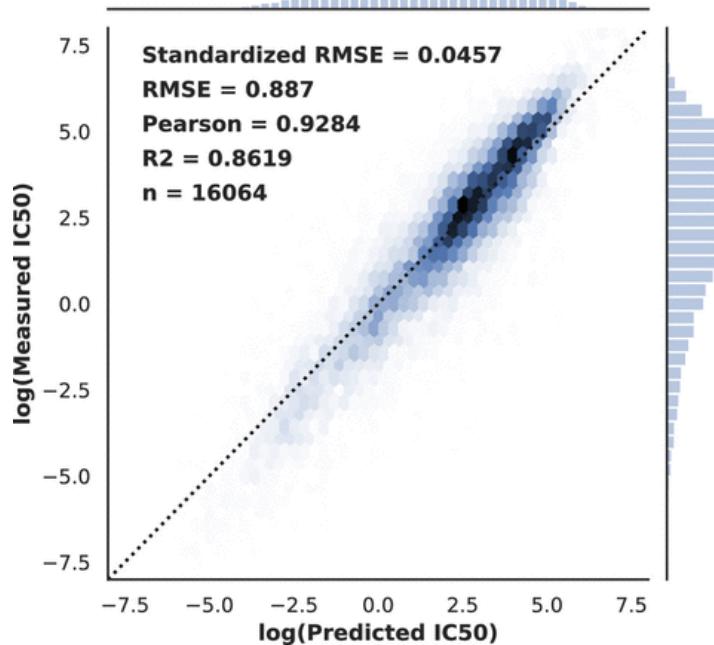
Data:
985 cell lines
208 drugs
= 175 603 pairs



