



universität
wien

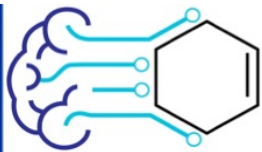
Bologna – Industrial Chemistry Department “Toso Montanari”
02.12.2022

*Drug Discovery and Cheminformatics:
discovering new drugs in the
Big Data era*



Vincenzo Palmacci

PhD Candidate MSCA



My studies at UNIBO

BSc Industrial Chemistry



- Organic synthesis
- Physical Chemistry
- Industrial chemistry

2013-2018

Bachelor thesis
Oct 2016 – Nov 2017

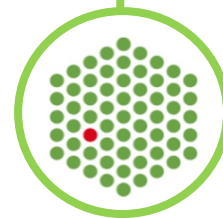
MSc Bioinformatics



- Protein design
- Analysis of Biological data
- Machine Learning

2018-2020

Master thesis
Erasmus + Internship



European Molecular
Biology Laboratory

BIOSAXS

Artificial neural networks for protein's
shape prediction

After UNIBO: Advanced Machine Learning for Drug Discovery



Machine Learning Research

- Artificial Intelligence for Drug Discovery
- Representation Learning
- Explainable AI



COMP3D group

- Computational Drug Design
- Prediction of assay interference

2021 - 2024?

Molecular AI



- Retrosynthesis planning
- De novo molecule generation

Today's Lecture

Drug Discovery

- The Drug Discovery Pipeline
 - Early Drug Discovery

Cheminformatics

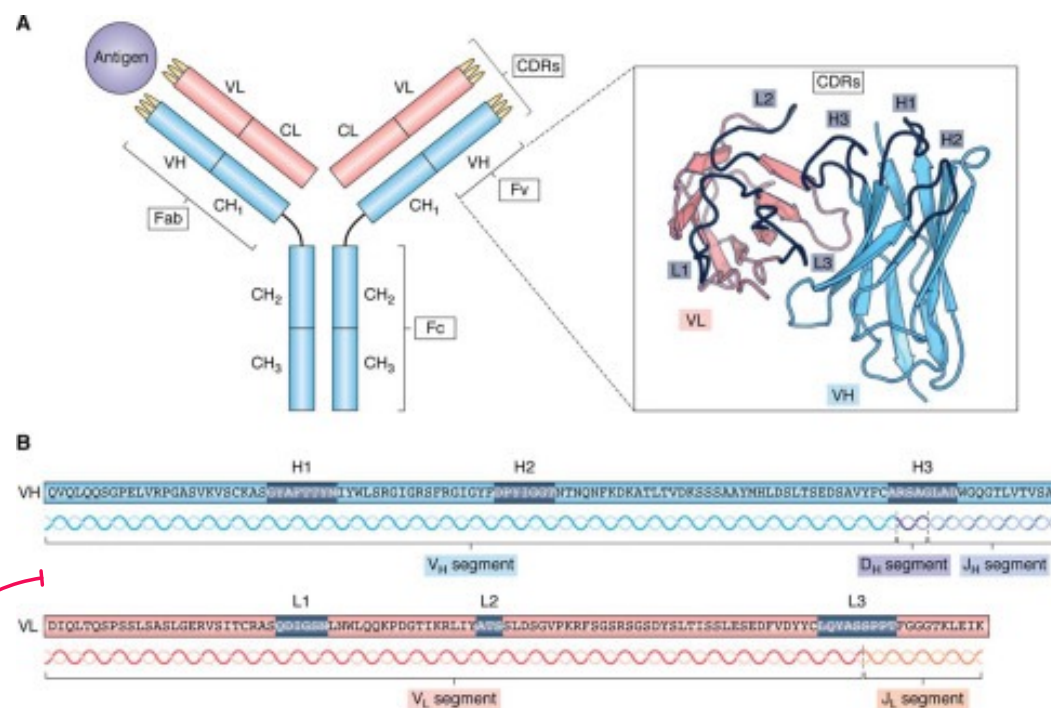
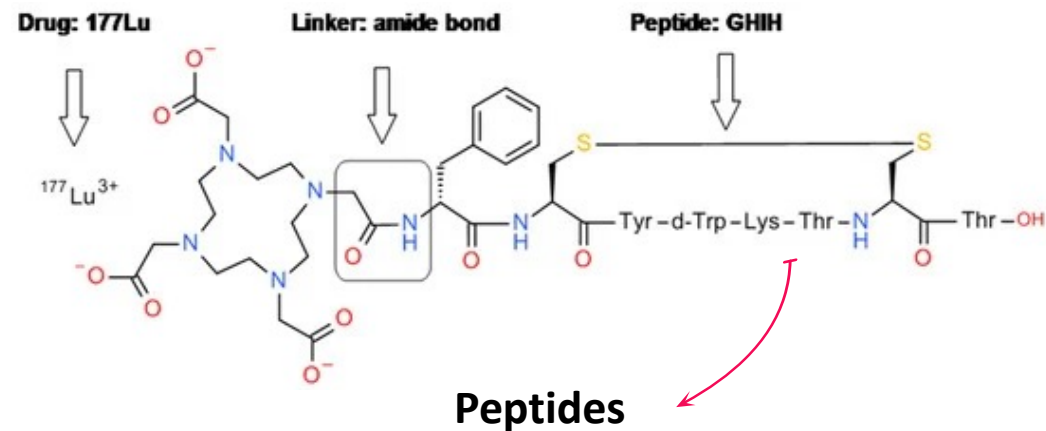
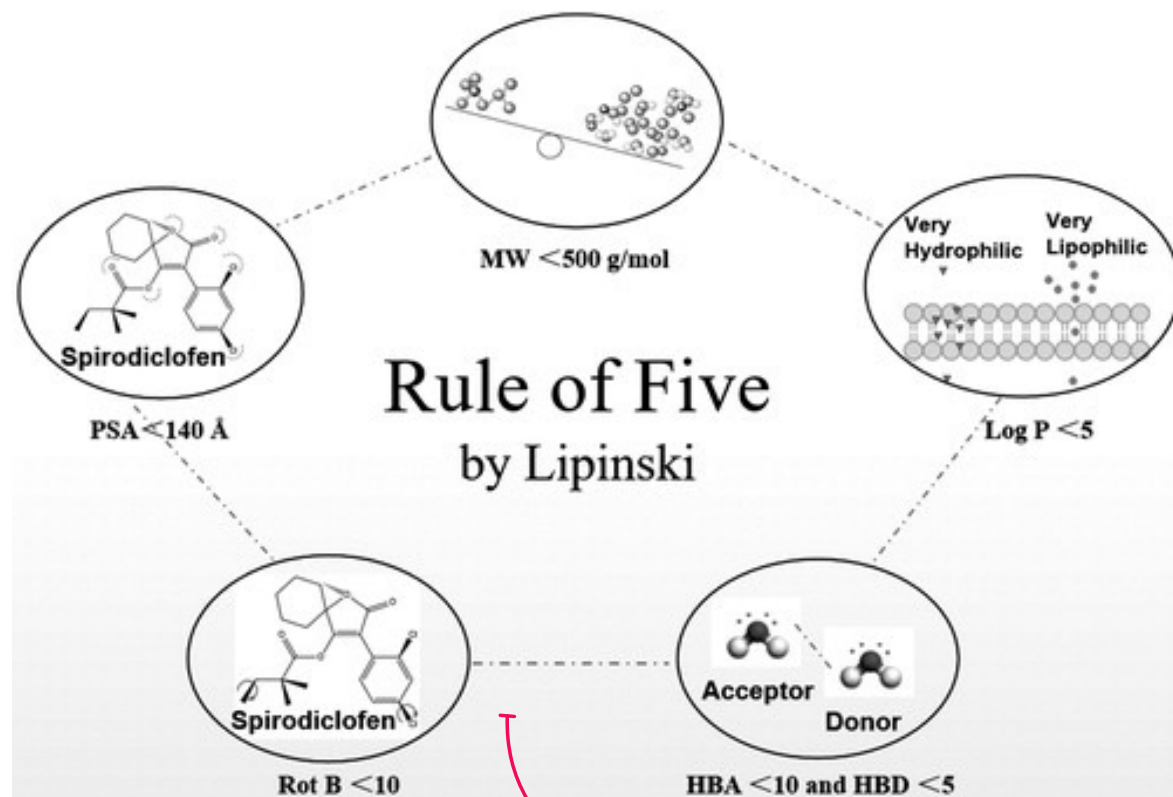
- Big Data in Chemistry
- Introduction to Cheminformatics: meaning and fundamentals
- Introduction to Machine Learning and Deep Learning
- Example: Toxicity prediction

Computational Drug Discovery

- Cool stuff: How Artificial Intelligence helps Drug Discovery



What is a drug?

