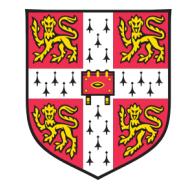


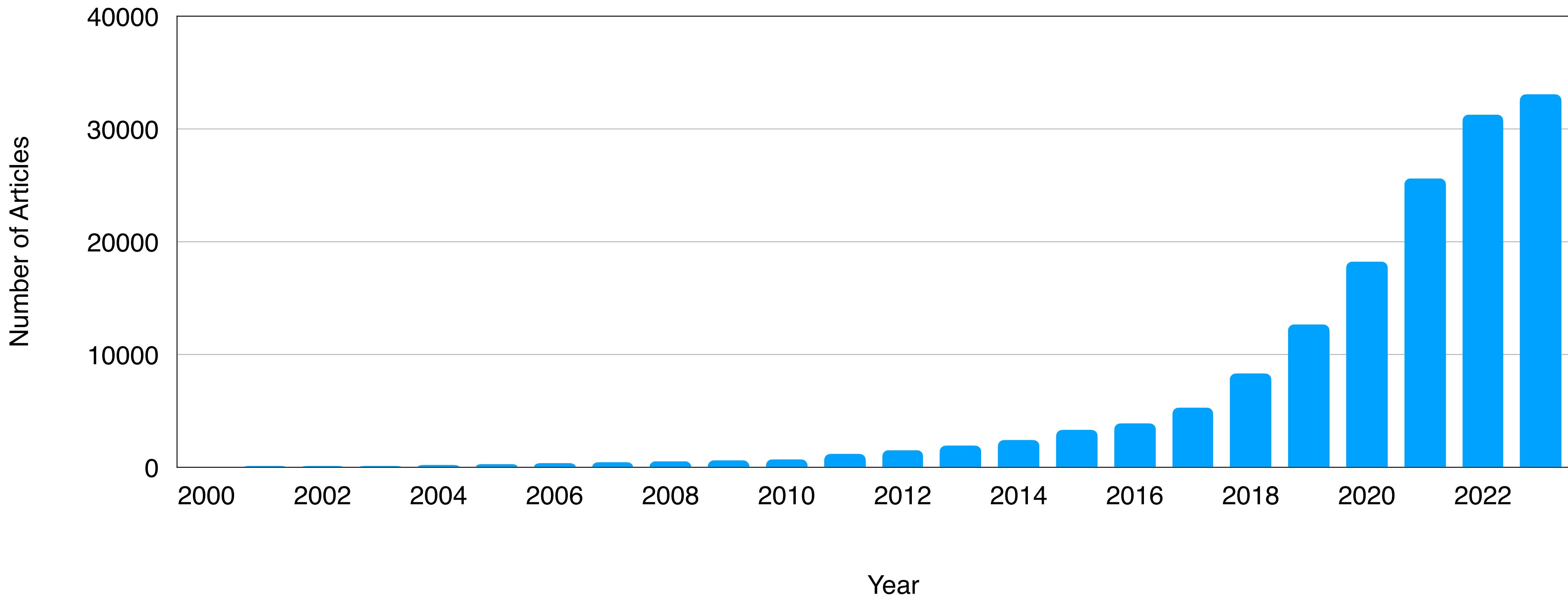
Machine Learning for Chemists

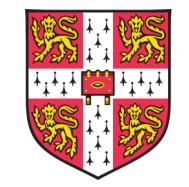
– Introduction –

5 February 2024

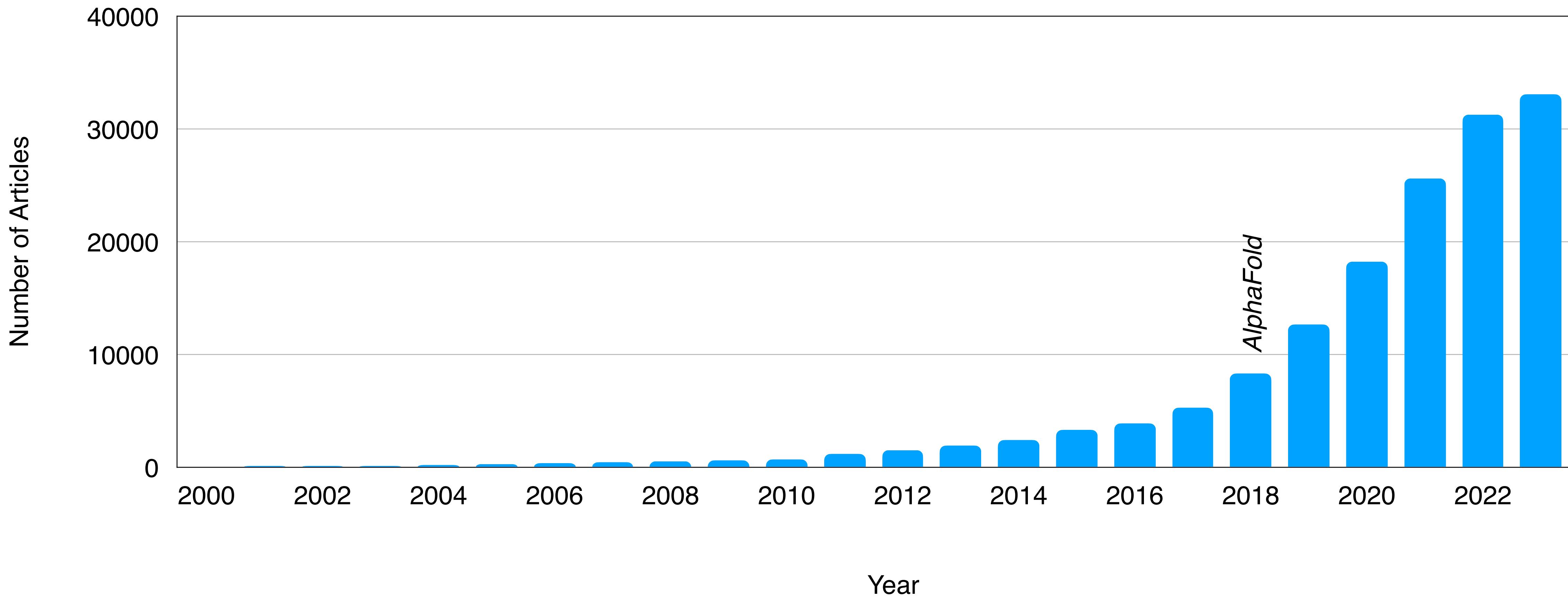


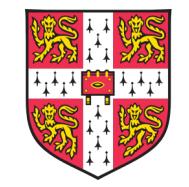
“Machine Learning” PubMed Articles



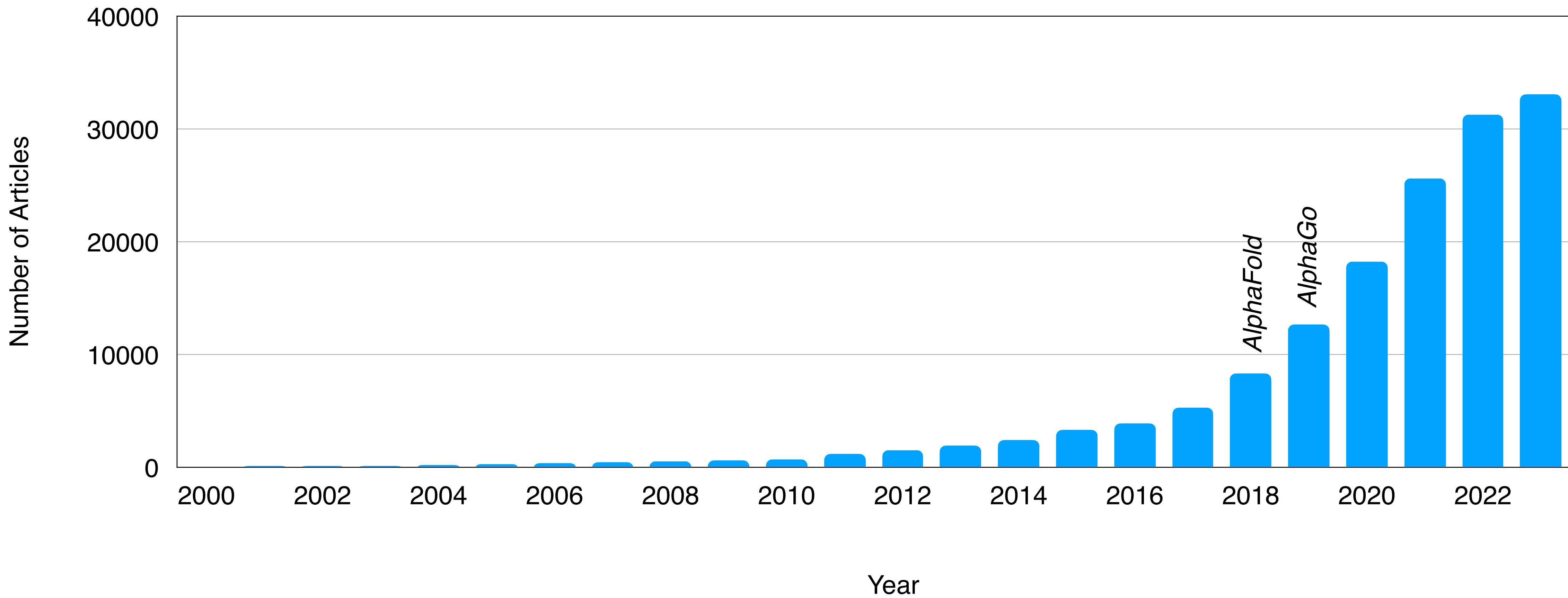