Treatment prediction in cell lines

Best and most recent try

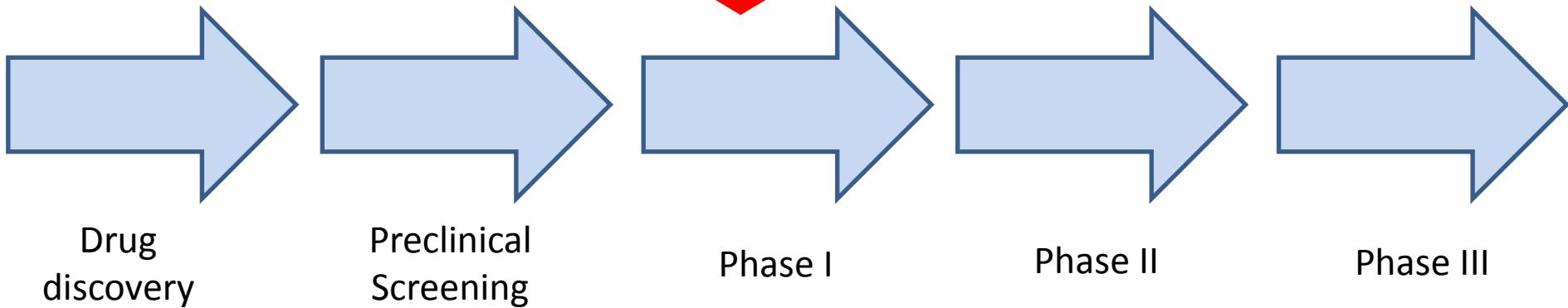
Test set



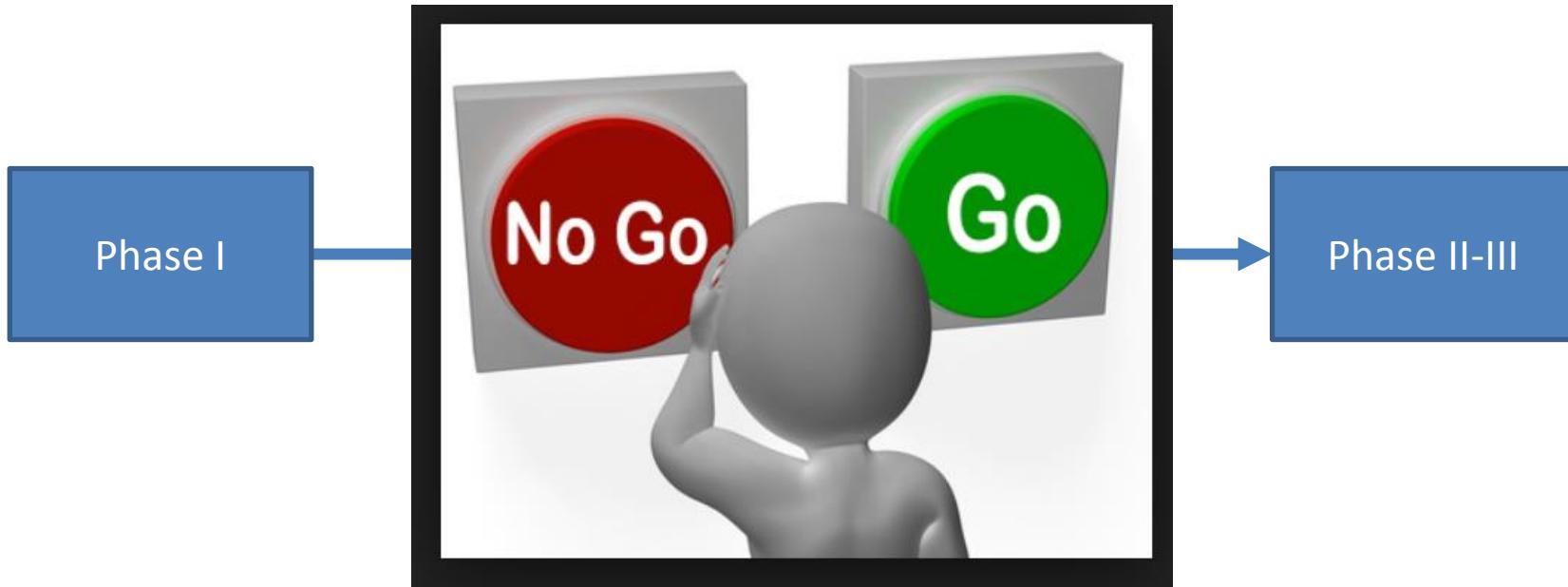
