



# Linux Day 2023

Giornata nazionale a favore della diffusione del software libero e del sistema operativo GNU/Linux



## Applicazioni di Intelligenza artificiale alla scoperta di nuove terapie



Dr. Ugo Perricone

Group Leader in Molecular Informatics @ FONDAZIONE Ri.MED

# Chi Sono

- Laurea in Chimica e Tecnologia Farmaceutica (UNIPA)
- Master in General Management (ISIDA-AFOR)-Specializzazione in Quality Management
- PhD in Scienze Molecolari e Biomolecolari (UNIPA-UNIVIE)



universität  
wien



Università  
degli Studi  
di Palermo

# Molecular Informatics



Ugo Perricone  
Group Leader  
in Molecular  
Informatics

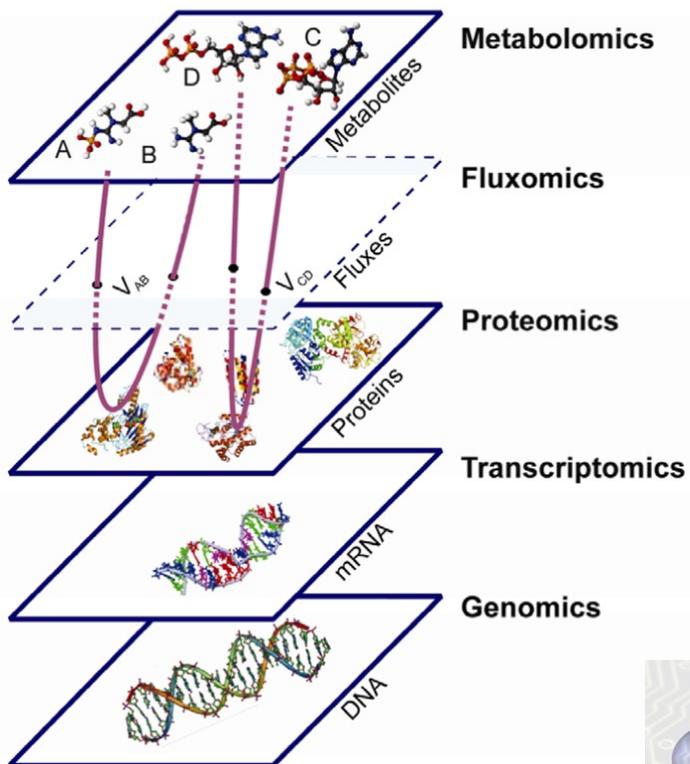


4 FELLOWSHIPS  
4 PhD STUDENTS

# L'interdisciplinarietà di un problema scientifico oggi



## Omics revolution



Metabolomics

Fluxomics

Proteomics

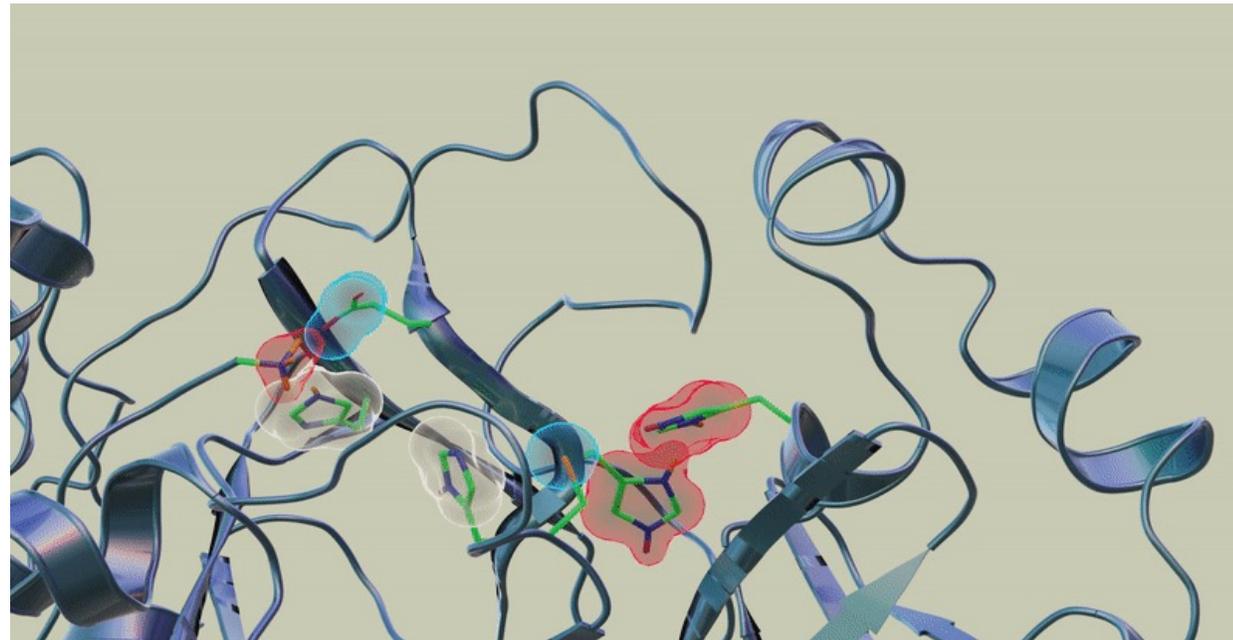
Transcriptomics

Genomics

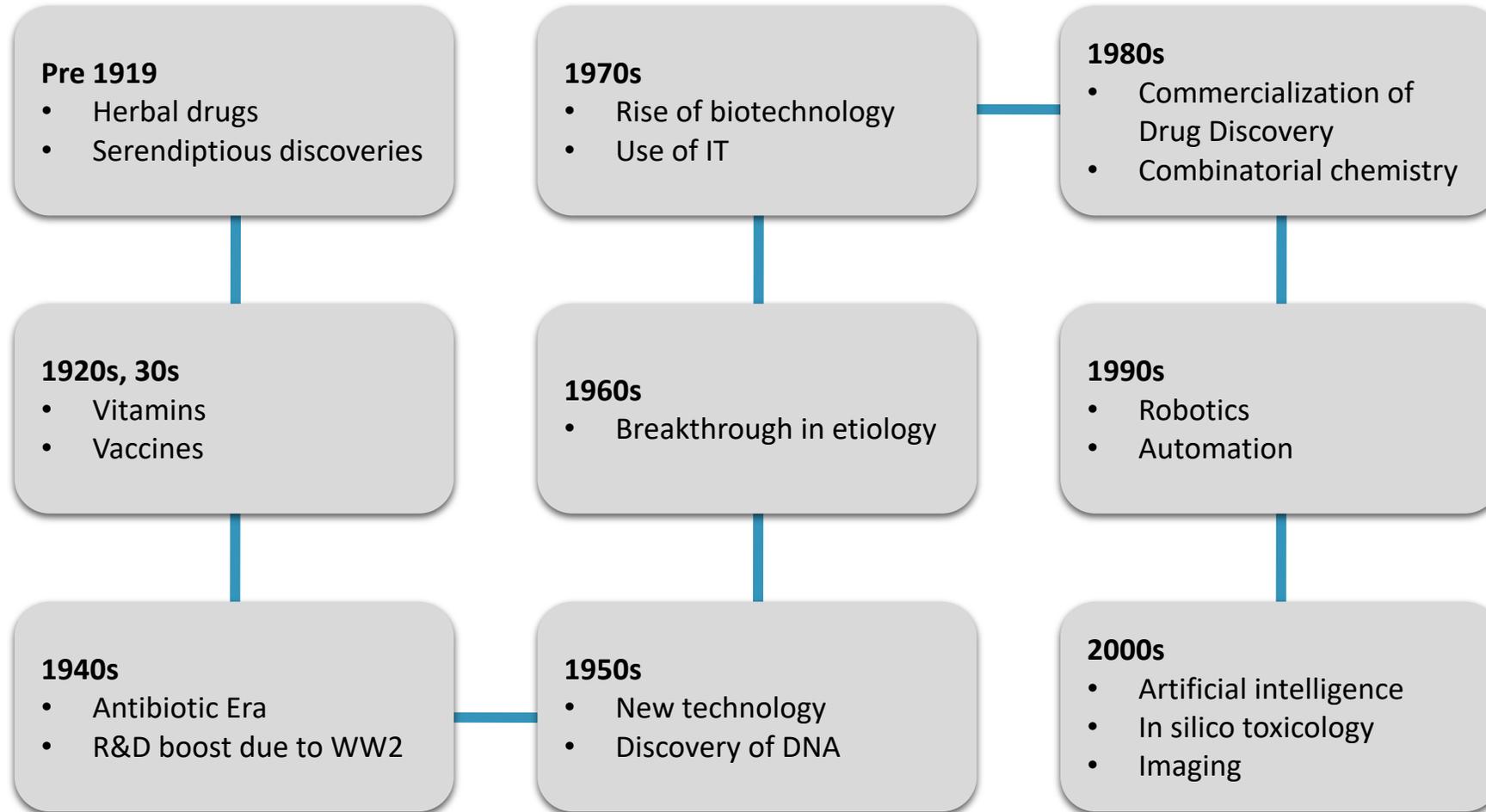
System biology  
 Integrative physiology  
 System medicine  
 System pharmacology  
 Regenerative medicine  
 Integrated biomarkers  
 Human disease  
 Prediction  
 Diagnostics  
 Treatment efficacy

- ✓ Organic Chemistry
- ✓ Inorganic Chemistry
- ✓ Medicinal Chemistry
- ✓ Physical Chemistry
- ✓ Physics (QM & MM)
- ✓ Math
- ✓ Informatics
- ✓ Structural Biology

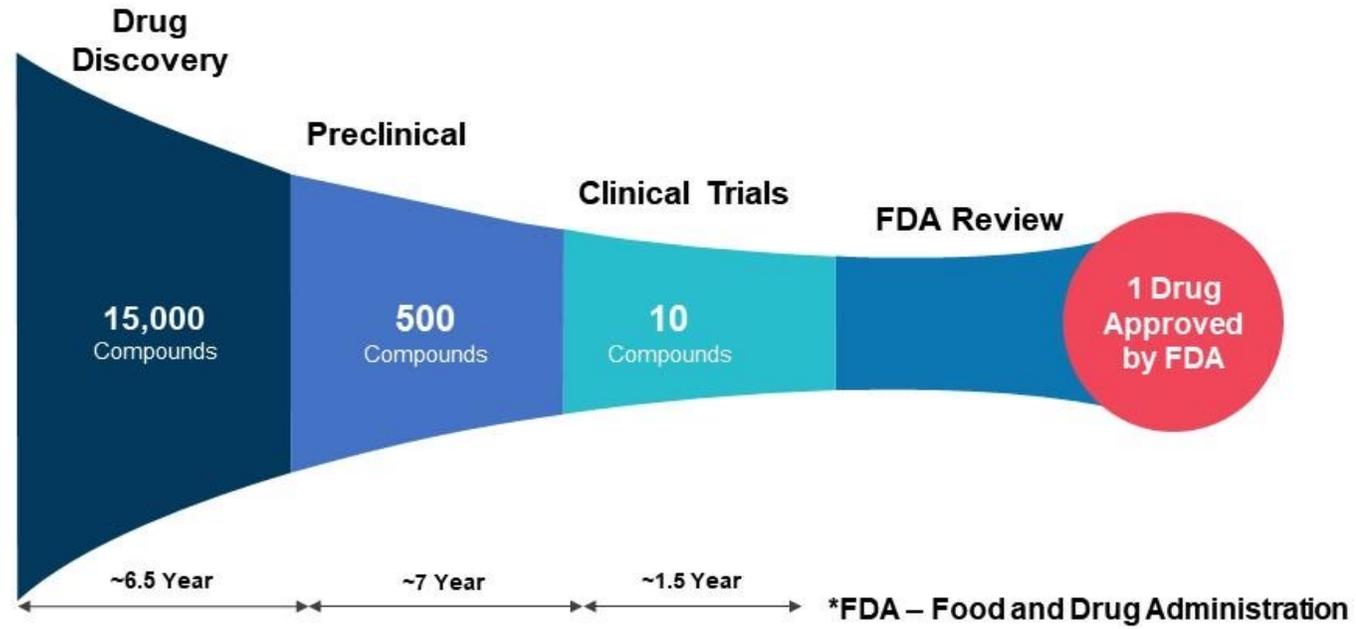
# Un esempio di Big Data nelle scienze... La scoperta di un nuovo farmaco



# Processo di scoperta di nuovi farmaci



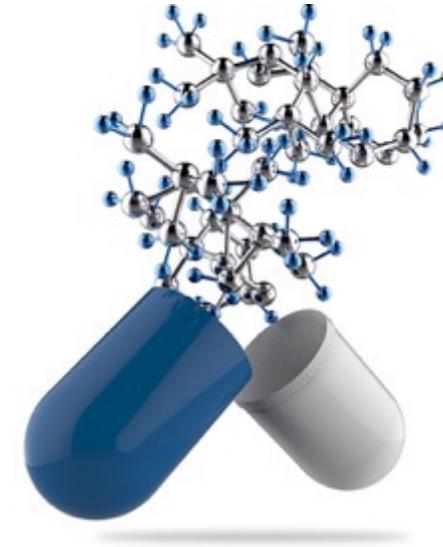
# Drug Discovery & Development - Timeline