“Machine Learning” PubMed Articles



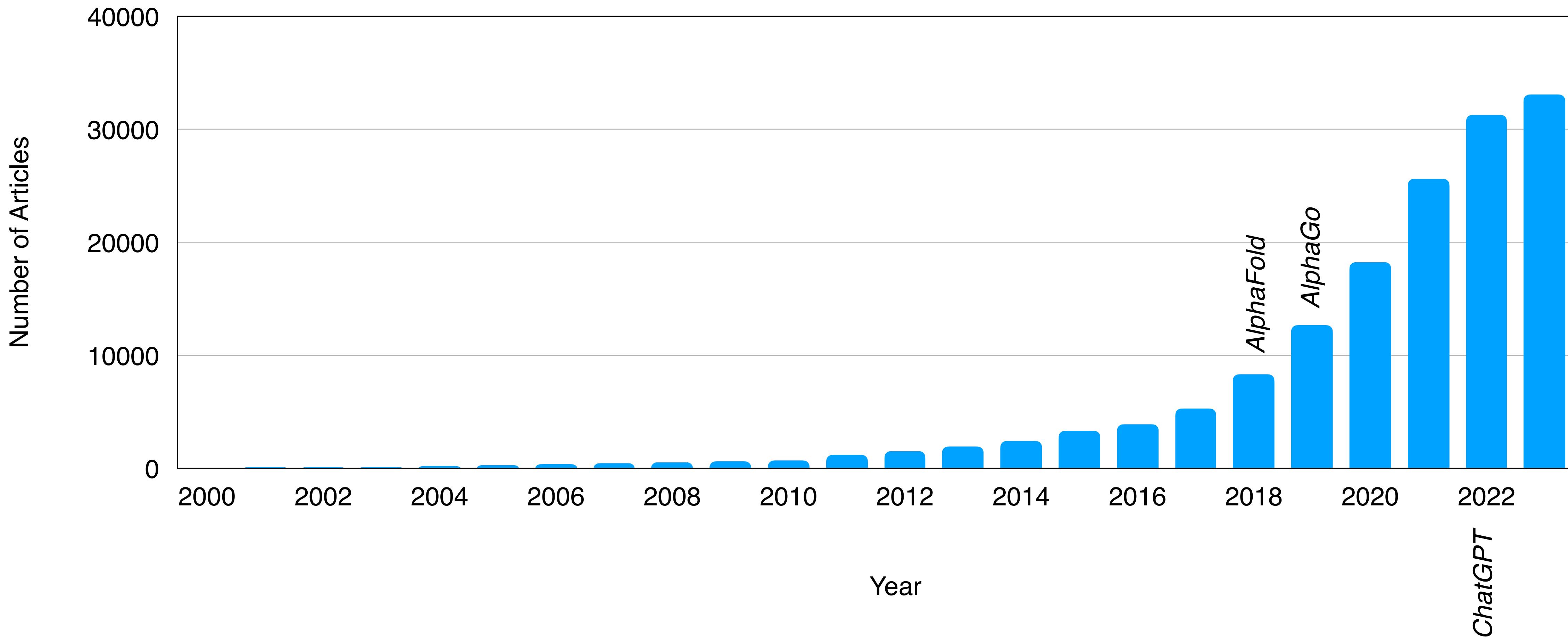


“Machine Learning” PubMed Articles





“Machine Learning” PubMed Articles

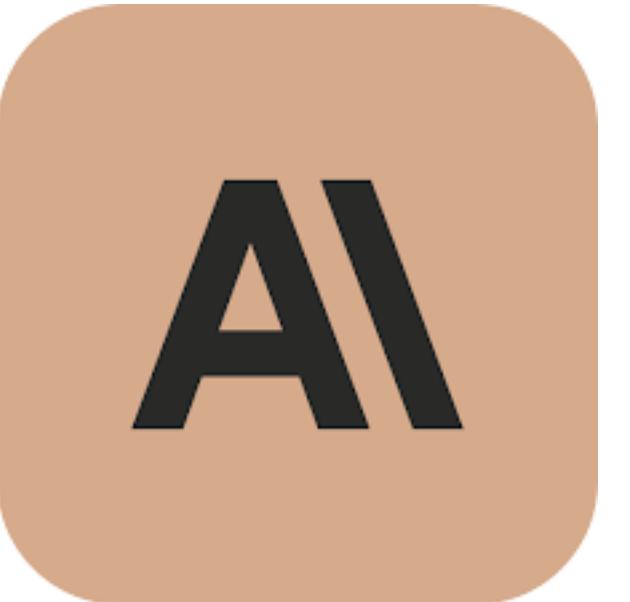


The Current New Toys



ChatGPT

OpenAI



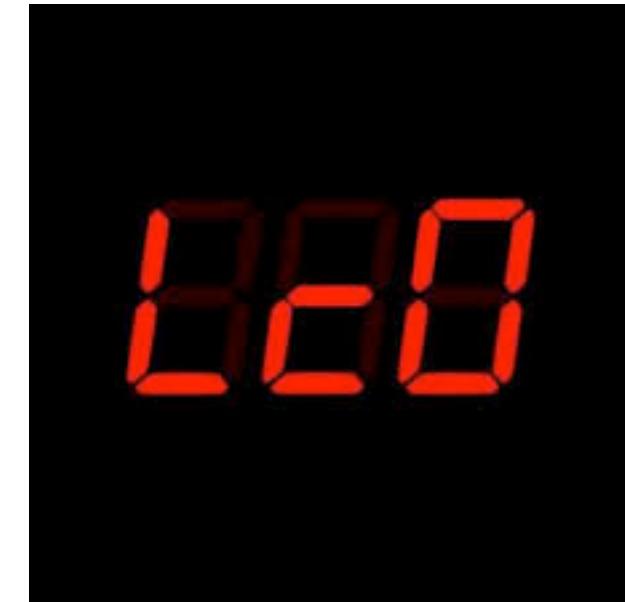
Claude

Anthropic



Midjourney