You are here

DITEP

Drug Development Department

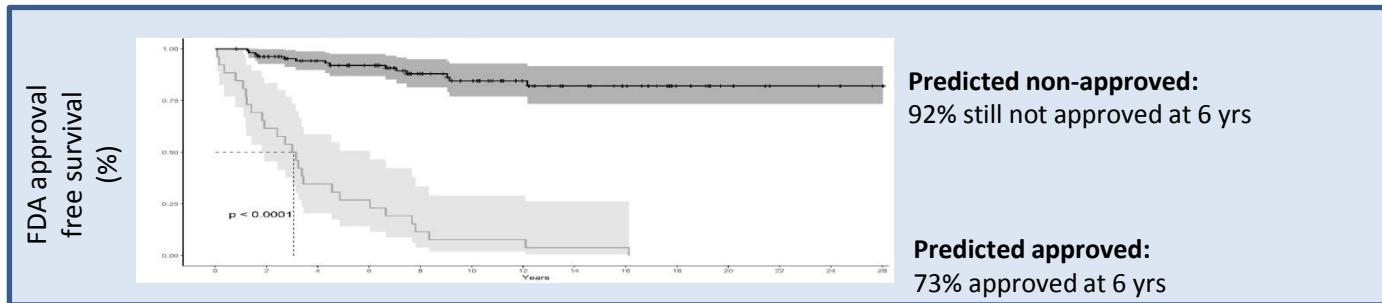
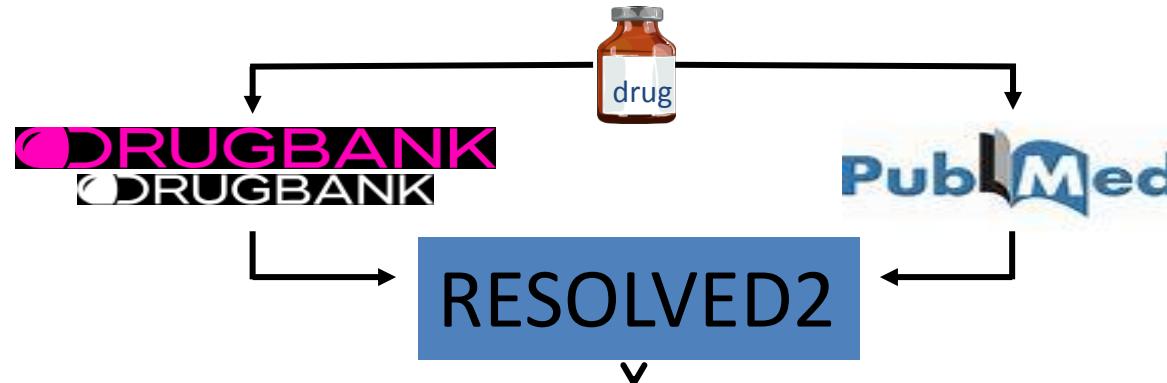