## DRUG DISCOVERY

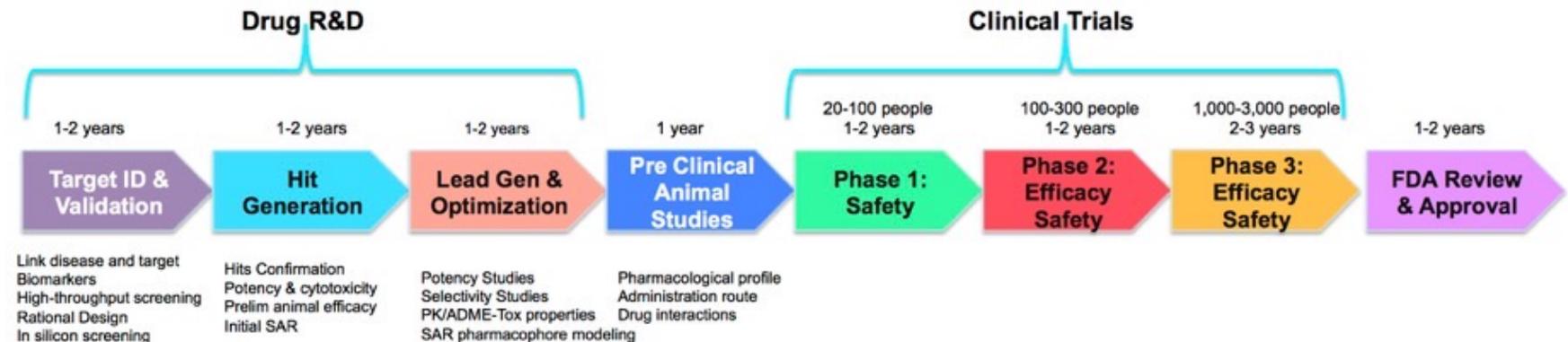
Advanced Data Analysis	Structural Biology and biophysics	Identification of therapeutic targets and screening	Medicinal Chemistry	Molecular Informatics	Proteomics
 <b>Claudia Coronello</b> Principal Investigator in Advanced Data Analysis	 <b>Caterina Alfano</b> Group Leader in Structural Biology and Biophysics	 <b>Chiara Cipollina</b> Group Leader in Experimental Lung Research	 <b>Maria De Rosa</b> Principal Investigator in Medicinal Chemistry	 <b>Ugo Perricone</b> Group Leader in Molecular Informatics	 <b>Simone Dario Scilabra</b> Principal Investigator in Proteomics
1 RESEARCHER 1 FELLOWSHIP 6 PhD STUDENTS	1 RESEARCHER 2 LABORATORY TECHNICIANS 2 PhD STUDENTS 1 TRAINEE	1 RESEARCHER 1 LABORATORY TECHNICIAN 2 FELLOWSHIPS 3 PhD STUDENTS	1 SENIOR SCIENTIST 1 POST-DOC 1 PhD STUDENT	4 FELLOWSHIPS 4 PhD STUDENTS	1 RESEARCHER 3 PhD STUDENTS
 UNIVERSITÀ DEGLI STUDI DI PALERMO	 UNIVERSITÀ DEGLI STUDI DI PALERMO	 Consiglio Nazionale delle Ricerche	 UNIVERSITÀ DEGLI STUDI DI PALERMO	 UNIVERSITÀ DEGLI STUDI DI PALERMO	 ISMETT Istituto di Ricerche e Cura in Genetica Sperimentale

# L'informatica Molecolare nella progettazione farmaceutica (Drug Design)

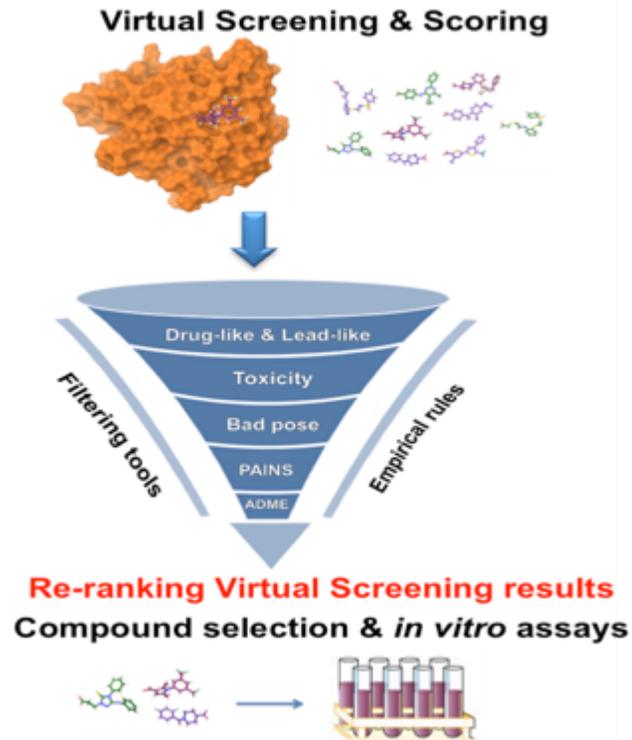


5 Ottobre 1981

*“Next Industrial Revolution: Designing Drugs by Computer at Merck”*



**SCOPO PRINCIPALE: Trovare molecole sicure (non tossiche) e selettive per una patologia e per un target proteico (recettore) → SCREENING VIRTUALE**



(1) Filtrare grandi librerie di molecole e trovare quelle che meglio colpiscono il bersaglio

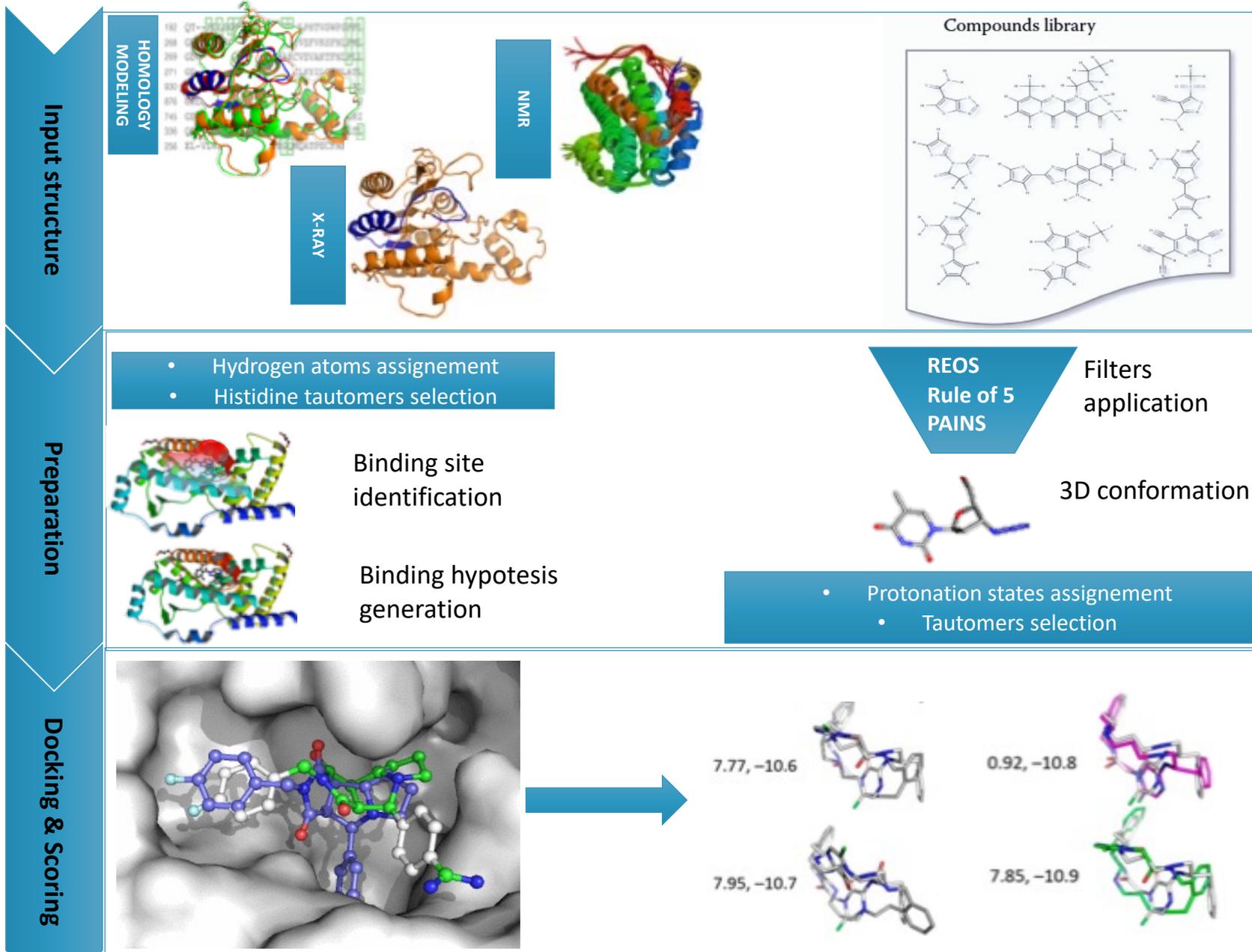
(2) Approfondire il modo in cui il farmaco interagisce con il recettore

(3) Ottimizzare le molecole attive trovate e confermate dai saggi biologici per renderle ancora più sicure ed efficaci

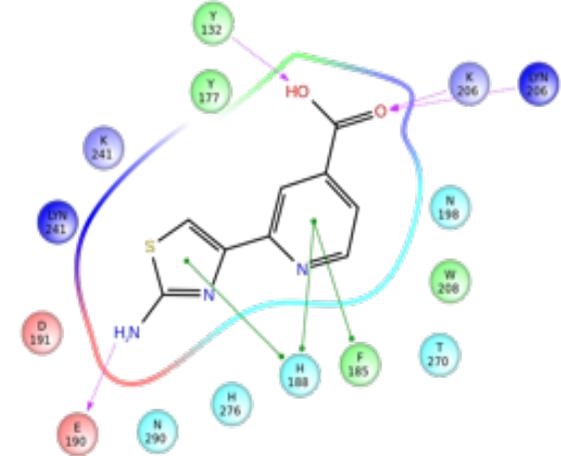
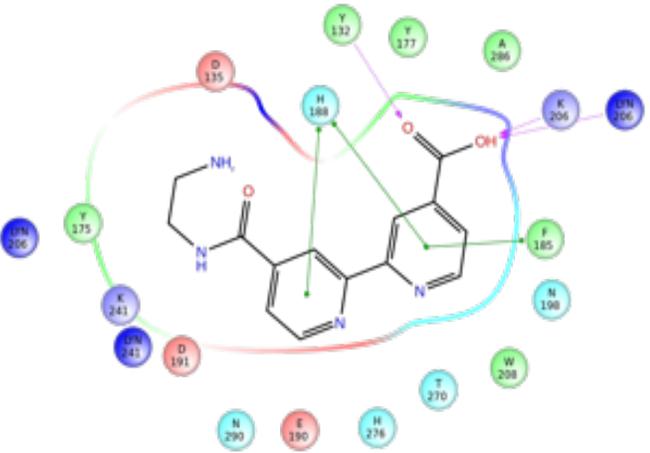
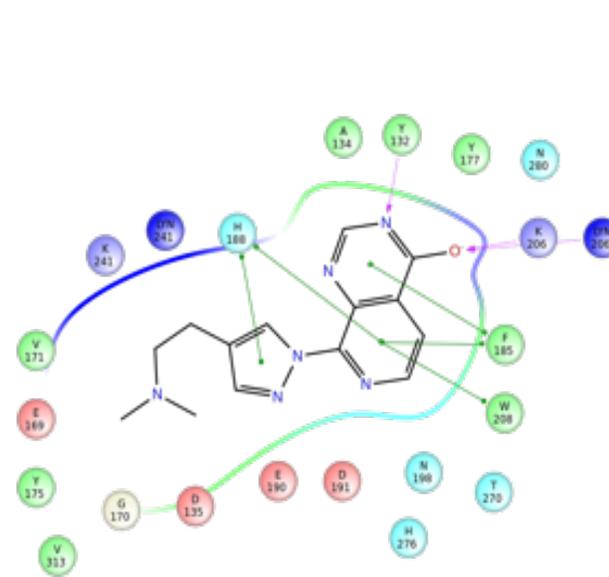
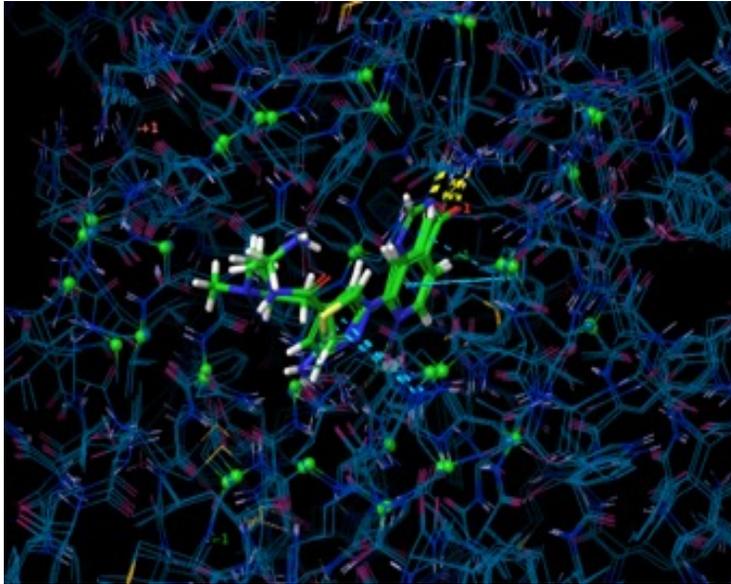
(4) Guidare nella progettazione razionale di nuove molecole