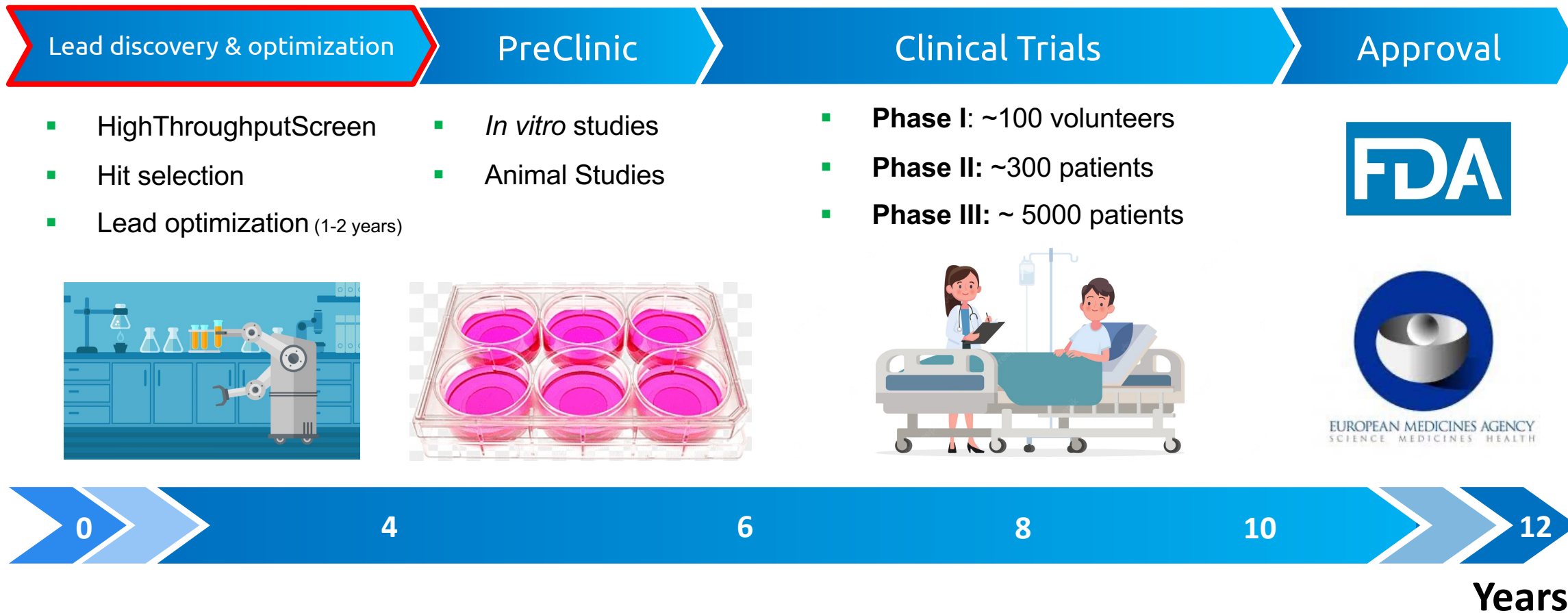


Antibodies

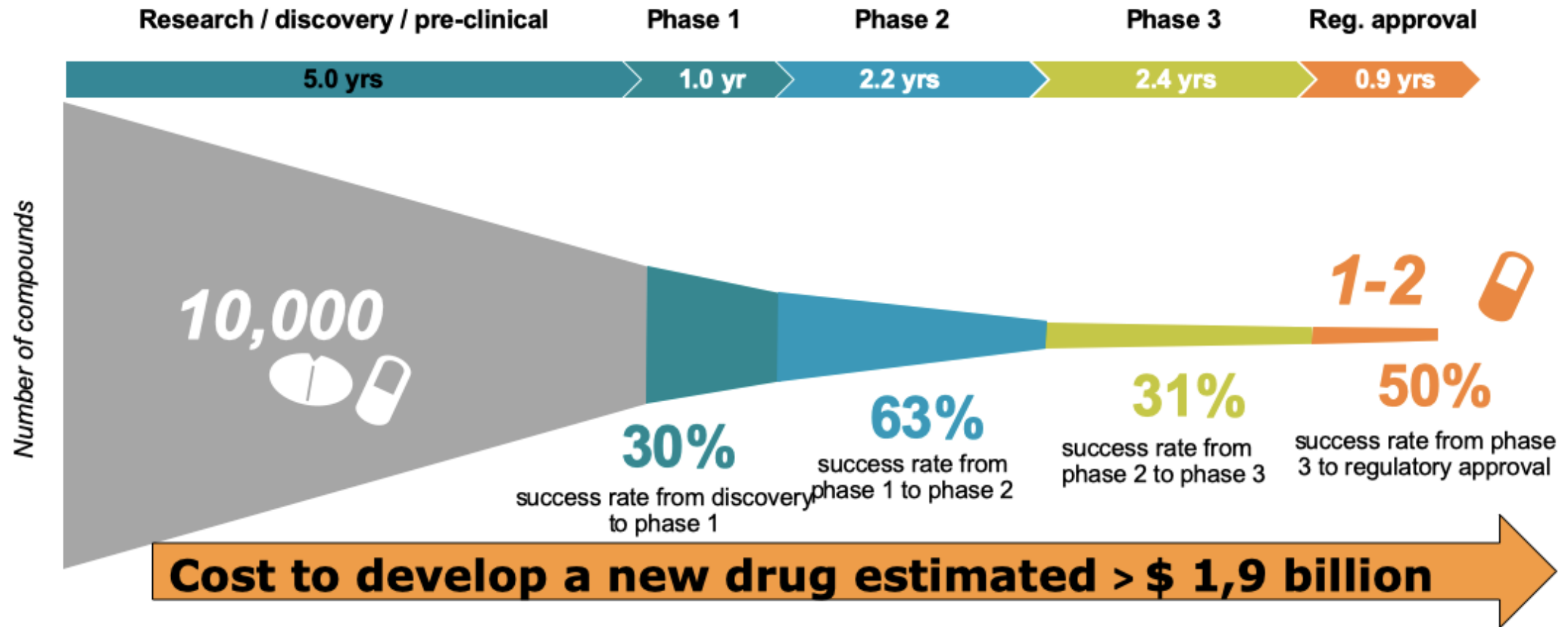
Drug Discovery: from target identification to drug approval

Start: 5 Millions molecules

End: 1 molecule



Pipeline steps success rate



Target identification and validation

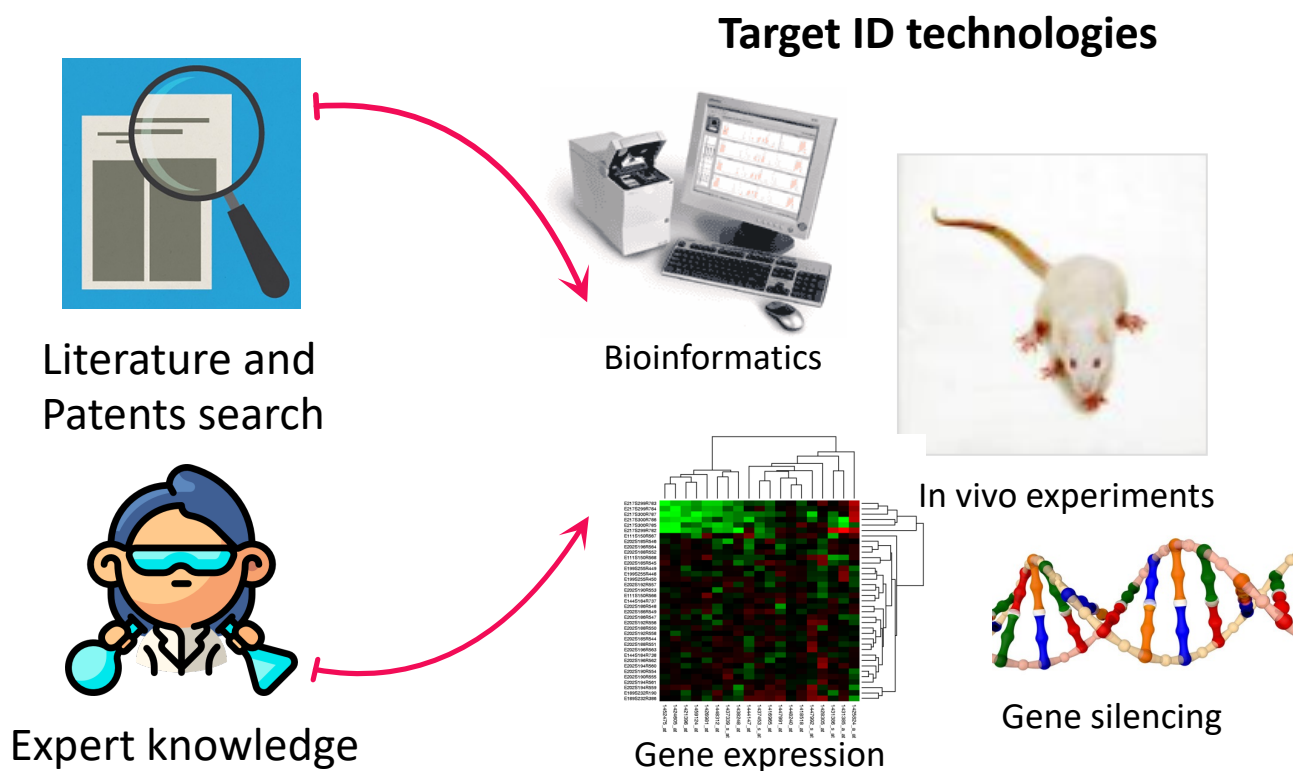
What is a Target?