Les enjeux à la fin de la phase I: la décision



Prediction of Drug Approval After Phase I Clinical Trials in Oncology: The RESOLVED2 model

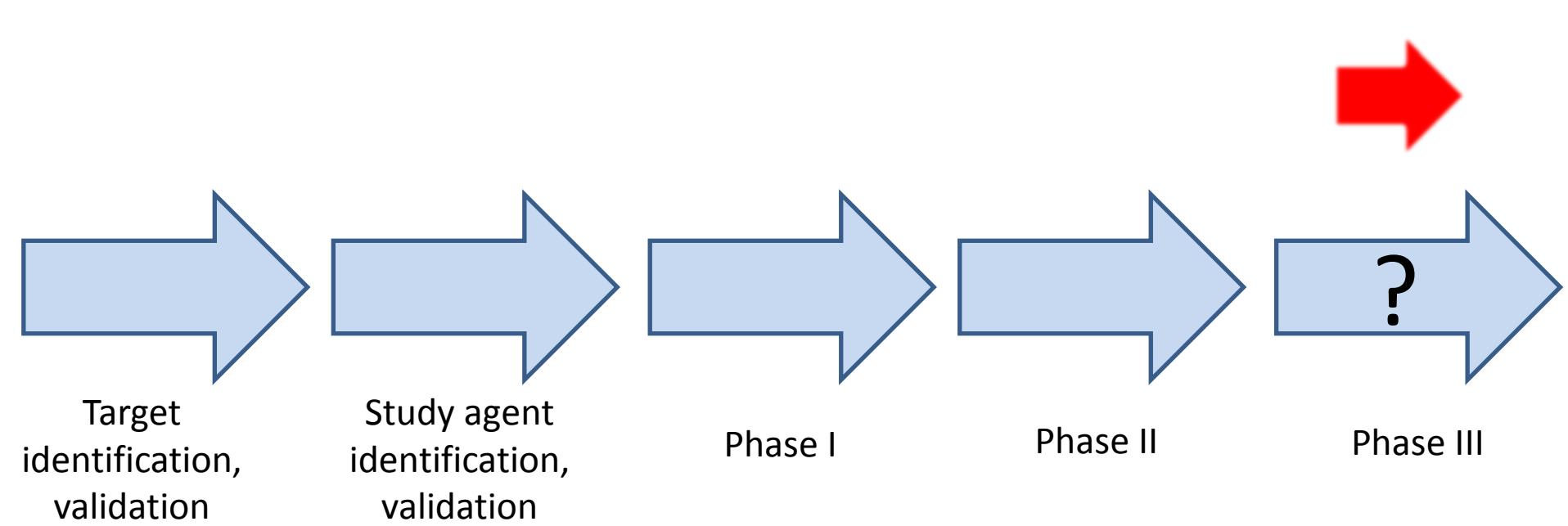
Is this new drug likely to be approved in the future?