(5) Supportare nella terapia personalizzata

# Molecular Modeling (Un esempio)

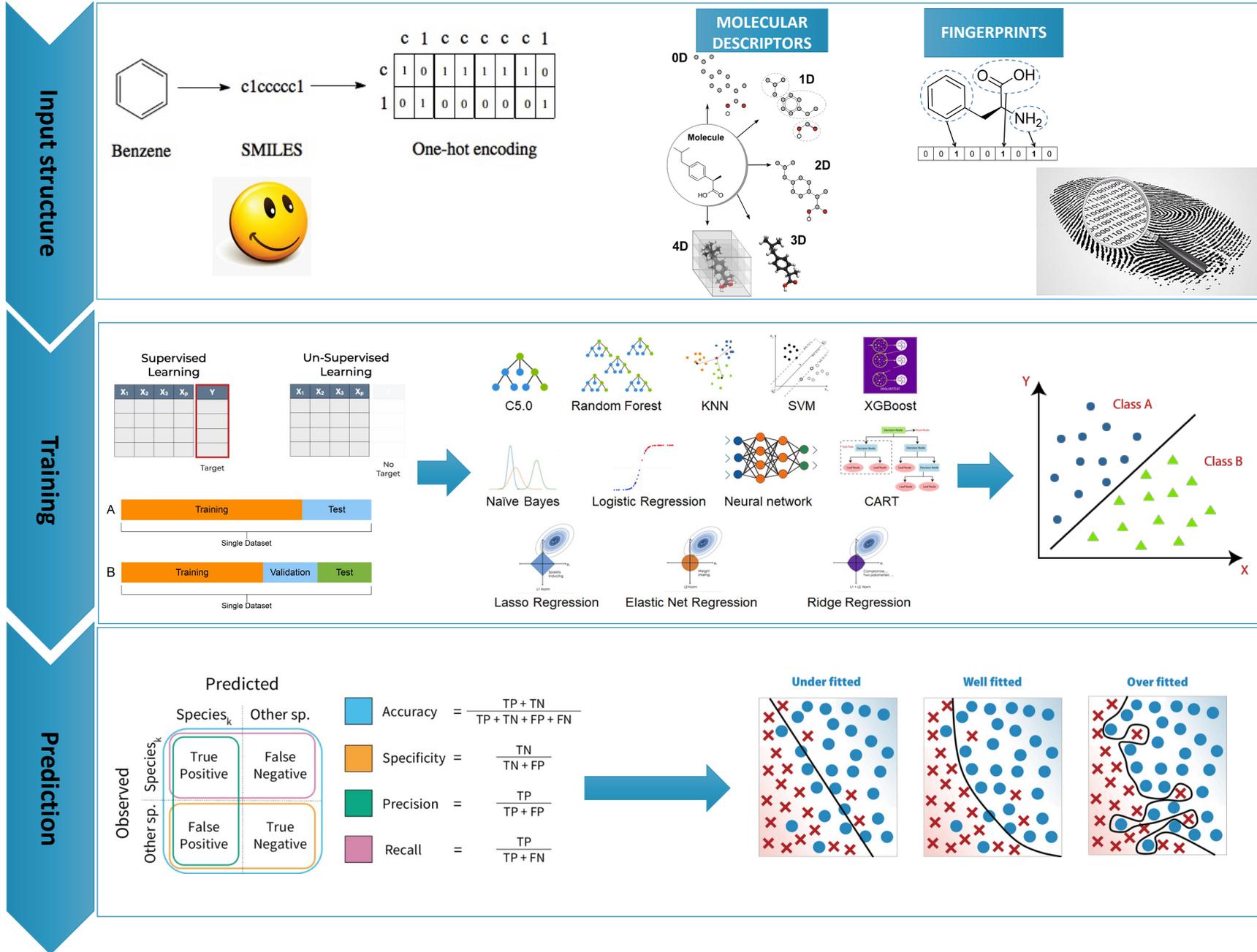


# Structures Analysis



PDB Nes	Residues										
	LYS206:A	LYS241:A	TYR132:A	ASP135:A	PHE185:A	GLU190:A	CYS11:C	SER196:A	SER288:A	ASN198:A	HIS188:A
5F37	✓	✓	✓	✗	✓	✗	✗	✗	✗	✗	✓
3PDQ	✓	✓	✓	✓	✓	✓	✗	✗	✗	✗	✓
3U4S	✓	✗	✓	✗	✗	✗	✓	✓	✓	✓	✗
4URA	✗	✓	✓	✗	✓	✗	✗	✗	✗	✗	✗
4V2V	✓	✗	✓	✗	✗	✗	✗	✓	✓	✓	✗
5A7N	✓	✗	✓	✗	✓	✓	✗	✗	✗	✗	✓
5A7O	✓	✓	✓	✗	✓	✓	✗	✗	✗	✗	✓
5A7P	✓	✗	✓	✗	✓	✓	✗	✗	✗	✗	✗
5A7Q	✓	✗	✓	✗	✓	✓	✗	✗	✗	✓	✓
5F2S	✓	✗	✓	✗	✓	✓	✗	✗	✗	✗	✓
5F3C	✓	✗	✓	✓	✓	✗	✗	✗	✗	✗	✓
5F3E	✓	✗	✓	✗	✓	✗	✗	✗	✗	✗	✓
5F3I	✓	✗	✓	✗	✓	✗	✗	✗	✗	✗	✓
5F32	✓	✗	✓	✗	✓	✓	✗	✗	✗	✗	✓

# Cheminformatics (Un esempio)





## Problema:

- Progettare molecole selettive su specifici target biologici (Kinasi)
- Metodologie classiche non sono in grado di guidare la progettazione selettiva → Il problema va risolto ad un livello superiore

## Soluzione:

- ✓ Sfruttare dati strutturali disponibili
- ✓ Reti neurali in grado di cogliere caratteristiche fondamentali e discriminanti la selettività su un target o su un altro considerando contemporaneamente caratteristiche multiparametriche

**scientific reports**

**OPEN KUALA: a machine learning-driven framework for kinase inhibitors repositioning**

Giada De Simone<sup>1,2,3,†</sup>, Davide Stefano Sardinia<sup>1,2,3,†</sup>, Maria Rita Gulotta<sup>1,2,3</sup> & Ugo Perricone<sup>1,2,3</sup>

**Abstract**  
The family of protein kinases comprises more than 500 genes involved in numerous functions. Hence, their physiological dysfunction has paved the way toward drug discovery for cancer, cardiovascular, and inflammatory diseases. As a matter of fact, Kinase binding sites high similarity has a double role. On the one hand it is a critical issue for selectivity, on the other hand, according to polypharmacology, a synergistic controlled effect on more than one target could be of great pharmacological interest. Another important aspect of binding similarity is the possibility of exploit it for repositioning of drugs on targets of the same family. In this study, we propose our approach called Kinase Drugs Machine Learning Network (KUALA) to automatically identify kinase active ligands by using specific sets of molecular descriptors and provide a multi-target priority score and a repositioning threshold to suggest the best reposable and non-reposable molecules. The comprehensive list of all Kinase-ligand pairs and their scores can be found at <http://github.com/marinofinelli/kuala>.

**Introduction**  
The kinase protein family is one of the most studied in literature because of the key role to many crucial biological processes such as cell division, signaling, and growth. Therefore, physiological dysfunction of the kinase activity have been associated with human disease. Given the importance of these proteins, it is not surprising that their biological role and the selectivity of their modulators are extensively studied. Indeed, this protein family has entered several drug discovery campaigns to treat cancer, cardiovascular, and inflammatory diseases. Selectivity is carefully supervised when designing new drugs, in order to minimize adverse effects on off-targets and consequently to reduce the potential toxicity. However, due to the high similarity of kinase binding sites, the design of novel selective inhibitors for a specific target still remains a challenge today. A low selectivity may influence the clinical trial progress due to high off-target toxicity. An example was the effect of imatinib, a CCRK inhibitor, as the attempt to reach phase III clinical trial. Although the selectivity of a drug towards a specific target should be strongly considered in order to achieve the right balance between the success rate and the possible toxicity on the organism, on the other hand, a multi-target effect might have some interesting applications. In fact, recent polypharmacology studies suggest that the efficacy of a drug can be improved by specifically modulating multiple targets. In other words, a drug that "hits" several targets belonging to one or more pathways (network of interacting proteins) in some cases may represent a more effective therapeutic approach, by limiting the drawbacks generally arising from the use of a single-target drug or from a combination of several drugs. Indeed, a certain rate of drug promiscuity is even sought to improve drugs to new therapeutic targets. As a matter of fact, the repositioning process of known drugs, however, the right balance between protein binding site similarity and the number of targets that a known ligand has had should be seriously considered in drug repositioning studies. In this context, the large amount of data available on public databases like ChEMBL<sup>1</sup> and KEGG<sup>2</sup> has allowed the scientific community to direct efforts towards the use of computational methodologies, such as machine learning (ML) models, to improve the repositioning capabilities of known ligands. In fact, in the last several ML and artificial intelligence (AI) approaches for bioactive ligands identification have been reviewed for protein families, such as G-protein coupled receptors,<sup>3</sup> and human diseases, like the Coronavirus Disease 2019 (COVID-19)<sup>4</sup> and Alzheimer disease.<sup>5</sup> Furthermore, training of ML algorithms has also been reported to develop accurate models for specific targets based on different fingerprint representations of compounds.<sup>6</sup>

**Conclusion**  
The proposed architecture outperforms state-of-the-art ML approaches, and some interesting insights on molecular fingerprints are revealed.

**Keywords**  
Deep learning, Drug design, Molecular fingerprints, Bioactivity prediction, Virtual screening

**BMC Bioinformatics**

**RESEARCH** Open Access

**Conventional architectures for virtual screening**

Isabella Mendolia<sup>1</sup>, Salvatore Contino<sup>1,4,5</sup>, Ugo Perricone<sup>1,2,3</sup>, Edoardo Ardizzone<sup>1</sup> and Roberto Pinna<sup>1</sup>

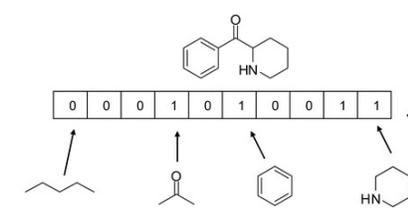
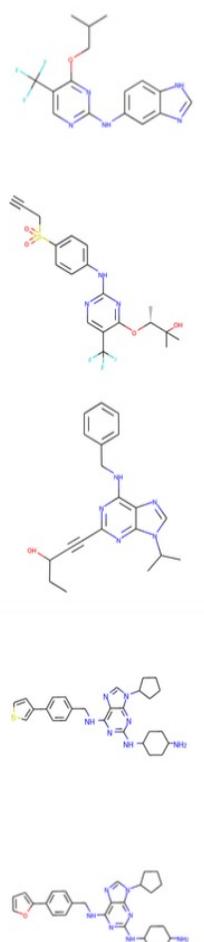
From Annual Meeting of the Bioinformatics Italian Society (BIS 2019)  
Palermo, Italy; 26-28 June 2019

**Background**  
Virtual Screening (VS) is a routinely applied computational technique used for drug design. However, some issues remain uncertain due to the complexity of the algorithms used behind the screening campaign, and this leads to generate models with different prediction reliability. Clinical candidate molecules selected by drug detection must have a profile responding to different criteria, that are based not only on the effect potency but also on the selectivity, safety as well as the so called ADMET properties (Absorption, Distribution, Metabolism, Excretion and Toxicity). Therefore, the design of the optimal compound is a multidimensional challenge involving different aspects of Chemistry and Biology, which can be faced using Machine Learning (ML).

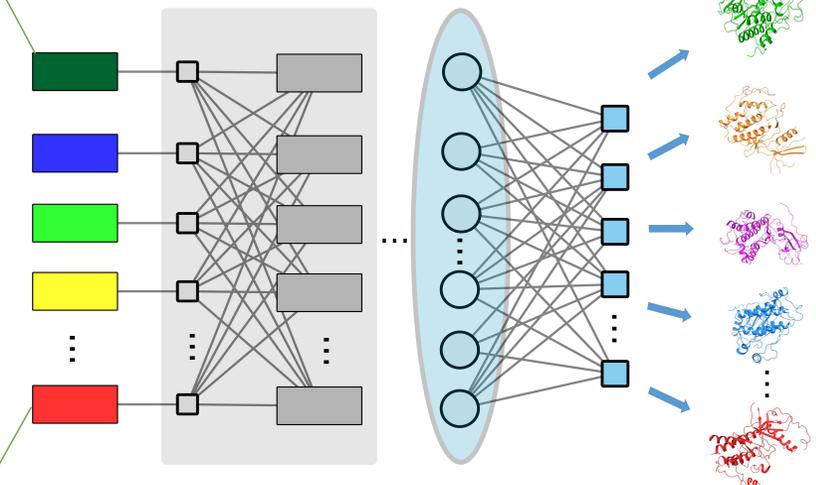
**One key aspect for ML approaches gaining success in property prediction, is the possibility to access and mining large data sets that contain heterogeneous information. Until recent years, the best performing ML techniques were "shallow" ones [1] that is**

**Introduction**  
Drug discovery is a very long and expensive process that includes many stages such as drug target identification, target validation, virtual screening (VS), hit-to-lead generation, lead optimization, and so on [1]. Moreover, developing a new drug has a mean project expenditure above 2 billion USD and takes about 10-15 years [2,3]. Despite the huge investment of time and money, the estimated clinical approval success of innovative small molecules during the drug discovery process is about 13%, that is, the overall risk of failure is very high. Drug design is naturally supported by computational methods in almost every stage. Yu and MacKerell [4] report a review that describes the drug discovery process and the emerging computational drug design methods. Computational methods can also guarantee a systematic assessment of molecular characteristics (e.g., bioactivity, ADMET properties, selectivity, and physicochemical properties) for general lead molecules with favorable properties in silico.