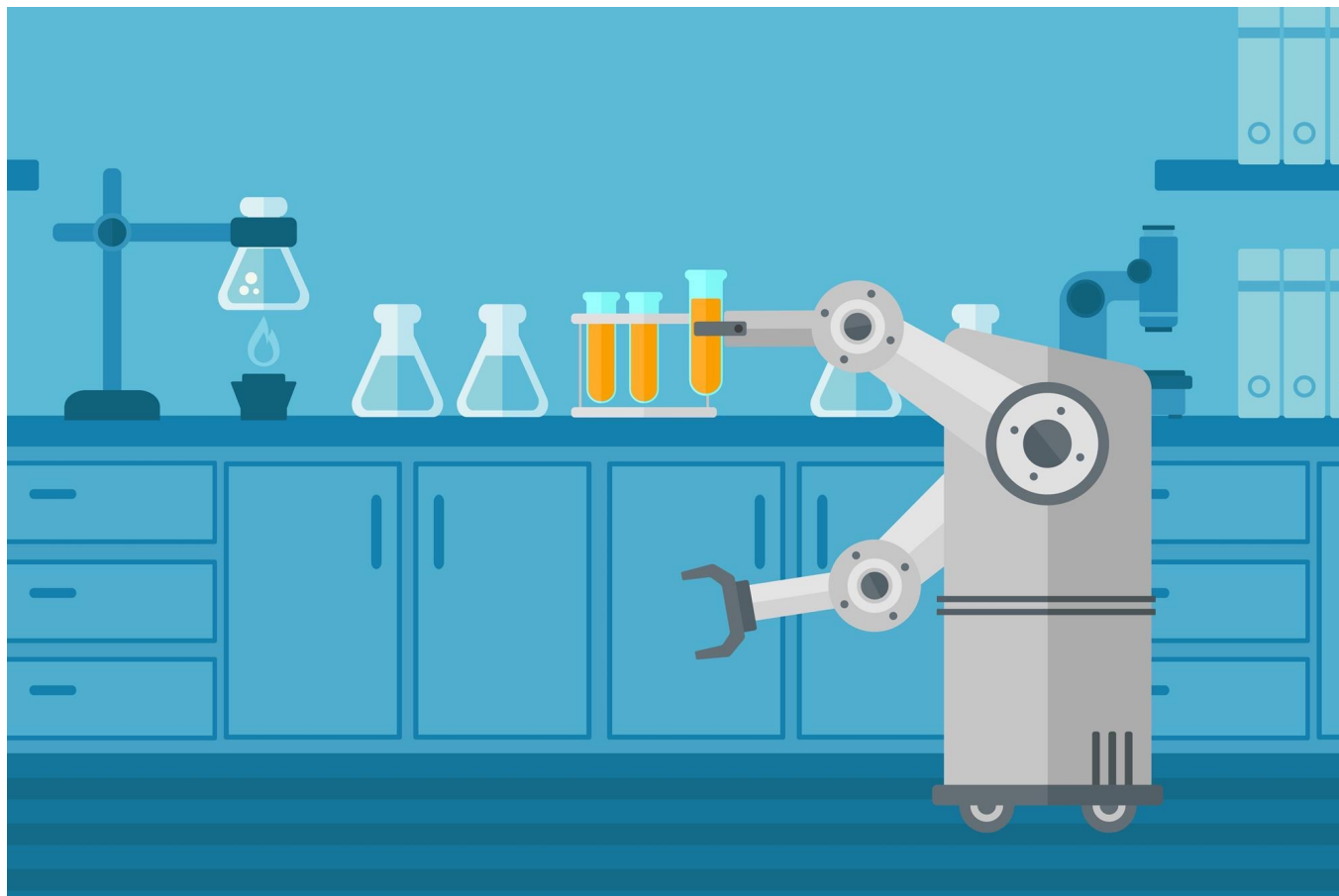
PDB ID: 1C3S

- **Biological target:** anything in a living organism that changes behaviour upon binding with other entities
- Common targets: Proteins, RNA or DNA
- Must have a role in the development of a disease

Identification and validation



Lead Discovery: High Throughput screening (HTS)



High Troughput Screening (HTS):

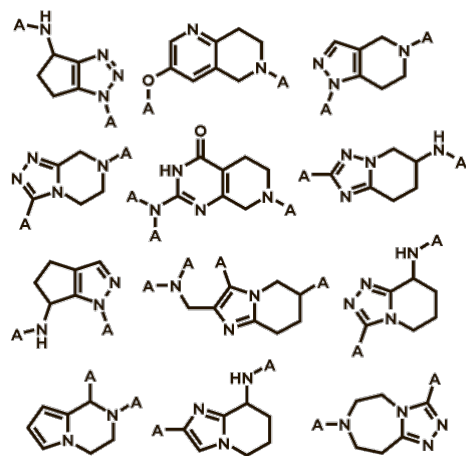
- Key step in the early Drug Discovery pipeline
- Main technology for hits identification
- At full capacity hundreds of thousands of compounds tested daily

High Troughput Screening advantages:

- Fully automatized workflow
- Extremely cost effective
- Suggest good starting structures that will be optimized in later steps

HTS: A closer look

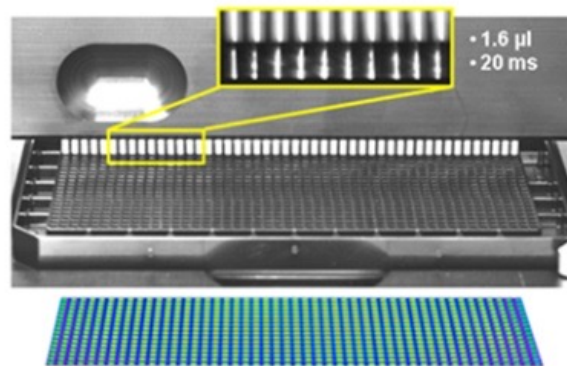
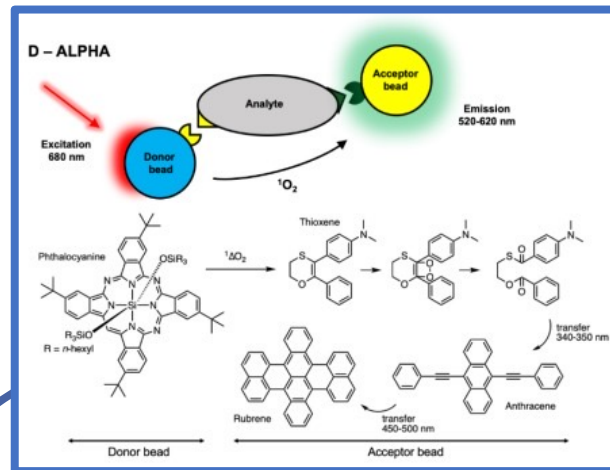
Biological assay



Compound Library:

Set of chemical substances to test.
Each pharma company has its own library which have been developed throughout the year

Millions of Molecules!!!