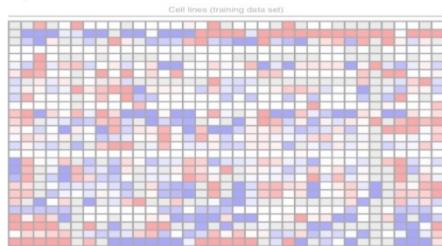
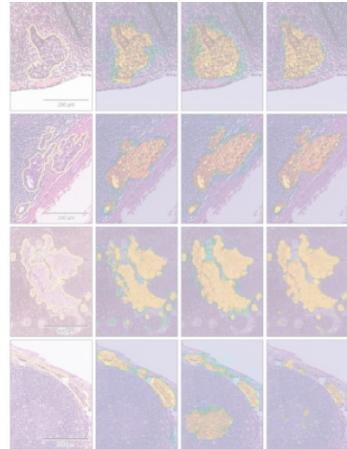
Under submission

DITEP 
Drug Development Department





ML/DL in oncology



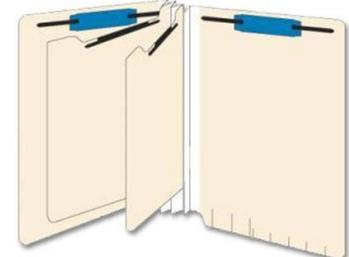
Diagnostics

Molecular biology
& pronostic



"Here's my sequence"
The New Yorker

Drug development
& prediction



Free text &
Monitoring patients

Tasks with EHR

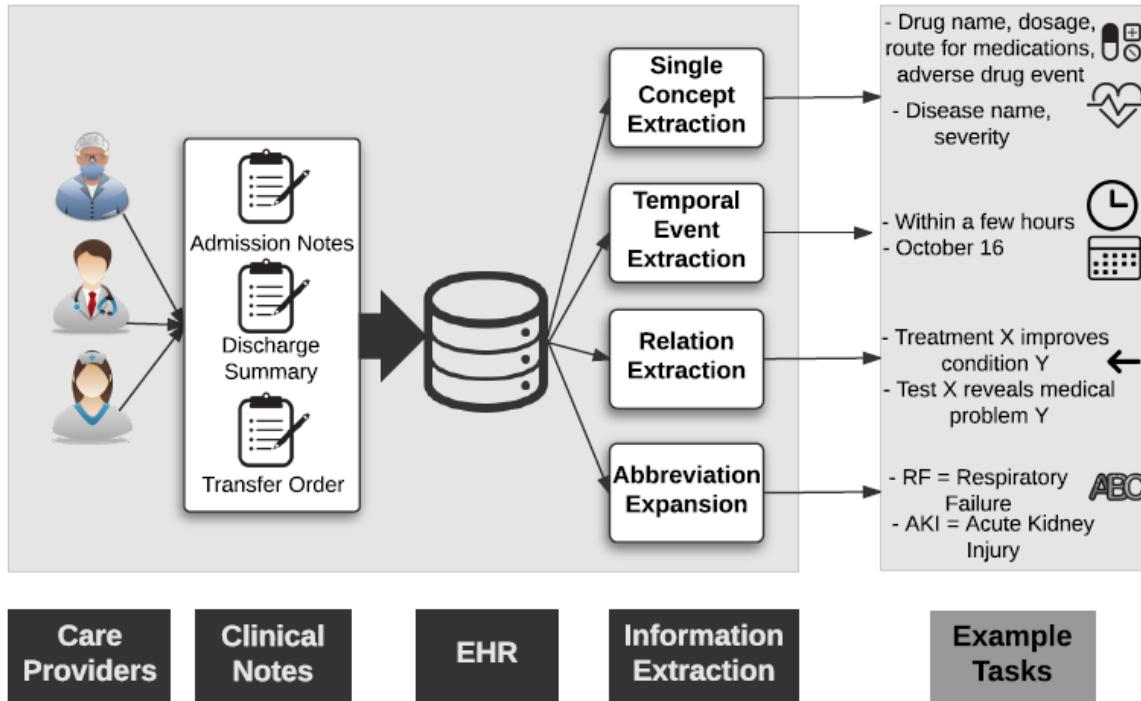
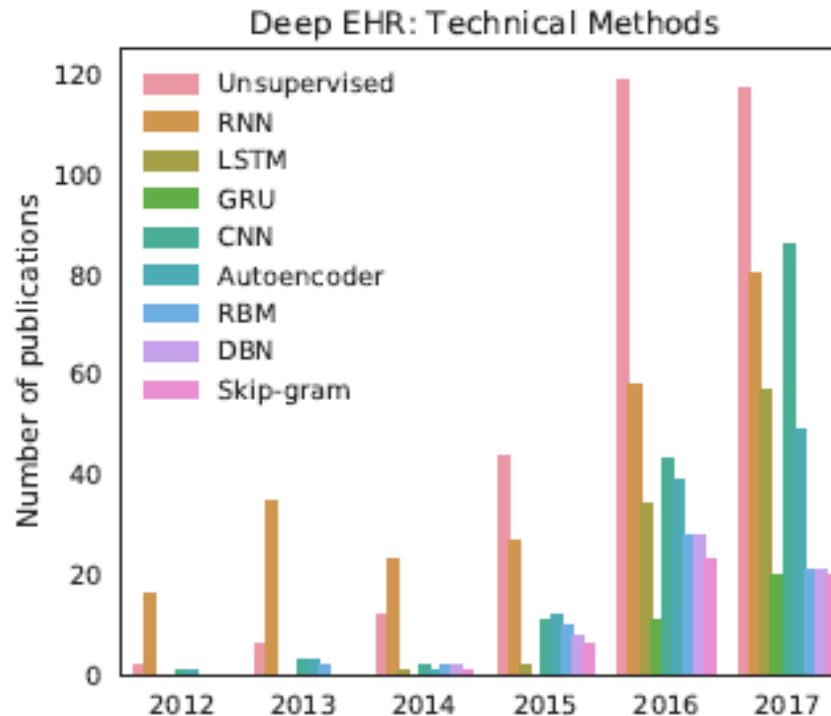


Fig. 7. EHR Information Extraction (IE) and example tasks.

Deep EHR: A Survey of Recent Advances in Deep Learning Techniques for Electronic Health Record (EHR) Analysis.
Benjamin Shickel, Patrick J. Tighe, Azra Bihorac, and Parisa Rashidi. arXiv:1706.03446v2

Methods for EHR



Taux d'attrition des patients



Consultation d'inclusion	Phase de screening	C1J1	Période de DLT	Sortie d'essai
100%	85% - 75% - (15%, 25%)	~65%	0%	
	Screen fails	Mortalité		

- Mckane A, et al: Determinants of patient screen failures in Phase 1 clinical trials. Invest New Drugs
- Kempf E et al: A Case-Control Study Brings to Light the Causes of Screen Failures in Phase 1 Cancer Clinical Trials. PLoS ONE

Olmos D, A'hern RP, Marsoni S, et al: Patient selection for oncology phase I trials: a multi-institutional study of prognostic factors. J Clin Oncol

Softwares : prediction of successful screening and DLT period completion (SSD)