In particular, Virtual Screening (VS) is an often discussed topic in Cheminformatics and Medicinal Chemistry and is widely applied in pharmaceutical research. VS consists of screening large small-molecule databases searching for bioactive molecules with respect to the target under investigation. This enables the researcher to cut the cost of experimentally testing thousands of compounds through a severe reduction in the number of candidate molecules. Research in the field of VS gained increasing importance in the last decade with Deep Learning (DL) becoming a natural discipline [1]. In this field, the scientific challenge is very rich with respect to the proper method for representing molecular structures that are handled by DL models. The very first methods used classical representations such as molecular fingerprints [5] and SMILES notation [6]. Recently, molecular fingerprints have been investigated along with neural embeddings, essentially a learned low-dimensional vector



0	1	0	1	0
0	0	1	0	0
0	1	1	0	0
1	0	1	1	1
0	0	0	1	...
0	1	1	0	0
1	1	1	0	1



## RESEARCH Open Access

### Convolutional architectures for virtual screening

Isabella Mendola<sup>1\*</sup>, Salvatore Contino<sup>1</sup>, Ugo Pericono<sup>2</sup>, Edoardo Andrazze<sup>3</sup> and Roberto Pirrone<sup>4</sup>  
 From Annual Meeting of the Bioinformatics Italian Society BIMS 2019  
 Palermo, Italy, 26-28 June 2019

**Background:** A Virtual Screening algorithm has to adapt to the different stages of this process. Early screening needs to ensure that all bioactive compounds are ranked in the first positions despite of the number of false positives, while a second screening round is aimed at increasing the prediction accuracy.

**Results:** A novel CNN architecture is presented to this aim, which predicts bioactivity of candidate compounds on CDR1 using a combination of molecular fingerprints as their vector representation, and has been trained suitably to achieve good results in regards both enrichment factor and accuracy in different screening modes (KOSM accuracy in active-only selection, and 98.88% in high precision discrimination).

**Conclusions:** The proposed architecture outperforms state-of-the-art ML approaches, and some interesting insights on molecular fingerprints are derived.

**Keywords:** Deep learning, Drug design, Molecular fingerprints, Bioactivity prediction, Virtual screening

**Background**  
 Virtual Screening (VS) is a routinely applied computational technique useful for drug design. However, some issues remain uncertain due to the complexity of the algorithms used behind the screening campaigns, and this leads to generate models with different prediction reliability. Clinical candidate molecules selected by drug detection must have a profile responding to different criteria, that are based not only on the effect potency but also on the selectivity, safety as well as the so called ADMET properties (Absorption, Distribution, Metabolism, Excretion and Toxicity). Therefore, the design of the optimal compound is a multidimensional challenge involving different aspects of Chemistry and Biology, which can be faced using Machine Learning (ML).

One key aspect for ML approaches gaining success is property prediction, is the possibility to access and mining large data sets that contain heterogeneous information.

Until recent years, the best performing ML techniques were "Shallow" ones [1] that is, in the learning stage, neural networks with a limited number of hidden layers (one or two layers) are being developed, although the introduction of a large number of hidden layers (up to the order of 100) and the usage of the so-called "ReLU" (Rectified Linear Unit) as the activation function, led to the emergence of "Deep" learning [2].

Deep learning (DL) is a subset of the machine learning (ML) that uses multiple layers of artificial neural networks to approximate complex functions. DL is a subset of ML that uses multiple layers of artificial neural networks to approximate complex functions. DL is a subset of ML that uses multiple layers of artificial neural networks to approximate complex functions.

International Journal of Molecular Sciences  
 Article  
**EMBER—Embedding Virtual Screening**  
 Isabella Mendola<sup>1\*</sup>, Salvatore Contino<sup>1</sup>

**BMC**

\* Correspondence: [isabella.mendola@unipa.it](mailto:isabella.mendola@unipa.it) (I.M.)  
 † These authors contributed equally to this work.

**Abstract:** In recent years, the debate in the field of applications of Deep Learning to Virtual Screening has focused on the use of neural embeddings with respect to classical descriptors in order to encode both structural and physical properties of ligands and/or targets. The attention on embeddings with the increasing use of Graph Neural Networks aimed at overcoming molecular fingerprints that are short range embeddings for atomic neighborhoods. Here, we present EMBER, a novel molecular embedding made by seven molecular fingerprints arranged as different "spectra" to describe the same molecule, and we prove its effectiveness by using deep convolutional architecture that assesses ligands' bioactivity on a dataset containing twenty protein-ligand complexes with similar binding sites to CDK2. The data set itself is presented, and the architecture is explained in detail along with its training procedure. We report experimental results and an explainability analysis to assess the contribution of each fingerprint to different targets.

**Keywords:** deep learning; drug design; virtual screening; embedding

**1. Introduction**  
 Drug discovery is a very long and expensive process that includes many stages such as drug target identification, target validation, virtual screening (VS), hit-to-lead generation, lead optimization, and so on [1]. Moreover, developing a new drug has a mean price expenditure above 2 billion USD and takes about 10–15 years [2–5]. Despite the huge investment of time and money, the estimated clinical approval success rate of innovative small molecules during the drug discovery process is about 13%. Thus, the overall rate of failure is very high. Drug design is supported by computational methods in almost every stage. Yu and Mackerrell [6] report a review that describes the drug discovery process and the corresponding computer-aided drug design methods. Computational methods do not guarantee a systematic assessment of molecular characteristics (e.g. bioactivity, ADMET properties, selectivity, and physicochemical properties) but generate lead molecules with favorable properties in silico.

In particular, Virtual Screening (VS) is an often discussed topic in Chemoinformatics and Medicinal Chemistry and is widely applied in pharmaceutical research. VS consists of screening large small-molecule databases searching for bioactive molecules with respect to the target under investigation. This enables the researcher to cut the cost of experimentally testing thousands of compounds through a severe reduction in the number of candidate molecules. Research in the field of VS gained increasing importance in the last decade when Deep Learning (DL) became a mature discipline [7]. In this field, the scientific debate is very rich with respect to the proper method for representing molecular structures that are learned by the network. The very first architectures used classical representations such as molecular fingerprints [8] and SMILES notations [9]. Recently, molecular graphs have been investigated along with neural embeddings, essentially a learned low-dimension vector

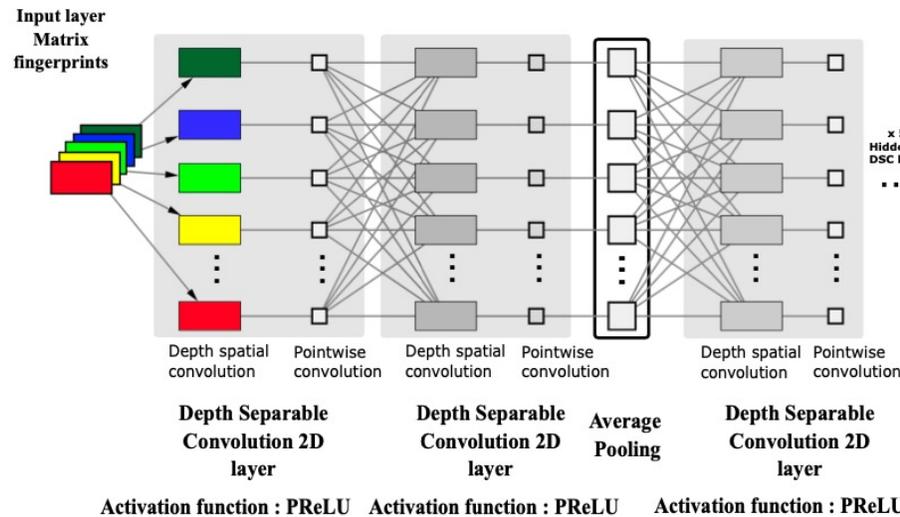
# AI nella progettazione di nuovi antitumorali



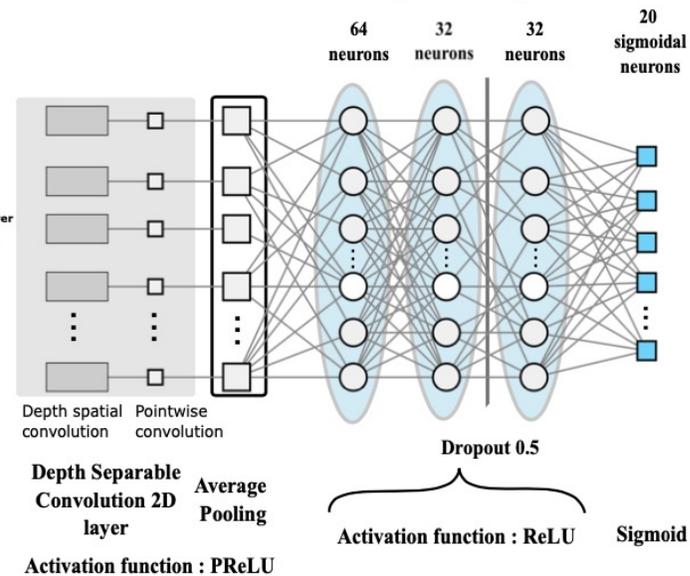
Target	PDB ID	Ligand Code *	Actives	Inactives
ACK	5ZXB	9KO	746	159,775
ALK	6E0R	HKJ	1665	227,247
CDK1	6GU2	F9Z	1241	124,473
CDK2	6INL	AJR	1924	225,087
CDK6	5L2S	6ZV	646	256,561
INSR	5E1S	5JA	1423	195,990
ITK	4RFM	3P6	1001	135,007
JAK2	6M9H	J9D	5526	577,409
JNK3	2B1P	AIZ	658	95,252
MELK	6GVX	TAK	1215	246,662
CHK1	6FC8	D4Q	2175	21,763
CK2a1	6JWA	5ID	1053	10,534
CLK2	6FYL	3NG	671	6800
DYRK1A	4YLK	4E2	1126	11,274
EGFR	5GNK	80U	4757	47,541
ERK2	6OPH	6QB	3525	35,237
GSK3B	5F94	3UO	2578	25,768
IRAK4	6EG9	OLI	2131	21,282
MAPK2K1	4AN9	ACP; 2P7	1254	12,508
PDK1	3NAX	MP7	1117	11,166

Target	Acc.	Loss	Sensitivity	MCC	AUC	F1-Score
ACK	0.9957	0.0226	0.5000	0.6742	0.9834	0.6463
ALK	0.9930	0.0402	0.6575	0.7913	0.9904	0.7804
CDK1	0.9910	0.0314	0.4537	0.6397	0.9850	0.6059
CDK2	0.9859	0.0431	0.5281	0.6338	0.9845	0.6287
CDK6	0.9966	0.0210	0.5865	0.7523	0.9895	0.7305
INSR	0.9893	0.0329	0.3779	0.5830	0.9858	0.5342
ITK	0.9945	0.0232	0.5886	0.7302	0.9905	0.7154
JAK2	0.9898	0.0472	<b>0.8474</b>	<b>0.9090</b>	<b>0.9950</b>	<b>0.9114</b>
JNK3	<b>0.9967</b>	<b>0.0154</b>	0.5905	0.7610	0.9901	0.7381
MELK	0.9957	0.0229	0.7081	0.8270	0.9897	0.8188
CHK1	0.9895	0.0512	0.6385	0.7650	0.9846	0.7565
CK2A1	0.9942	0.0253	0.5166	0.6944	0.9857	0.6667
CLK2	0.9936	0.0259	0.2255	0.4137	0.9771	0.3485
DYRK1A	0.9916	0.0321	0.4080	0.5987	0.9776	0.5591
EGFR	0.9845	0.0604	0.7536	0.8331	0.9874	0.8357
ERK2	0.9881	0.0563	0.7295	0.8292	0.9886	0.8272
GSK3	0.9843	0.0554	0.5827	0.6892	0.9762	0.6856
IRAK4	0.9936	0.0287	0.7611	0.8611	0.9938	0.8571
MAP2K1	0.9931	0.0319	0.5497	0.7184	0.9795	0.6954
PDK1	0.9945	0.0271	0.6310	0.7757	0.9875	0.7613