Hits evaluation: find the **LEADS**

- Look at compounds with positive signal (usually ~10% of the library)
- Compare active compounds with results of older assays
- Select most promising hits

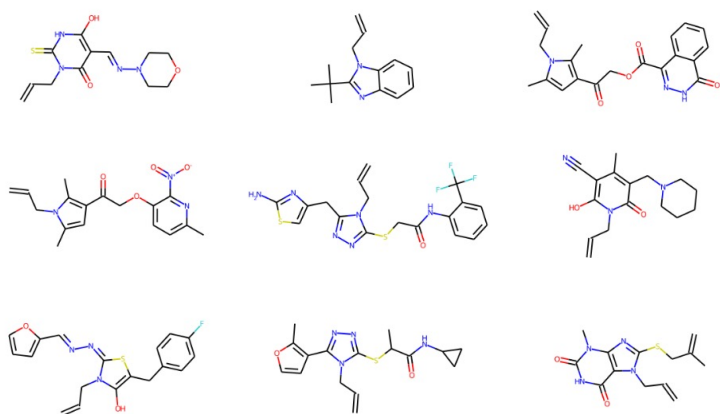


Expert knowledge



Cheminformatics

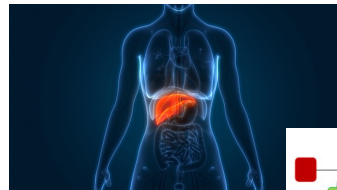
Lead Discovery: Lead optimization



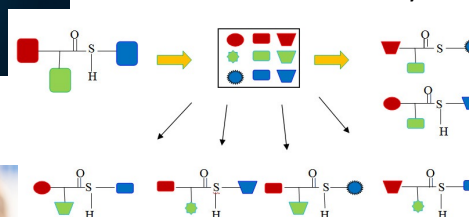
List of promising hits*

Molecules with showing bioactivity in HTS phase that passed statistical and expert-rule filters.

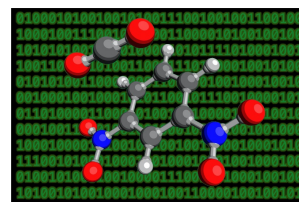
ADMET prediction



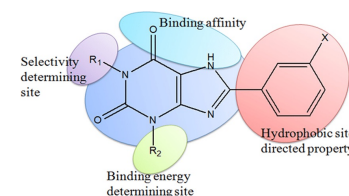
Combinatorial Chemistry



Medicinal Chemistry



Computational Chemistry



Structure Activity Relationship

**~300
molecules**

Optimized molecules

- Improved potency
- Better stability
- Increased solubility in water
- Greater selectivity

Optimization

Thousands of variations generated by classical lab chemists and automated synthesis

PreClinic: *In-vitro* and Animal studies

In vivo studies

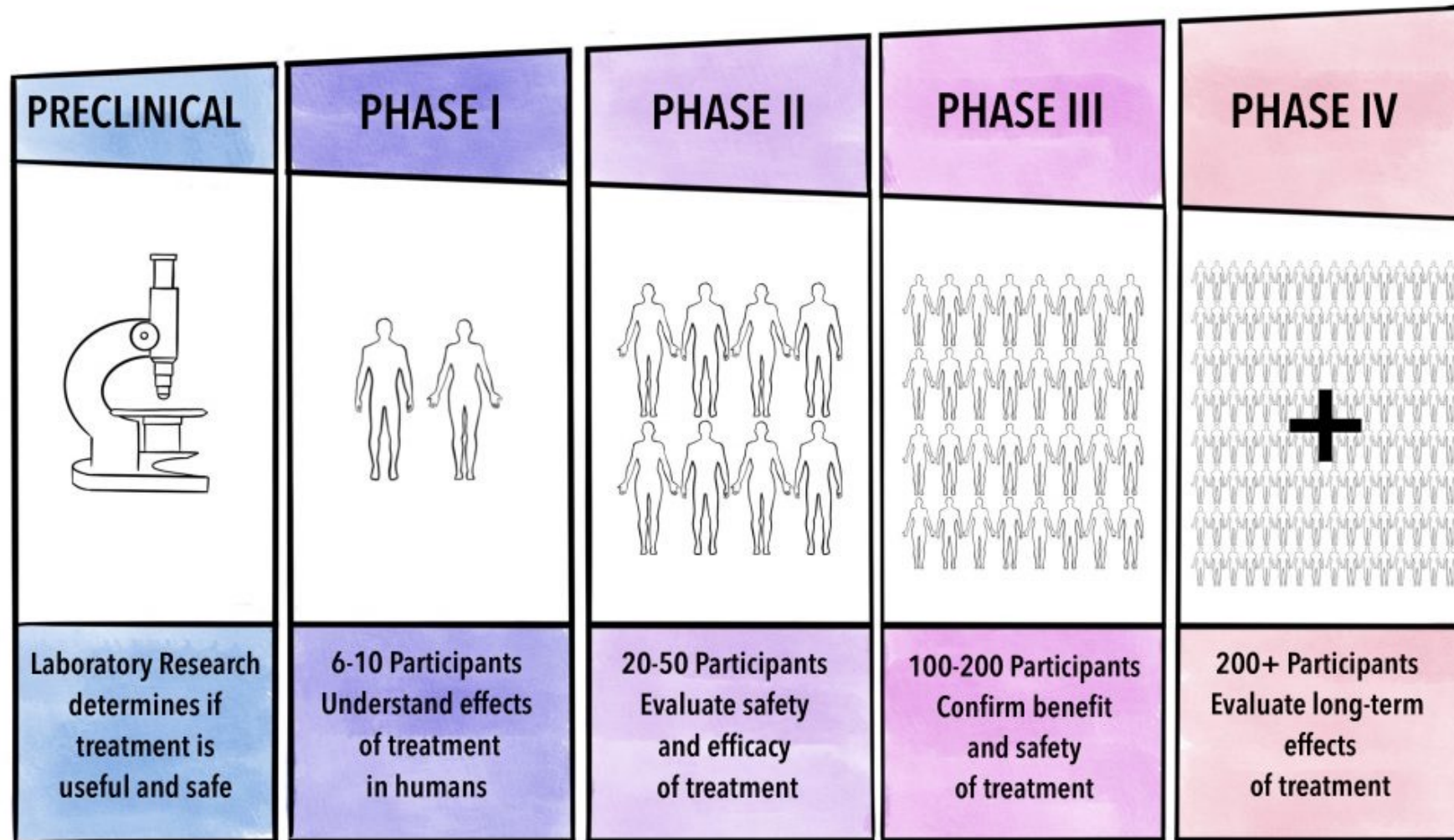


In vitro studies



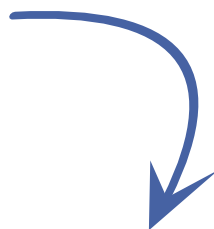
- **Pharmacodynamics:** What the drug does in the body?
- **Pharmacokinetics:** What the body does to the drug?
- **ADME and Toxicology**
- Estimate a safe starting dose for clinical trials in humans

Clinical Trials: Drug safety and effectiveness



The Drug Discovery pipeline: RECAP

- Long and expensive process: 12 years and ca. 2 Billion \$ invested
- Involves scientists from all the fields: Chemists (medicinal, organic, analytic, etc...), Biologists, Mathematician/Physicists
- Lot of improvement possible:
 - Speed up the discovery of new drugs
 - Increase the success rate of drug candidates



How? With Chemoinformatics!!!
(and Machine Learning)

“Big Data” in Chemistry



Early drug discovery

Huge amount of **noisy** data

- HTS: millions of chemicals with preliminary activity values on multiple targets.
- Lead optimization: thousands of molecules with IC_{50} values and ADMET properties



PreClinic and Clinical Trials

Few, very **reliable** data

- Toxicity studies: hundreds of compounds with known toxicity
- Pharmacodynamics: known metabolism in humans

“Big Data” in Chemistry: public data sources and databases



<https://pubchem.ncbi.nlm.nih.gov/>

- Large collection of chemical informations:
 - compounds
 - biological assays
- **112 Millions** compounds
- **301 Millions** Bioactivity
- **42 Millions** Patents



<https://go.drugbank.com/>

- Online database containing information on drugs and drug targets
- **15 Thousands** drugs:
 - **5 K** approved
 - **7 K** experimental

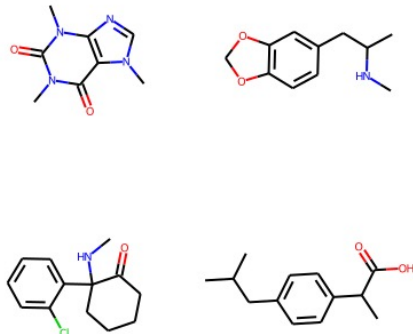


<https://www.ebi.ac.uk/chembl/>

- Manually curated database of bioactive molecules:
 - In vitro and in vivo assays
- **2.3 Millions** Compounds
- **15 K** Assays

What's the deal with all those data?