En cours de développement

DITEP Data Science Team

DTNSI



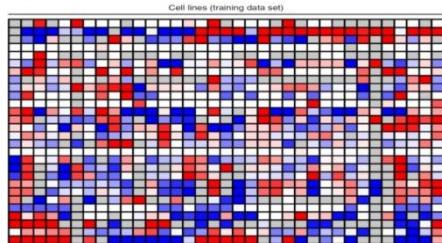
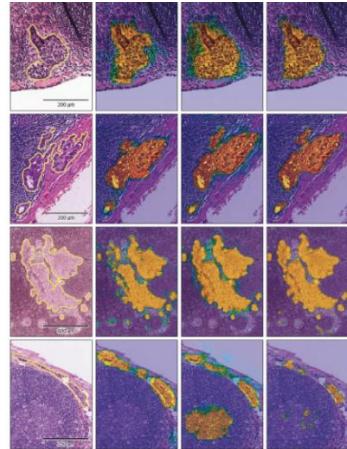
CANCER CAMPUS
GRAND PARIS

Username

Password

Log in

ML/DL in oncology



"Here's my sequence"

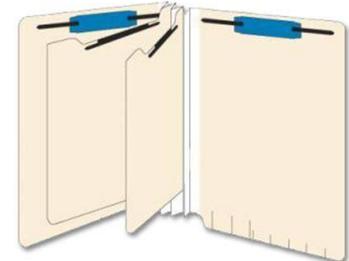
The New Yorker

Diagnostics

Molecular biology
& pronostic

Drug development
& prediction

Monitoring patients



Review for AI in medicine

REVIEW ARTICLE | FOCUS

<https://doi.org/10.1038/s41591-018-0300-7>

nature
medicine

High-performance medicine: the convergence of human and artificial intelligence

Eric J. Topol 

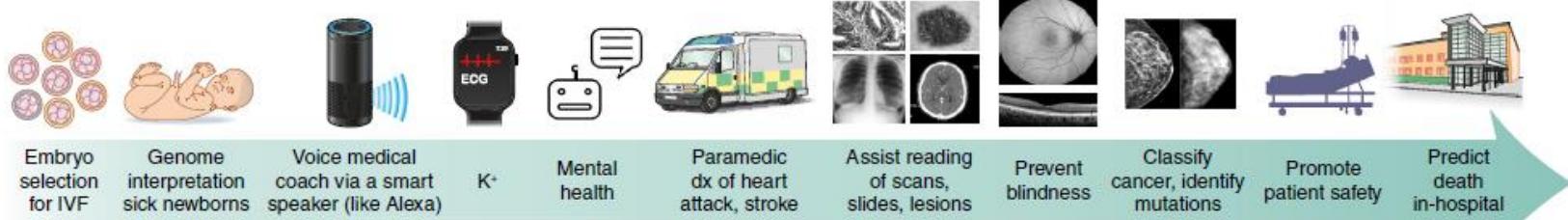


Fig. 2 | Examples of AI applications across the human lifespan. dx, diagnosis; IVF, in vitro fertilization K⁺, potassium blood level. Credit: Debbie Maizels/
Springer Nature

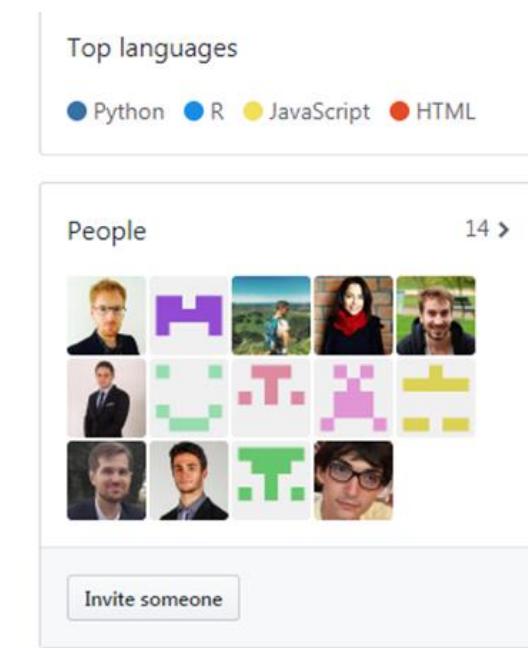
Data Science & DITEP

Medical team



loic.verlingue@gustaveroussy.fr

Data Science team



Info <https://github.com/DITEP>