## Depth Separable Convolution

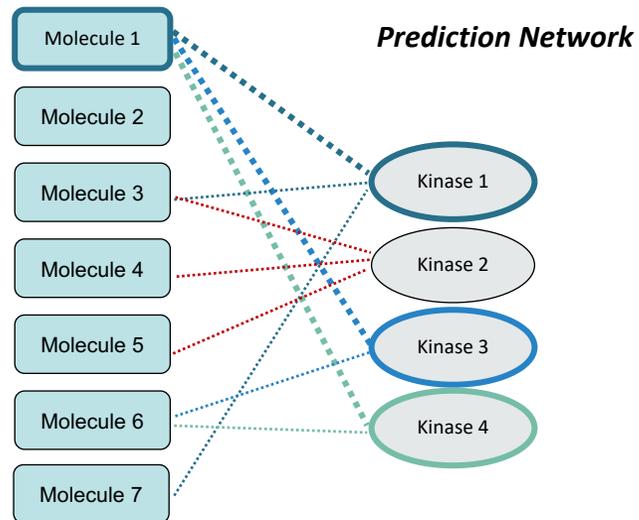
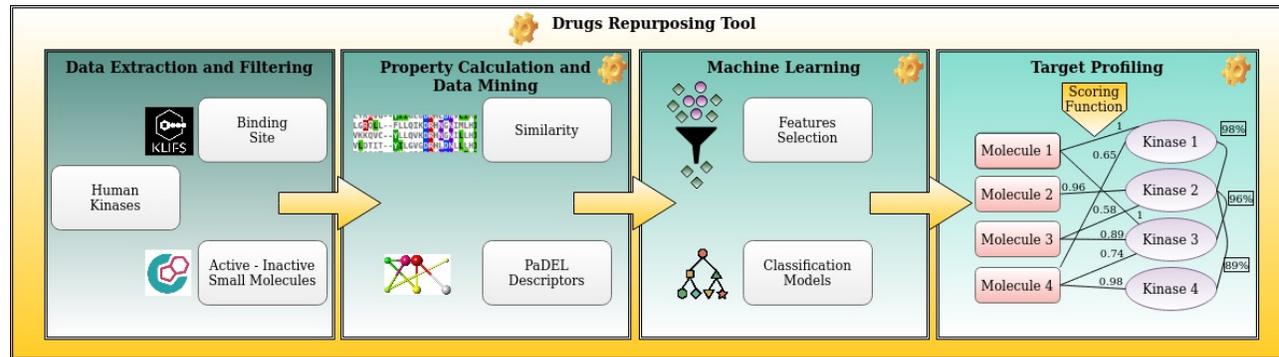


## Fully connected Multi-layer Perceptron



## Drug Repurposing: Kinase drUGs mACHine Learning frAmework (KUALA)- Training

Drug repurposing involves the investigation of existing drugs for new therapeutic purposes

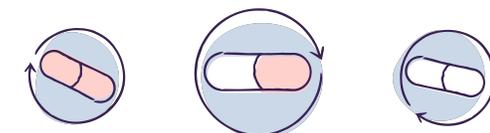


## Repurposing Threshold

 **Similarity** Binding sites alignment

 **Targets** Number of known targets for the ligand

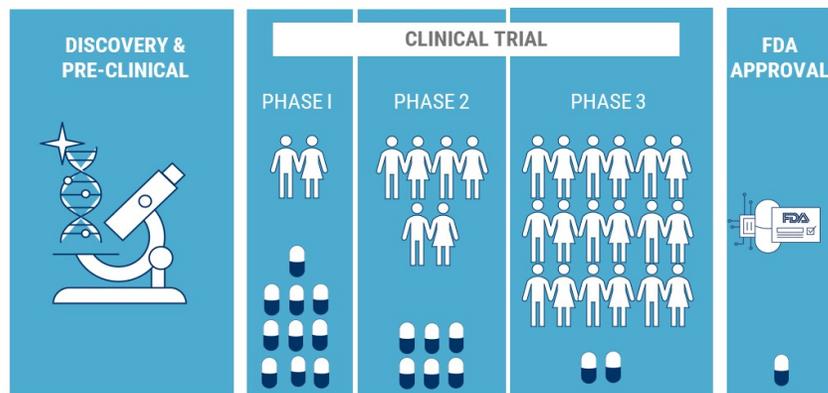
 **Performance** F1-measure of the predictive model



# AI nella progettazione di nuovi antitumorali



## Drug Repurposing: Kinase drUGs mACHine Learning frAmework (KUALA)- Validation



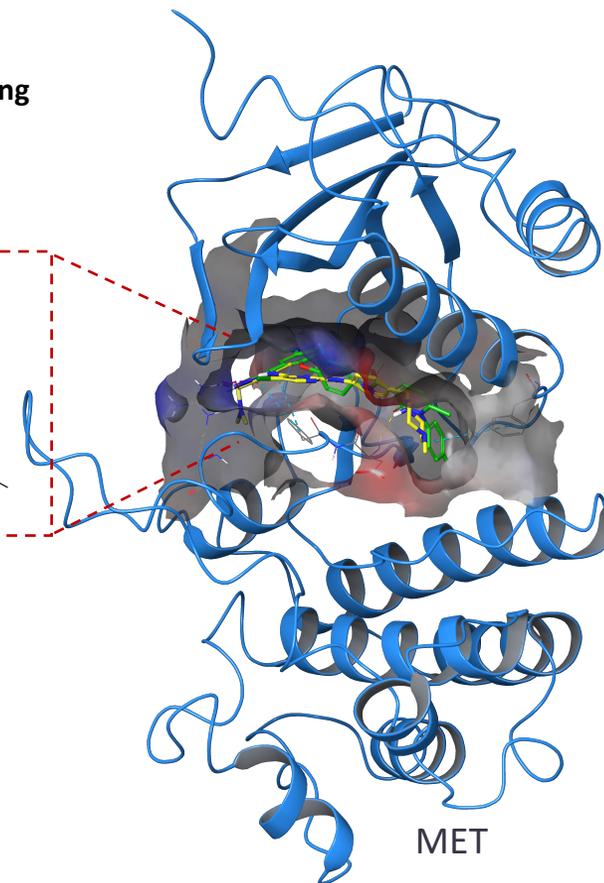
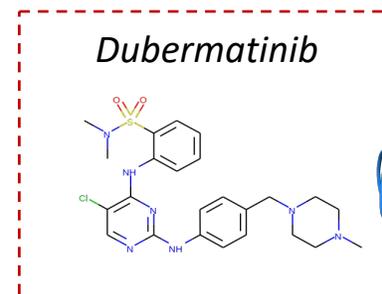
Source: cbinsights.com

### Model Validation with Kinase Inhibitors in Clinical Trials

Kinase models can correctly classify **high percentages** of active molecules



### Example of repurposing



### Predicted inhibitor

DUBERMATINIB (CHEMBL2022968)

- AXL inhibitor, highly effective in inducing apoptosis
- Phase I/II in Patients with Previously Treated CLL

MET

## Problema:

- Progettare molecole selettive sulla proteasi virale di Sars-COV-2 (Mpro)
- Le proteasi hanno similarità dei siti di binding
- Prioritizzazione molecole con metodiche MM classiche insufficiente se usata singolarmente

## Soluzione:

- ✓ Sfruttare dati strutturali disponibili
- ✓ Creazione di un modello di classificazione binaria (Attivo/Inattivo) per 'blindare' la predizione di attività

Article

## Support Vector Machine as a Supervised Learning for the Prioritization of Novel Potential SARS-CoV-2 Main Protease Inhibitors

Nedra Mekni <sup>1,2,\*</sup>, Claudia Coronello <sup>2</sup>, Thierry Langer <sup>1</sup>, Maria De Rosa <sup>2,†</sup> and Ugo Perricone <sup>2,\*†</sup>

<sup>1</sup> Department of Pharmaceutical Chemistry, University of Vienna, 1090 Vienna, Austria;

thierry.langer@univie.ac.at

<sup>2</sup> Drug Discovery Unit, Fondazione RiMED, 90128 Palermo, Italy; ccoronello@fondazionerimed.com (C.C.);

mderosa@fondazionerimed.com (M.D.R.)

\* Correspondence: nmekni@fondazionerimed.com (N.M.); uperricone@fondazionerimed.com (U.P.)

† These authors contributed equally to this work.

**Abstract:** In the last year, the COVID-19 pandemic has highly affected the lifestyle of the world population, encouraging the scientific community towards a great effort on studying the infection molecular mechanisms. Several vaccine formulations are nowadays available and helping to reach immunity. Nevertheless, there is a growing interest towards the development of novel anti-covid drugs. In this scenario, the main protease (Mpro) represents an appealing target, being the enzyme responsible for the cleavage of polypeptides during the viral genome transcription. With the aim of sharing new insights for the design of novel Mpro inhibitors, our research group developed a machine learning approach using the support vector machine (SVM) classification. Starting from a dataset of two million commercially available compounds, the model was able to classify two hundred novel chemo-types as potentially active against the viral protease. The compounds labelled as actives by SVM were next evaluated through consensus docking studies on two PDB structures and their binding mode was compared to well-known protease inhibitors. The best five compounds selected by consensus docking were then submitted to molecular dynamics to deepen binding interactions stability. Of note, the compounds selected via SVM retrieved all the most important interactions known in the literature.

**Keywords:** machine learning; classification; main protease; COVID-19; molecular docking



Citation: Mekni, N.; Coronello, C.; Langer, T.; Rosa, M.D.; Perricone, U. Support Vector Machine as a Supervised Learning for the Prioritization of Novel Potential SARS-CoV-2 Main Protease Inhibitors. *Int. J. Mol. Sci.* **2021**, *22*, 7714. <https://doi.org/10.3390/ijms22147714>

Academic Editor: Daeui Park

Received: 2 July 2021  
Accepted: 15 July 2021  
Published: 19 July 2021

**Publisher's Note:** MDPI stays neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Copyright: © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<https://creativecommons.org/licenses/by/4.0/>).

## 1. Introduction

The COVID-19 pandemic, also known as Severe Acute Respiratory Syndrome Coronavirus-2 (SARS-CoV-2) is afflicting the health and routines of billions of people worldwide.

During the last few months, we are witnessing a race against time to vaccinate as many people as possible; however, the disparities in vaccine distribution between countries and the new emerging variants represent a further public health concern, making it hard to reach a full immunization [1,2].

SARS-CoV-2 is a member of the betacoronavirus family, together with SARS-CoV and Middle East Respiratory Syndrome (MERS-CoV). The enormous scientific effort worldwide led to a better understanding of SARS-CoV-2 structure and the infection mechanism, spotting four main druggable targets, namely the Spike (S) protein, Papain-like protease (PLpro), RNA-dependent RNA polymerase (RdRp) and the main protease/3C-like protease (Mpro/3CLpro) [3,4]. In particular, SARS-CoV-2 Mpro leads a crucial role in the viral replication process. Mpro is a cysteine protease responsible for the cleavage of polypeptides during the viral genome transcription, promoting the generation of non-structural proteins, which can assemble to form new infectious virions. As shown in Figure 1, the Mpro catalytic site includes four subsites, namely S1, S2, S3 and S4, hosting the binding site of protease inhibitors. [5]. Of special importance, the catalytic dyad is enclosed into the

# AI nella progettazione di antivirali (anti Sars-COV-2)



**Abstract:** In the last year, the COVID-19 pandemic has highly affected the lifestyle of the world population, encouraging the scientific community towards a great effort in studying the infection molecular mechanisms. Several vaccine formulations are nowadays available and hoping to reach immunity. Nevertheless, there is a growing interest towards the development of novel anti-viral drugs. In this scenario, the main protease (Mpro) represents an appealing target, being the enzyme responsible for the cleavage of polypeptides during the viral genome transcription. With the aim of sharing new insights for the design of novel Mpro inhibitors, our research group developed a machine learning approach using the support vector machine (SVM) classification. Starting from a dataset of two million commercially available compounds, the model was able to classify two hundred novel chemical types as potentially active against the viral protease. The compounds labeled as active by SVM were next evaluated through consensus docking studies on two FDB structures and their binding mode was compared to well-known protease inhibitors. The best five compounds selected by consensus docking were then submitted to molecular dynamics to deepen binding interactions stability. Of note, the compounds selected via SVM retrieved all the most important interactions known in the literature.