Given a set of new molecules:

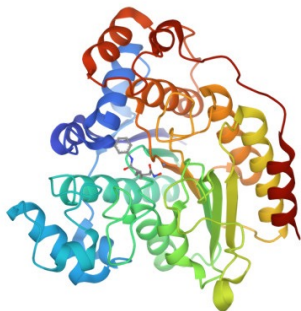


30L to everyone that can name all those molecules ;)

Knowledge from old experiments
can be used for:

- Predict which molecules can be bioactive
- Predict which molecules can be toxic
- Predict the affinity toward a certain protein

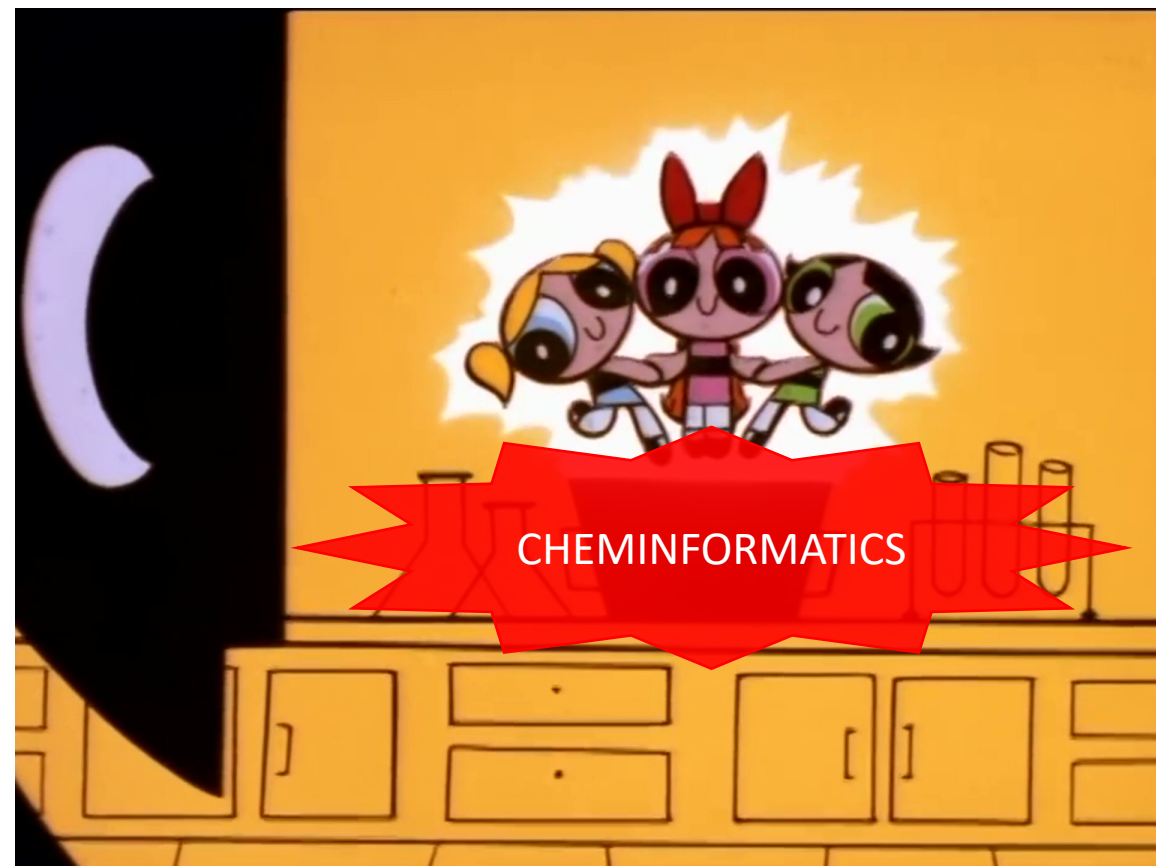
Given a protein target:



- Generate a set of structures with high affinity

How? With **Chemoinformatics!!!**
(and Machine Learning)

Cheminformatics: Chemistry as data science



Cheminformatics: What it's done in practice?

In a nutshell: with Cheminformatics we try to exploit all the available chemical knowledge to help experimental

Chemists:

- Visualize large amounts of different compounds in an informative way (e.g.: colored by bioactivity)

EXAMPLE: <https://peter-ertl.com/molecular/rings/magicrings.html>

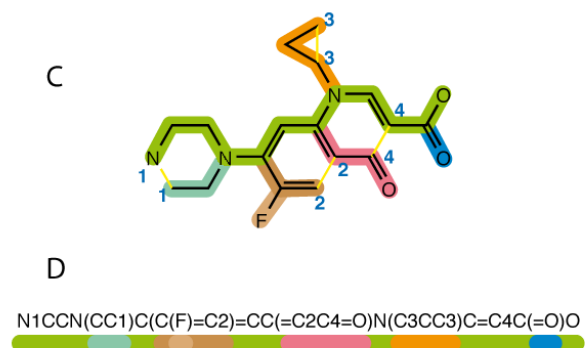
- Develop algorithms that can use the data to predict molecular properties
- Develop algorithms that can use the data to help with synthesis (e.g.: prediction of synthetic routes)
- Find new informative ways to represent chemical compounds

Molecular representations: SMILES, Fingerprints, 2D descriptors

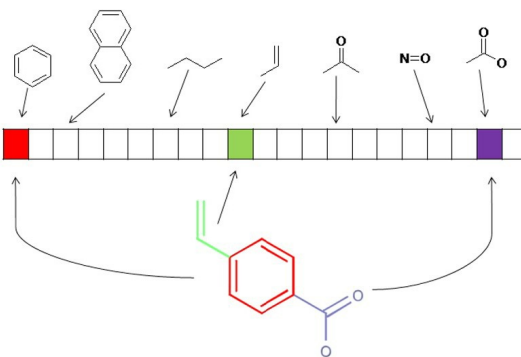
Molecular representations

- Need a computer readable representation of molecules
- Molecular representations must describe:
 - different types of structures (small molecules, peptides, polymers)
 - with different properties (stereochemistry, valence)
- Precise representations for specific structures are needed to optimize the process of AI-driven discovery