**Check for updates**

**Keywords:** machine learning; classification; main protease; COVID-19; molecular docking

**Academic Editor:** Christ Park

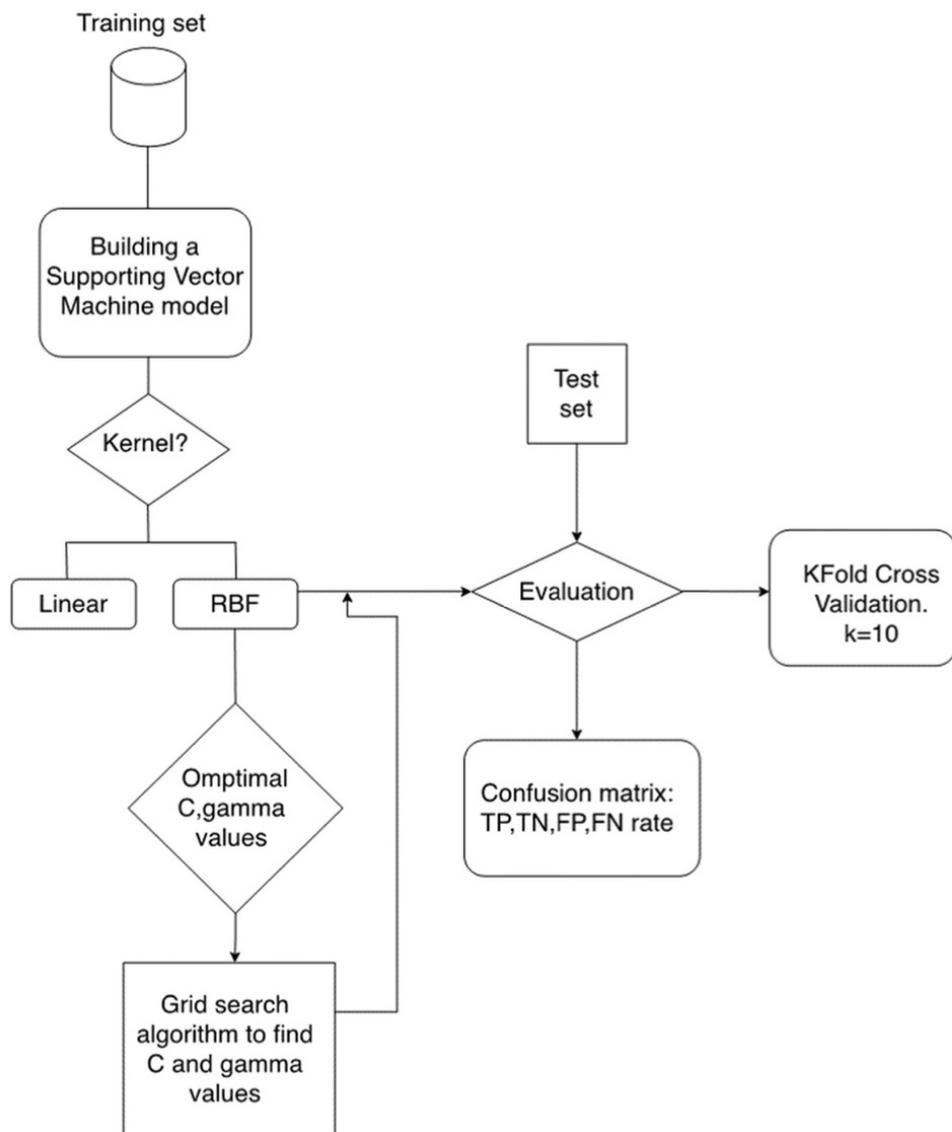
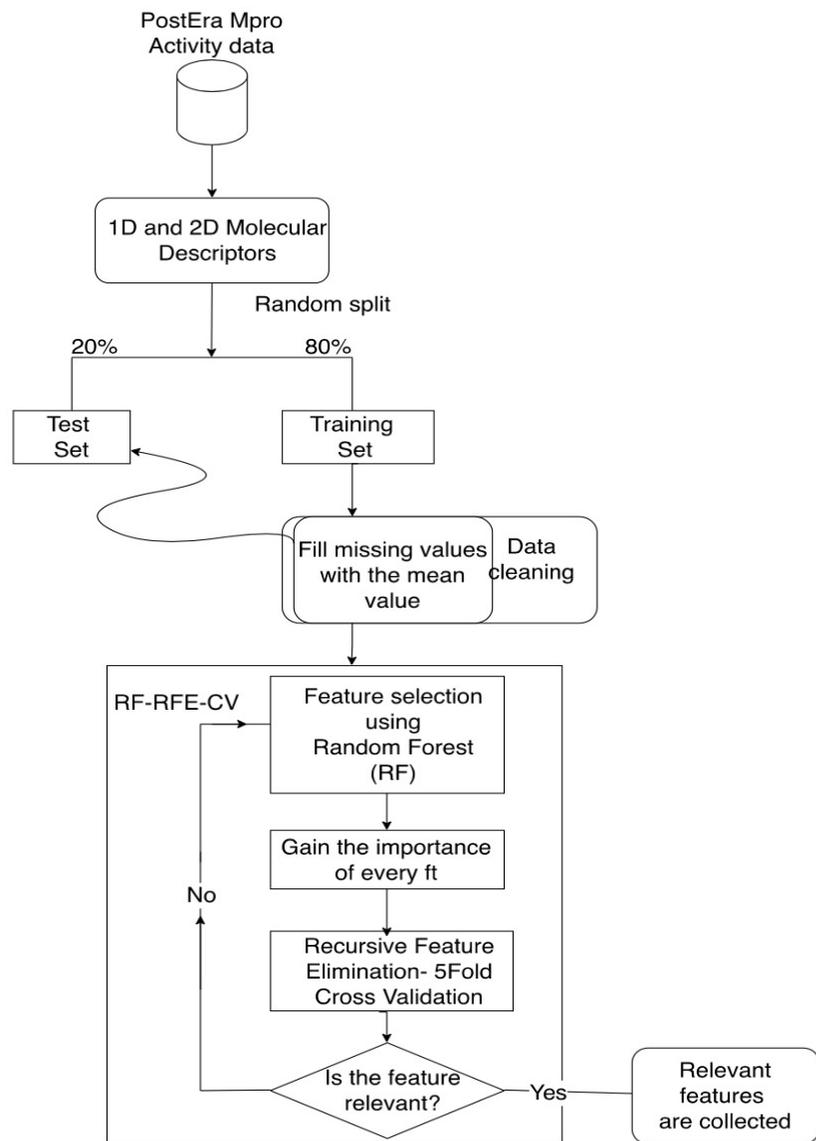
**Received:** 23 July 2021

**Accepted:** 17 July 2021

**Published:** 19 July 2021

**Copyright:** © 2021 by the authors. Licensee MDPI, Basel, Switzerland. This article is an open access article distributed under the terms and conditions of the Creative Commons Attribution (CC BY) license (<http://creativecommons.org/licenses/by/4.0/>).

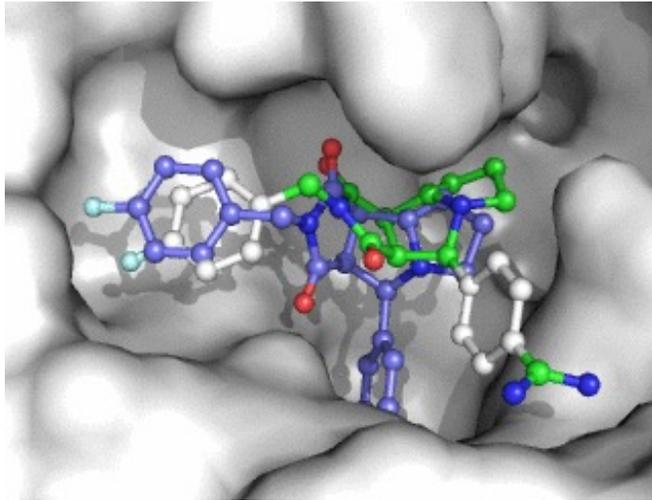
**1. Introduction**  
The COVID-19 pandemic, also known as Severe Acute Respiratory Syndrome Coronavirus-2 (SARS-CoV-2) is afflicting the health and routines of billions of people worldwide. During the last few months, we are witnessing a race against time to vaccinate as many people as possible; however, the disparities in vaccine distribution between countries and the new emerging variants represent a further public health concern, making it hard to reach a full immunization [1,2].  
SARS-CoV-2 is a member of the betacoronavirus family, together with SARS-CoV and Middle East Respiratory Syndrome (MERS-CoV). The enormous scientific effort worldwide led to a better understanding of SARS-CoV-2 structure and the infection mechanism, spotting four main druggable targets, namely the Spike (S) protein, Papain-like protease (PLpro), RNA-dependent RNA polymerase (RdRp) and the main protease (3C-like protease) (Mpro) (S1, S2, S3, S4) [3]. In particular, SARS-CoV-2 Mpro holds a crucial role in the viral replication process. Mpro is a cysteine protease responsible for the cleavage of polypeptides during the viral genome transcription, promoting the generation of non-structural proteins, which can assemble to form new infectious viruses. As shown in Figure 1, the Mpro catalytic site includes four subsites, namely S1, S2, S3 and S4, housing the binding site of protease inhibitors [1]. Of special importance, the catalytic dyad is enclosed into the



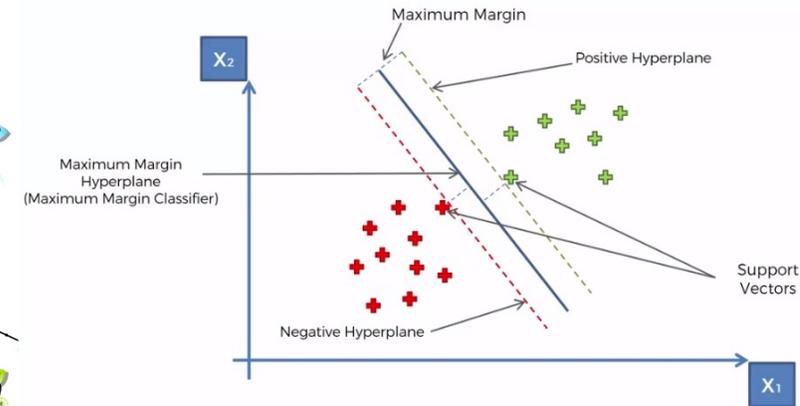
## Ranking /Prioritisation

Cmpd	Docking Pose	Ligand Interaction
I		
II		
III		

## Molecular Docking



## Modello ML per la classificazione



International Journal of Molecular Sciences

Article  
**Support Vector Machine as a Supervised Learning for the Prioritization of Novel Potential SARS-CoV-2 Main Protease Inhibitors**

Neda Mekki <sup>1,†</sup>, Claudia Comendoli <sup>1</sup>, Thierry Langer <sup>1</sup>, Maria De Rosa <sup>1,†</sup> and Ugo Perleone <sup>1,\*,†</sup>

<sup>1</sup> Department of Pharmaceutical Chemistry, University of Vienna, 1090 Vienna, Austria; \* [perleone@pharmazie.at](mailto:perleone@pharmazie.at); <sup>†</sup> [nedamekki@pharmazie.at](mailto:nedamekki@pharmazie.at); [claudia.comendoli@pharmazie.at](mailto:claudia.comendoli@pharmazie.at); [thierry.langer@pharmazie.at](mailto:thierry.langer@pharmazie.at); [maria.derosa@pharmazie.at](mailto:maria.derosa@pharmazie.at)

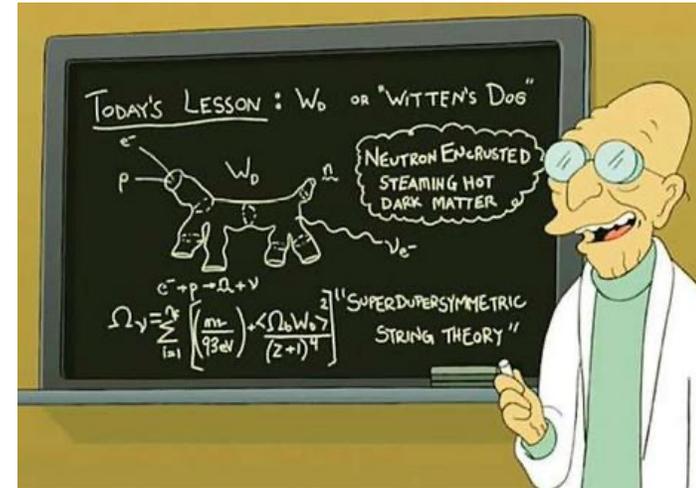
**Abstract:** In the last year, the COVID-19 pandemic has highly affected the lifestyle of the world population, encouraging the scientific community towards a goal effort on studying the infection molecular mechanism. Several vaccine formulations are nowadays available and helping to reach immunity. Nevertheless, there is a growing interest towards the development of novel antiviral drugs. In this scenario, the main protease (Mpro) represents an appealing target, being the enzyme responsible for the cleavage of polyproteins during the viral genome transcription. With the aim of sharing new insights for the design of novel Mpro inhibitors, our research group developed a machine learning approach using the support vector machine (SVM) classification. Starting from a dataset of two million commercially available compounds, the model was able to identify two hundred novel chemo-types as potentially active against the viral protease. The compounds labeled as active by SVM were most evaluated through consensus docking studies on PDB structures and their binding mode was compared to well-known protease inhibitors. The best five compounds selected by consensus docking were then subjected to molecular dynamics to deepen binding interactions stability. Of note, the compounds selected via SVM retrieved all the most important interactions known in the literature.

**Keywords:** machine learning; classification; main protease; COVID-19; molecular docking

**1. Introduction**  
 The COVID-19 pandemic, also known as Severe Acute Respiratory Syndrome Coronavirus 2 (SARS-CoV-2) affecting the health and mortality of billions of people worldwide. During the last few months, we are witnessing a race against time to vaccinate as many people as possible, however, the disparities in vaccine distribution between countries and the new emerging variants represent a further public health concern, making it hard to reach a full immunization [1].

SARS-CoV-2 is a member of the betacoronavirus family together with SARS-CoV and Middle East Respiratory Syndrome (MERS-CoV). The infectious scientific effort worldwide led to a better understanding of SARS-CoV-2 structure and the infection mechanism, spotting four main druggable targets, namely the Spike (S) protein, Papain-like protease (PLpro), RNA-dependent RNA polymerase (RdRp) and the main protease (3C-like protease) (Mpro/3CLpro) [1]. In particular, SARS-CoV-2 Mpro has a crucial role in the viral replication process. Mpro is a cysteine protease responsible for the generation of polyproteins during the viral genome transcription, promoting the generation of non-structural proteins, which can assemble to form new infectious viruses. As shown in Figure 1, the Mpro catalytic site includes four subsites, namely S1, S2, S3 and S4, hosting the binding site of protease inhibitors [1]. Of special importance, the catalytic dyad is embedded into

# L'importanza di modelli computazionali (*In silico*) sempre più evoluti



## ORIGINAL RESEARCH ARTICLE

Front. Physiol., 12 September 2017 | <https://doi.org/10.3389/fphys.2017.00668>



## Human *In Silico* Drug Trials Demonstrate Higher Accuracy than Animal Models in Predicting Clinical Pro-Arrhythmic Cardiotoxicity

Elisa Passini<sup>1\*</sup>, Oliver J. Britton<sup>1</sup>, Hua Rong Lu<sup>2</sup>, Jutta Rohrbacher<sup>2</sup>, An N. Hermans<sup>2</sup>, David J. Gallacher<sup>2</sup>, Robert J. H. Greig<sup>3</sup>, Alfonso Bueno-Orovio<sup>1</sup> and Blanca Rodriguez<sup>1</sup>

<sup>1</sup>Computational Cardiovascular Science Group, Department of Computer Science, University of Oxford, Oxford, United Kingdom

<sup>2</sup>Global Safety, Pharmacology, Discovery Sciences, Janssen Research and Development, Janssen Pharmaceutica NV, Beerse, Belgium

<sup>3</sup>Oxford Computer Consultants Ltd., Oxford, United Kingdom



News & Comment Research

News Opinion Research Analysis Careers Books & Culture

NEWS • 11 JULY 2018

## Software beats animal tests at predicting toxicity of chemicals

Machine learning on mountain of safety data improves automated assessments.

Richard Van Noorden

## Problema:

- La tossicità di un farmaco è un fattore multiparametrico
- Testare la tossicità comporta il sacrificio di molti animali e costi molto elevati
- Trattandosi di fattori multiparametrici, molti test in vivo sono poco affidabili perché il «sistema» animale è troppo differente da quello umano

## Soluzione:

- ✓ Applicare multiclassificatori che considerino contemporaneamente tutti i parametri disponibili
- ✓ Applicazione di modelli quali-quantitativi per permettere l'ottimizzazione molecolare e la riduzione della tossicità

## Vantaggi:

- Predire per molecole simili a quelle studiate la possibile tossicità e anticipare fallimenti nelle campagne di drug discovery
- Lavorare con dati umani opportunamente modellati e possibilità di trasferire i modelli da un comparto ad un altro

### Machine Learning of Toxicological Big Data Enables Read-Across Structure Activity Relationships (RASAR) Outperforming Animal Test Reproducibility

Thomas Luechtefeld,<sup>\*,†</sup> Dan Marsh,<sup>†</sup> Craig Rowlands,<sup>‡</sup> and Thomas Hartung<sup>\*,§,1</sup>

<sup>†</sup>Johns Hopkins University, Bloomberg School of Public Health, Center for Alternatives to Animal Testing (CAAT), Baltimore, Maryland 21205; <sup>‡</sup>ToxTrack, Baltimore, Maryland 21209; <sup>§</sup>UL Product Supply Chain Intelligence, Underwriters Laboratories (UL), Northbrook, Illinois 60062; and <sup>\*</sup>University of Konstanz, CAAT-Europe, Konstanz 78464, Germany

## nature

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

[nature](#) > [news](#) > [article](#)

NEWS | 11 July 2018

### Software beats animal tests at predicting toxicity of chemicals

Machine learning on mountain of safety data improves automated assessments.

[Richard Van Noorden](#)

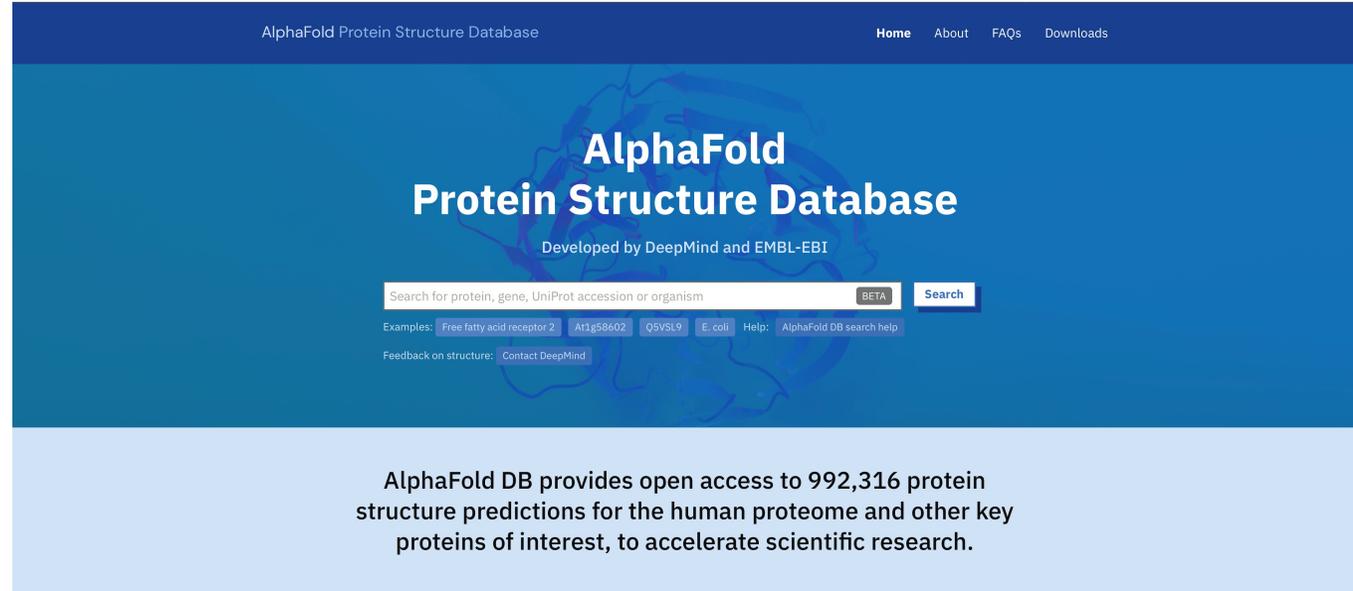


## Problema:

- Non tutte le strutture tridimensionali delle proteine sono disponibili (es. difficoltà legate ad isolamento e purificazione proteina)
- Gli studi strutturali stanno alla base della comprensione di processi biologici e patofisiologici

## Soluzione:

- ✓ Sfruttare tutti i dati disponibili in letteratura su sequenze e folding proteine (input : training/validation/test sets)
- ✓ Addestrare reti in grado di predire la struttura tridimensionale partendo dalla sequenza aminoacidica
- ✓ Predire strutture non disponibili



AlphaFold Protein Structure Database

Home About FAQs Downloads

## AlphaFold Protein Structure Database

Developed by DeepMind and EMBL-EBI

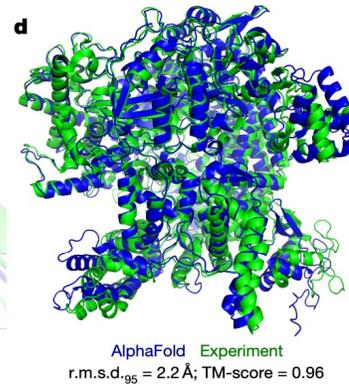
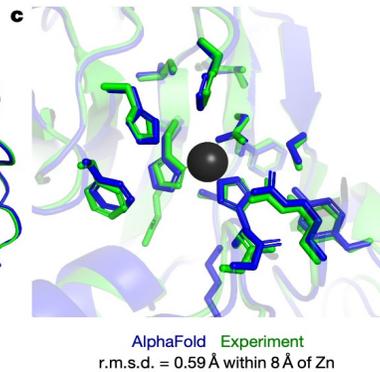
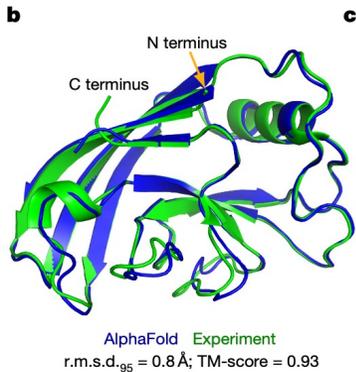
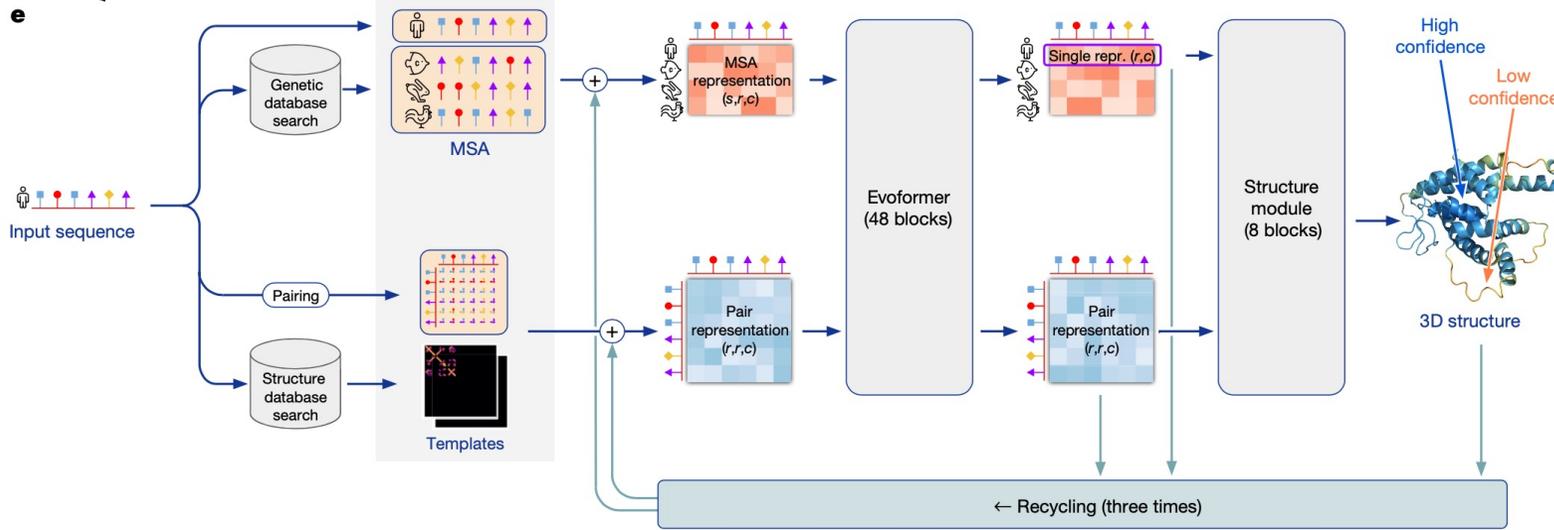
Search for protein, gene, UniProt accession or organism

Examples: [Free fatty acid receptor 2](#) [A1g58602](#) [Q5VSL9](#) [E. coli](#) [Help: AlphaFold DB search help](#)

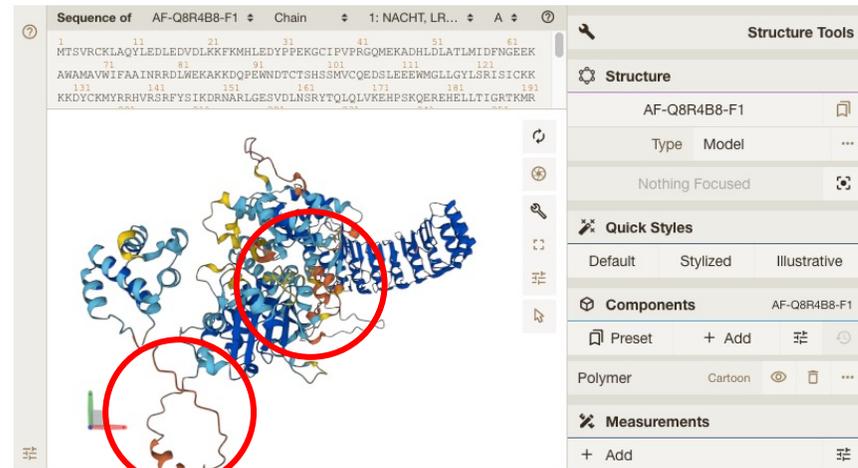
Feedback on structure: [Contact DeepMind](#)

AlphaFold DB provides open access to 992,316 protein structure predictions for the human proteome and other key proteins of interest, to accelerate scientific research.

# AI nella predizione di Strutture 3D...il caso AlphaFold



## 3D viewer

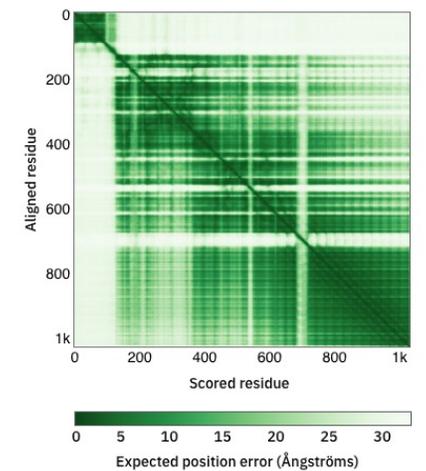


## Model Confidence

- Very high (pLDDT > 90)
- High (90 > pLDDT > 70)
- Low (70 > pLDDT > 50)
- Very low (pLDDT < 50)

AlphaFold produces a per-residue model confidence score (pLDDT) between 0 and 100. Some regions below 50 pLDDT may be unstructured in isolation.

## Predicted aligned error (PAE)



Click and drag a box on the PAE viewer to select regions of the structure and highlight them on the 3D viewer.

PAE data is useful for assessing inter-domain accuracy – go to [Help section](#) below for more information.

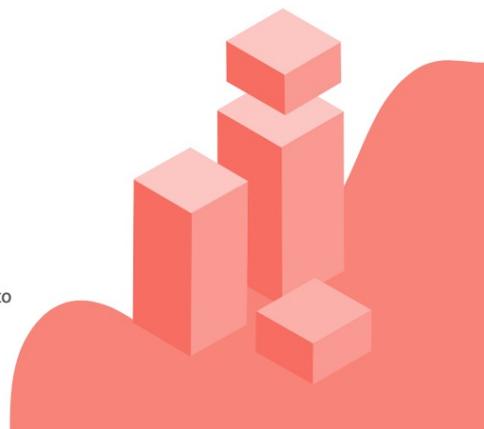
# Un delle più grandi rivoluzioni nel drug discovery moderno?