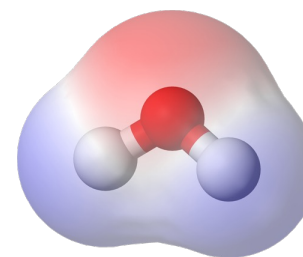
SMILES



Fingerprints



2D descriptors

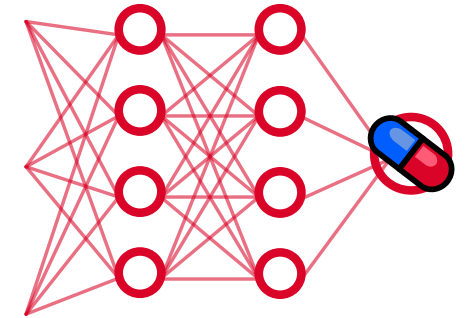
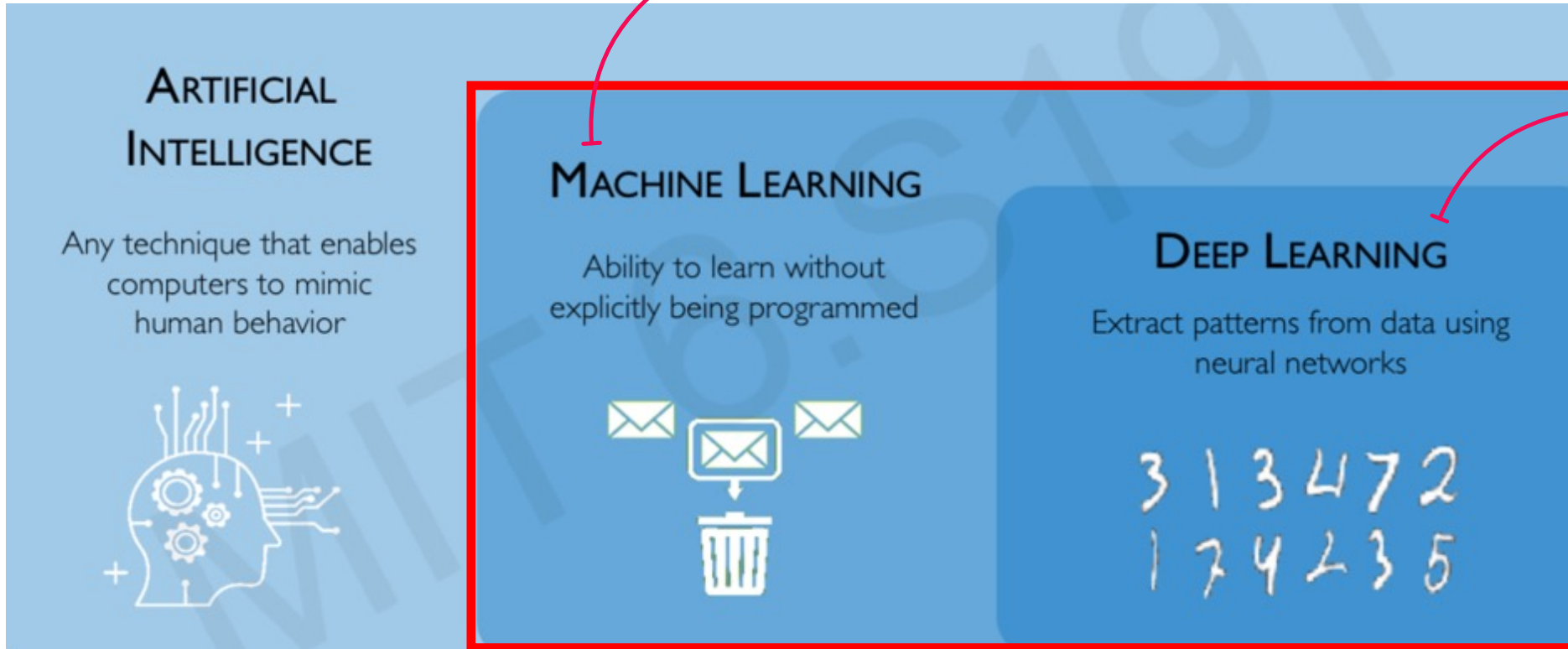


Xz	LogP
Yz	MW
Pz	Charge
Kz	TPSA
Hz	# Rings

Properties vector

A hint of Machine Learning and Deep Learning

- Extract common rules that can explain a dataset.
- Rules: set of characteristics that are connected to a property

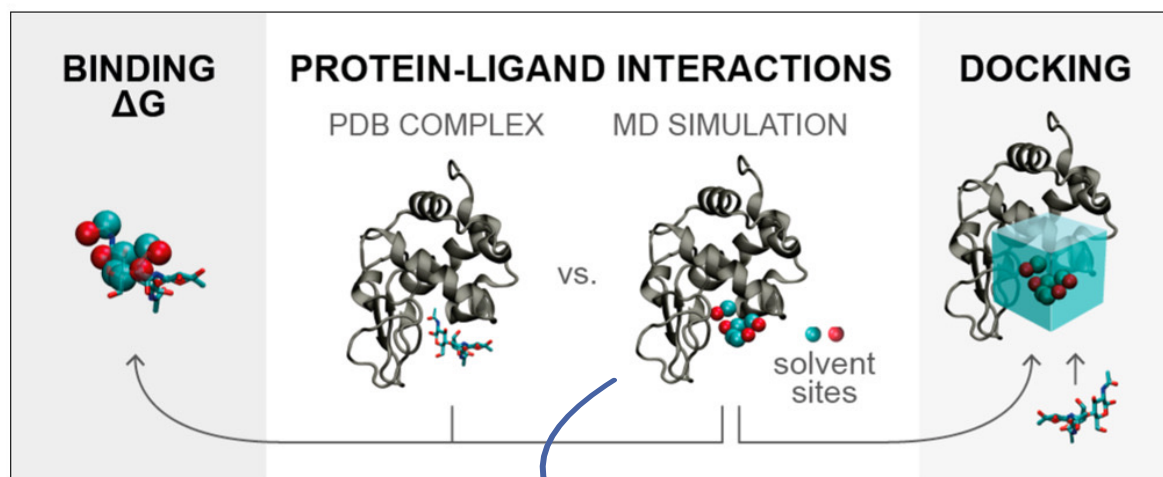


- Special type of algorithms called Artificial Neural Networks

Why not using standard Computational Chemistry?

Note: Results from Computational Chemistry experiments can be used as good starting points for modern approaches

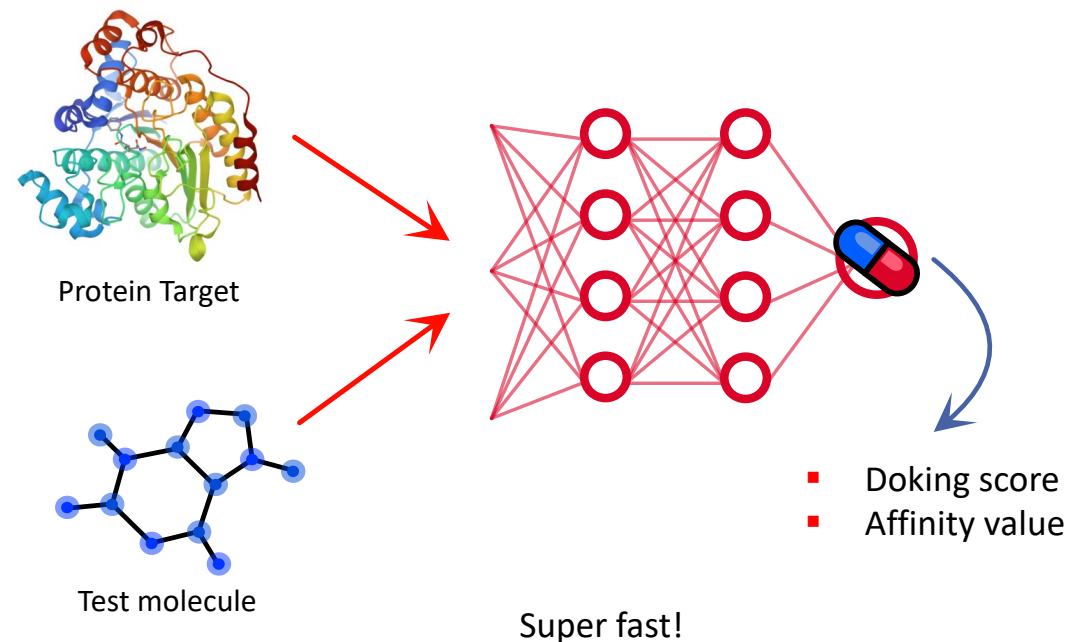
Classic computational chemistry approach



Extremely slow!

Days, even weeks depending on the initial conditions

Machine Learning approach



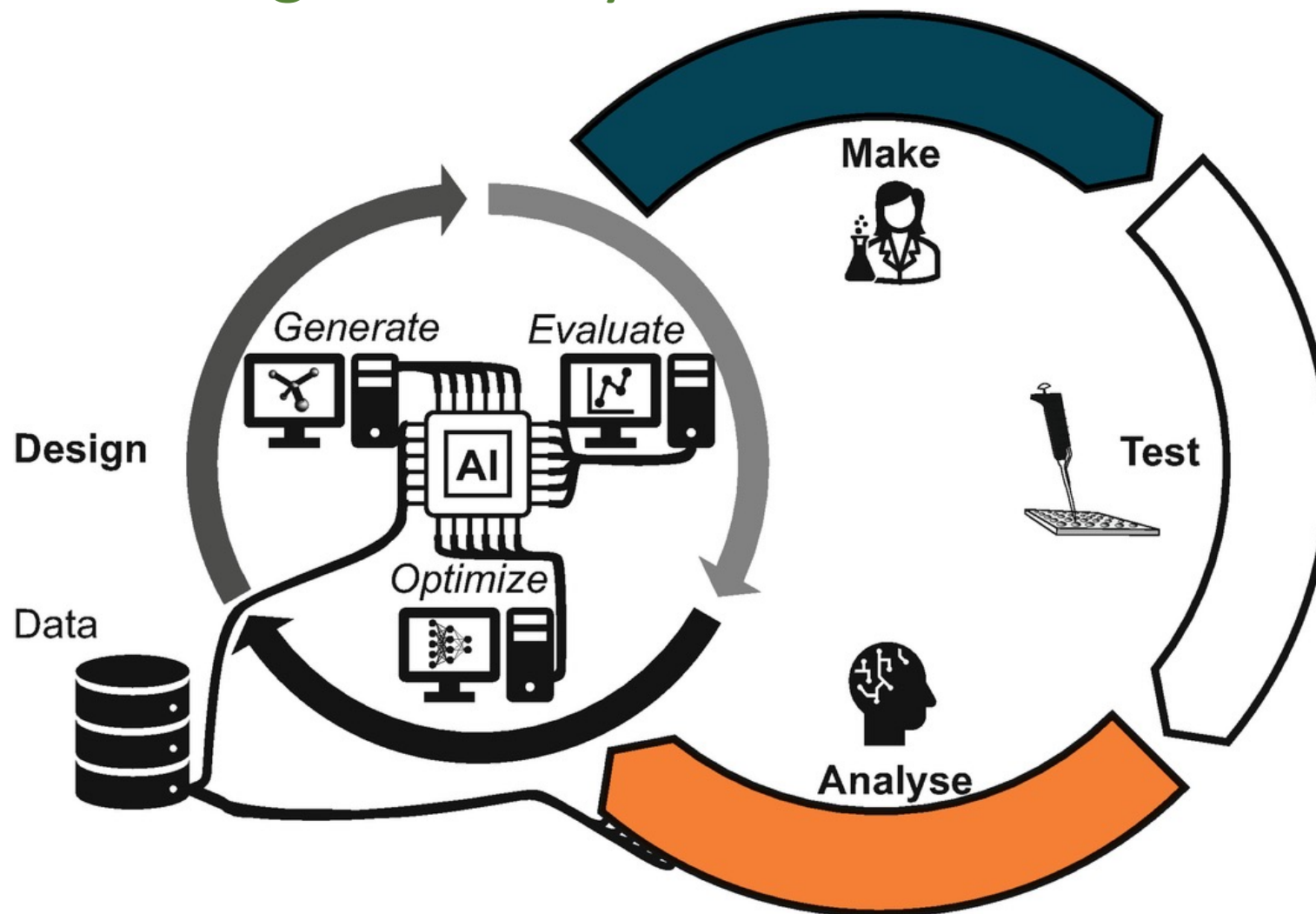
Super fast!

After training results for new molecules can be obtained in fractions of seconds

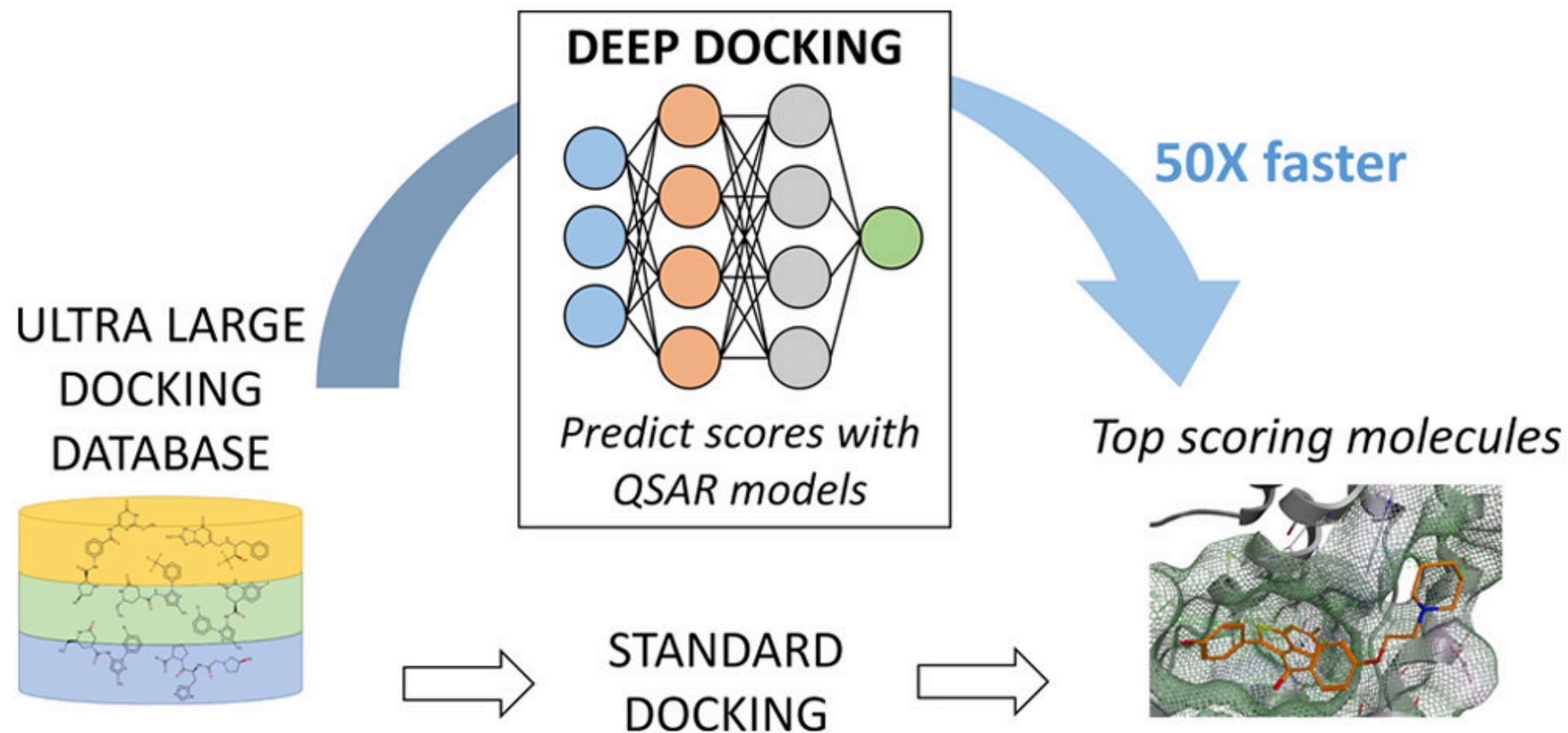
Cheminformatics: RECAP

- Goal: helping experimentalists with data science
- Very requested professional profile
- Requires knowledge from different fields: Chemistry, Data Science, Computer Science, lots of Patience
- Still room for improvement and many unsolved problems: great research field!
- Requires meaningful molecular representations
- Use of Machine Learning and Deep Learning: nowadays the two most looked forward technologies in Life Science