[Get Started](#) [Tutorials](#) [Docs](#) [Contribute](#) [Publications](#) [About](#) [GitHub](#)

## A powerful and flexible machine learning platform for drug discovery

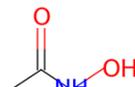
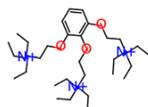
TorchDrug is a machine learning platform designed for drug discovery, covering techniques from graph machine learning (graph neural networks, geometric deep learning & knowledge graphs), deep generative models to reinforcement learning. It provides a comprehensive and flexible interface to support rapid prototyping of drug discovery models in PyTorch.



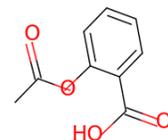
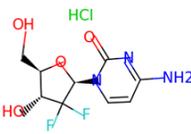
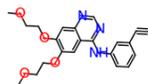
```
$ pip install torchdrug
```

[COPY](#)

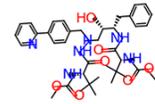
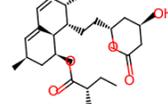
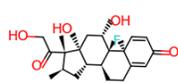
non-toxic  
approved



toxic  
not approved



toxic  
approved



## Key Features

### Minimal Domain Knowledge

Build and train machine learning models for drug discovery with minimal domain knowledge.

### Comprehensive Benchmarks

Benchmarks provide a systematic comparison of deep learning architectures for drug discovery.

### Datasets and Building Blocks

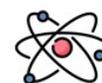
Empower fast iteration of ideas by a large collection of common datasets and building blocks.

### Scalable Training and Inference

Seamlessly scale models to multiple CPUs, multiple GPUs, or even distributed settings.

[VIEW ALL FEATURES](#)

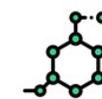
## Solutions to Drug Discovery Tasks



**Property Prediction**



**Pretrained Molecular Representations**



**De Novo Molecule Design & Optimization**



**Reaction Prediction & Retrosynthesis**



**Biomedical Knowledge Graph Reasoning**



**Protein Representation Learning**

## Pro e contro di AI in drug discovery



- Taglio drastico dei tempi di ricerca e dei costi
  - Maggiore precisione nella progettazione di terapie
  - Anticipo nella predizione di effetti tossici per salute umana e ambiente
  - Riduzione drastica di uso di modelli animali
  - Impiego per terapie personalizzate
- 
- Necessaria una profonda conoscenza di dominio medico-scientifico per applicazioni in questo campo...rischio alto di misinterpretazione dei risultati
  - Necessità di grandi moli di dati per l'addestramento...non sempre disponibili
  - Qualità dei dati presenti nei DB pubblici non sempre adeguata

## AI distinction

INS018\_055 is Insilico's first wholly owned program in which AI was used to identify a novel target and generate novel small molecules through the company's Pharma.AI platform.

Insilico began development of INS018\_055 in February 2021 using Pharma.AI. The platform incorporates a pair of specific-function platforms.



Zhavoronkov, PhD, Insilico Medicine founder and CEO; and Feng Ren, PhD, Co-CEO and CSO and Head of R&D, inside the company's AI-powered robotics lab. Based in Hong Kong, Insilico supports every one of its development programs through an autonomous robotics laboratory in Suzhou, China. [Insilico Medicine]

## Three clinical candidates

INS018\_055 is one of three Insilico candidates to have advanced into clinical trials. The others are ISM3312, an oral small molecule 3CLPro inhibitor designed to treat COVID-19; and ISM3091, a small molecule ubiquitin specific protease 1 (USP1) inhibitor being developed to treat BRCA-mutant breast cancer. ISM3312 is now in a Phase I trial and actively dosing patients, Zhavoronkov said.

# JEDI BILLION MOLECULES AGAINST COVID19 GRANDCHALLENGE

Screening billions of possible molecular compounds that can block SARS-CoV-2

54 billion molecules

130 teams

500 scientists

ONE CHALLENGE

## molecular informatics models – molecules – systems

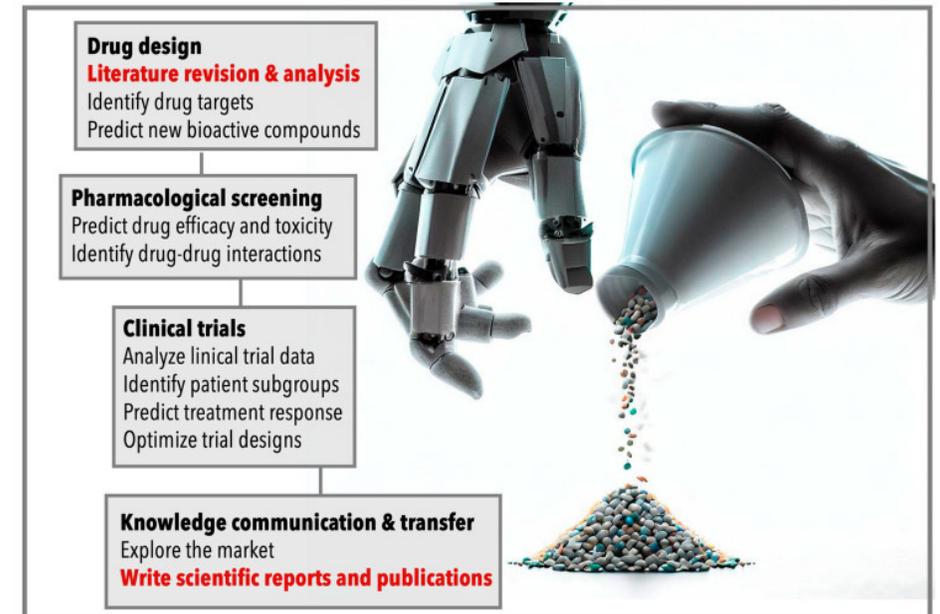
### Research Article

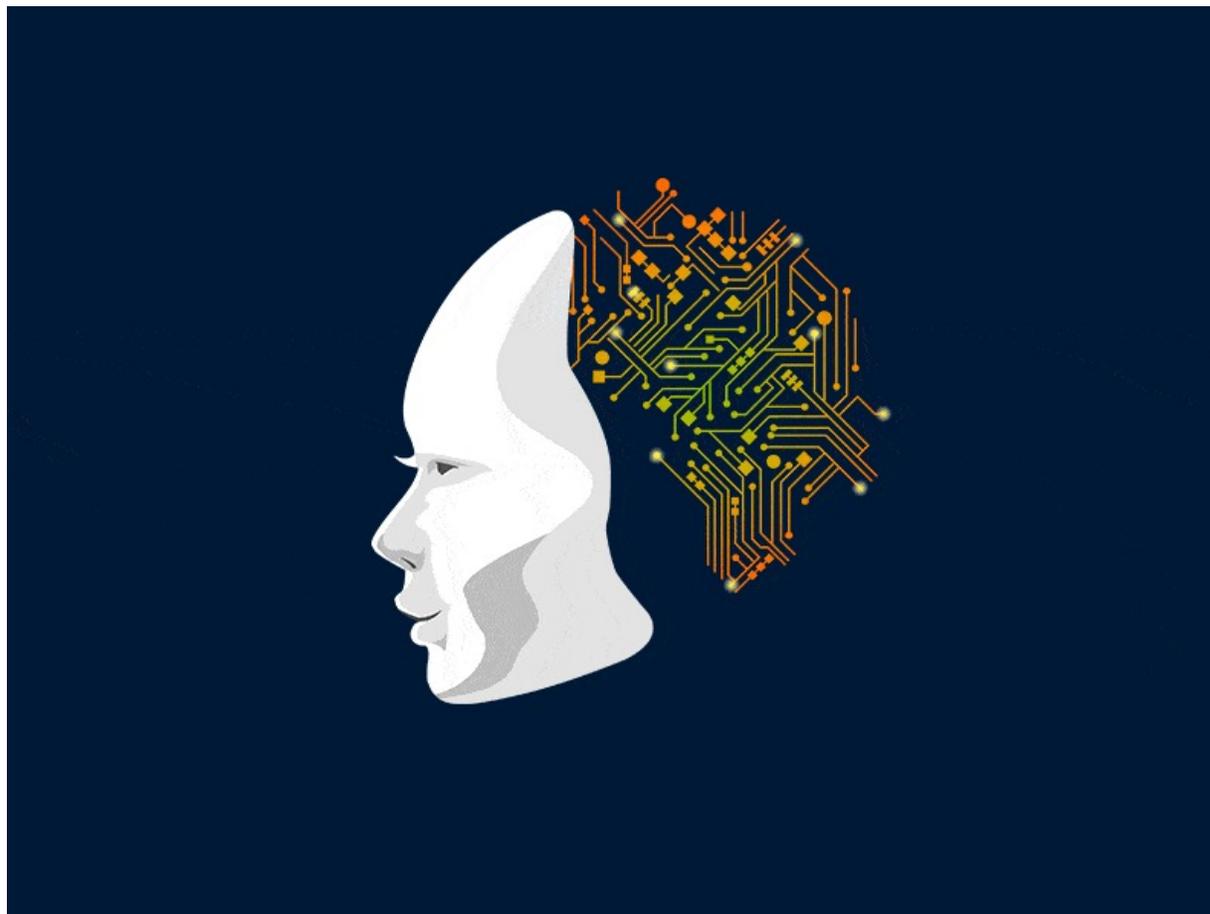
### A community effort in SARS-CoV-2 drug discovery

Johannes Schimunek, Philipp Seidl, Katarina Elez, Tim Hempel, Tuan Le, Frank Noé, Simon Olsson, Lluís Raich, Robin Winter, Hatice Gokcan, Filipp Gusev, Evgeny M. Gutkin, Olexandr Isayev, Maria G. Kurnikova, Chamali H. Narangoda, Roman Zubatyuk, Ivan P. Bosko, Konstantin V. Furs, Anna D. Karpenko, Yury V. Kornoushenko, Mikita Shuldau, Artsemi Yushkevich, Mohammed Benabderrahmane, Patrick Bousquet-Melou, Ronan Bureau, Beatrice Charton, Bertrand Cirou, Gérard Gil, William J. Allen, Suman Sirimulla, Stanley Watowich, Nick Antonopoulos, Nikolaos Epitropakis, Agamemnon Krasoulis, Vassilis Pitsikalis, Stavros Theodorakis, Igor Kozlovskii, Anton Maliutin, Alexander Medvedev, Petr Popov, Mark Zaretskii, Hamid Eghbal-zadeh, Christina Halmich, Sepp Hochreiter, Andreas Mayr, Peter Ruch, Michael Widrich, Francois Berenger, Ashutosh Kumar, Yoshihiro Yamanishi, Kam Zhang, Emmanuel Bengio, Yoshua Bengio, Moksh Jain, Maksym Korablyov, Cheng-Hao Liu, Marcous Gilles, Enrico Glaab, Kelly Barnsley, Suhasini M. Iyengar, Mary Jo Ondrechen, V. Joachim Haupt, Florian Kaiser, Michael Schroeder, Luisa Pugliese, Simone Albani, Christina Athanasiou, Andrea Beccari, Paolo Carloni, Giulia D'Arrigo, Eleonora Gianquinto, Jonas Goßen, Anton Hanke, Benjamin P. Joseph, Daria B. Kokh, Sandra Kovachka, Candida Manelfi, Goutam Mukherjee, Abraham Muñoz-Chicharro, Francesco Musiani, Ariane Nunes-Alves, Giulia Paiardi, Giulia Rossetti, S. Kashif Sadiq, Francesca Spyraakis, Carmine Talarico, Alexandros Tsengenes, Rebecca Wade, Conner Copeland, Jeremiah Gaiser, Daniel R. Olson, Amitava Roy, Vishwesh Venkatraman, Travis J. Wheeler, Haribabu Arthanari, Klara Blaschitz, Marco Cespugli, Vedat Durmaz, Konstantin Fackeldey, Patrick D. Fischer, Christoph Gorgulla, Christian Gruber, Karl Gruber, Michael Hetmann, Jamie E. Kinney, Krishna M. Padmanabha Das, Shreya Pandita, Amit Singh, Georg Steinkellner, Guilhem Tesseyre, Gerhard Wagner, Zi-Fu Wang, Ryan J. Yust, Dmitry S. Druzhilovskiy, Dmitry Filimonov, Pavel V. Pogodin, Vladimir Poroikov, Anastassia V. Rudik, Leonid A. Stolbov, Alexander V. Veselovsky, Maria De Rosa, Giada De Simone, Maria R. Gulotta, Jessica Lombino, Nedra Mekni, Ugo Perricone, Arturo Casini, Amanda Embree, D. Benjamin Gordon, David Lei, Katelin Pratt, Christopher A. Voigt, Kuang-Yu Chen, Yves Jacob, Tim Krischuns, Pierre Lafaye, Agnès Zettor, M. Luis Rodríguez, Kris M. White, Daren Fearon, Frank von Delft, Martin A. Walsh, Dragos Horvath, Charles L. Brooks, Babak Falsafi, Bryan Ford, Adolfo García-Sastre, Sang Yup Lee, Nadia Naffakh, Alexander M. M. Guenther Klambauer, Thomas M. Hermans  ... See fewer authors ^

# CONCLUSIONI

- AI è uno strumento sempre più potente...e come tale va gestito correttamente
- AI come strumento di collaborazioni multidisciplinari tra ricercatori in ambito medico scientifico ed esperti del settore informatico
- AI NON è una minaccia per il lavoro, NON rimpiazza gli umani, ma richiede competenze diverse dal passato
- Come qualsiasi novità AI non va temuta, ma va conosciuta e se gestita può portare ad un miglioramento globale!





GRAZIE PER LA VOSTRA ATTENZIONE

Per info o collaborazioni  
[uperricone@fondazionerimed.com](mailto:uperricone@fondazionerimed.com)