Computer aided Drug Discovery

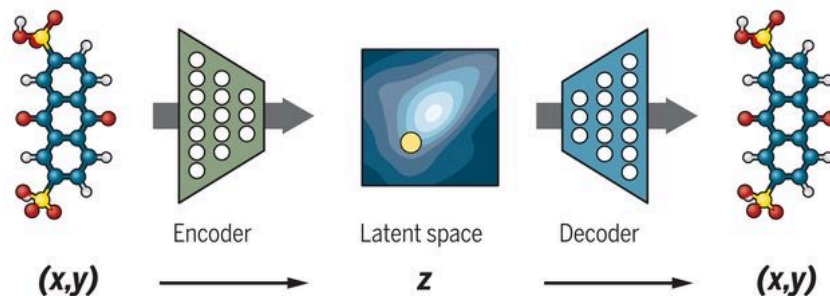


Molecular docking

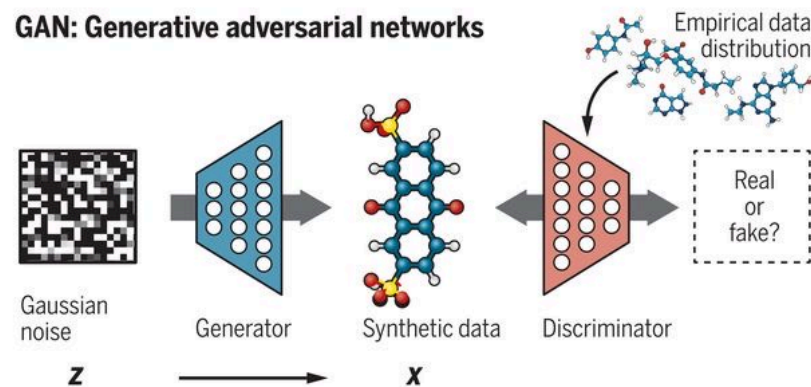


De novo molecules generation

VAE: Variational autoencoders

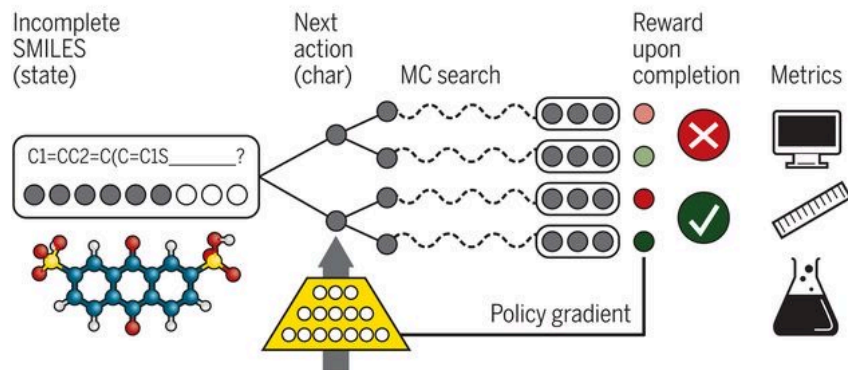


GAN: Generative adversarial networks

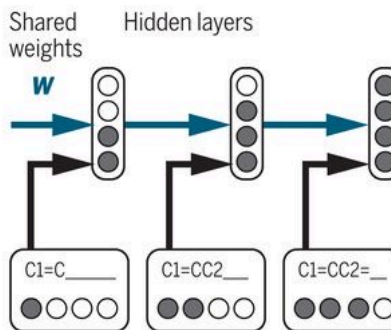


RL: Reinforcement learning

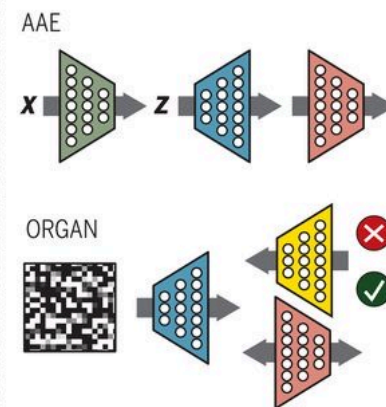
Policy gradient with Monte Carlo tree search (MCTS)



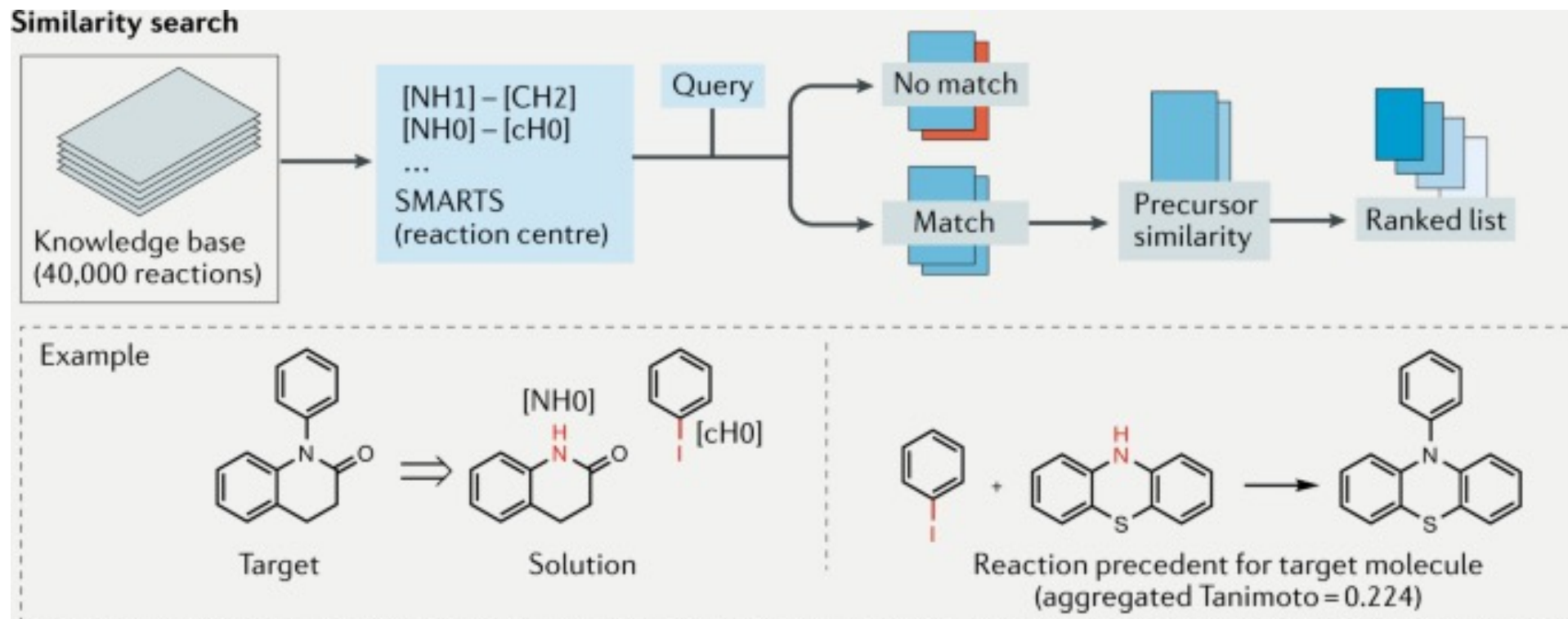
RNN: Recurrent neural network



Hybrid approaches



Retrosynthesis planning and synthesis prediction

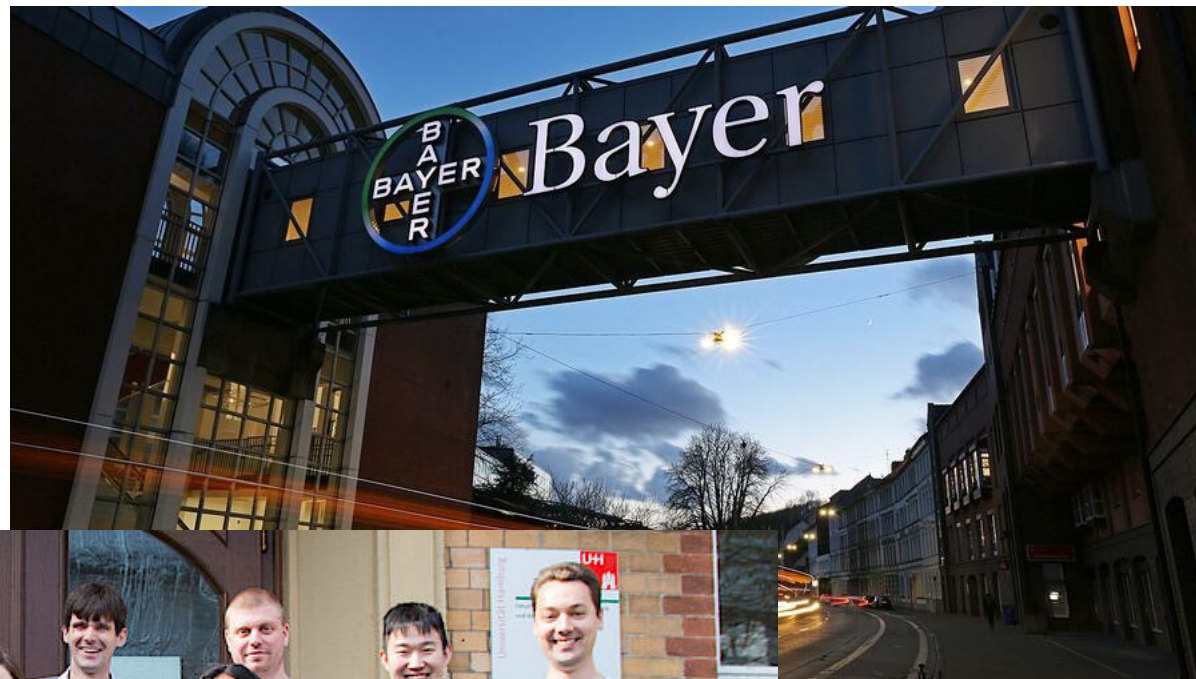


Thanks!

<https://github.com/bayer-science-for-a-better-life>



<https://ai-dd.eu/>



<https://comp3d.univie.ac.at/>

Take home message

- Drug Discovery is very time expensive and money consuming but opens up many opportunities to have an impact on society with science
- Standard laboratory Chemistry in 2022 requires computational approaches to be efficient
- Data Science tools allows a better understanding of Chemistry in the Drug Discovery context
- My suggestion: have a look to Deep Learning and Machine Learning
- If you want to do research choose a topic that fascinates you and have fun with it
- Cheminformatics starter pack: RdKit (<https://www.rdkit.org/>), Python (<https://www.python.org/>)