Midjourney, Inc.



Leela Chess Zero

Based on Deepmind's AlphaGo Zero

Open-source (free to use)



“Big picture” reasoning behind common ML techniques

PCA, SVM, Random Forest, Gaussian Process, Feed Forward Neural Networks

Basic steps in running the algorithms

Some helpful debugging tips

Feel free to email or come with questions!



Python - it's flaws and intricacies

Rigorous mathematical proofs behind the ML techniques

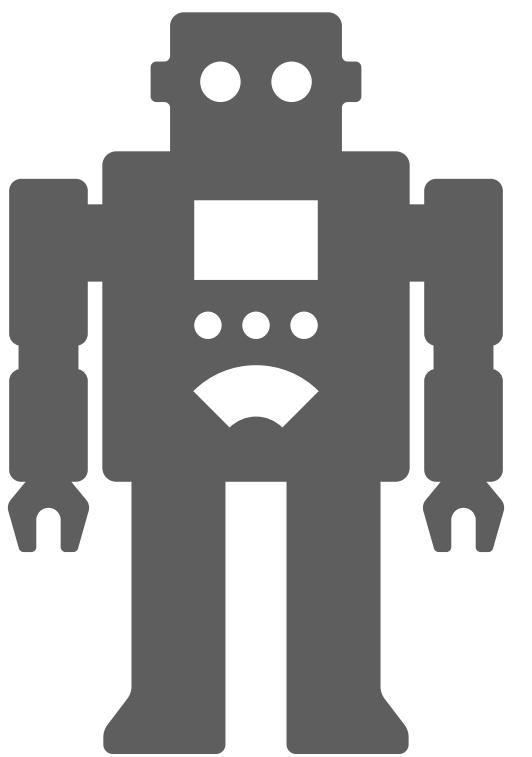
A comprehensive review of ML techniques

Feel free to ask questions about things not covered!

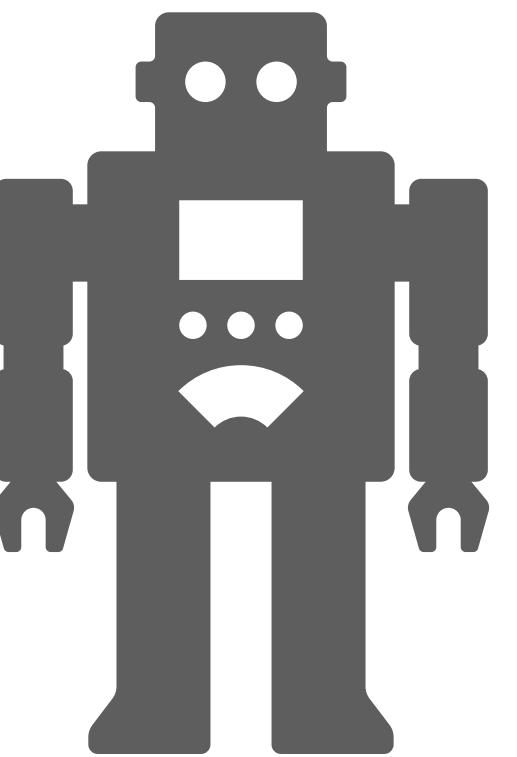
This is how we represent ML:



This is how we **should** represent ML:

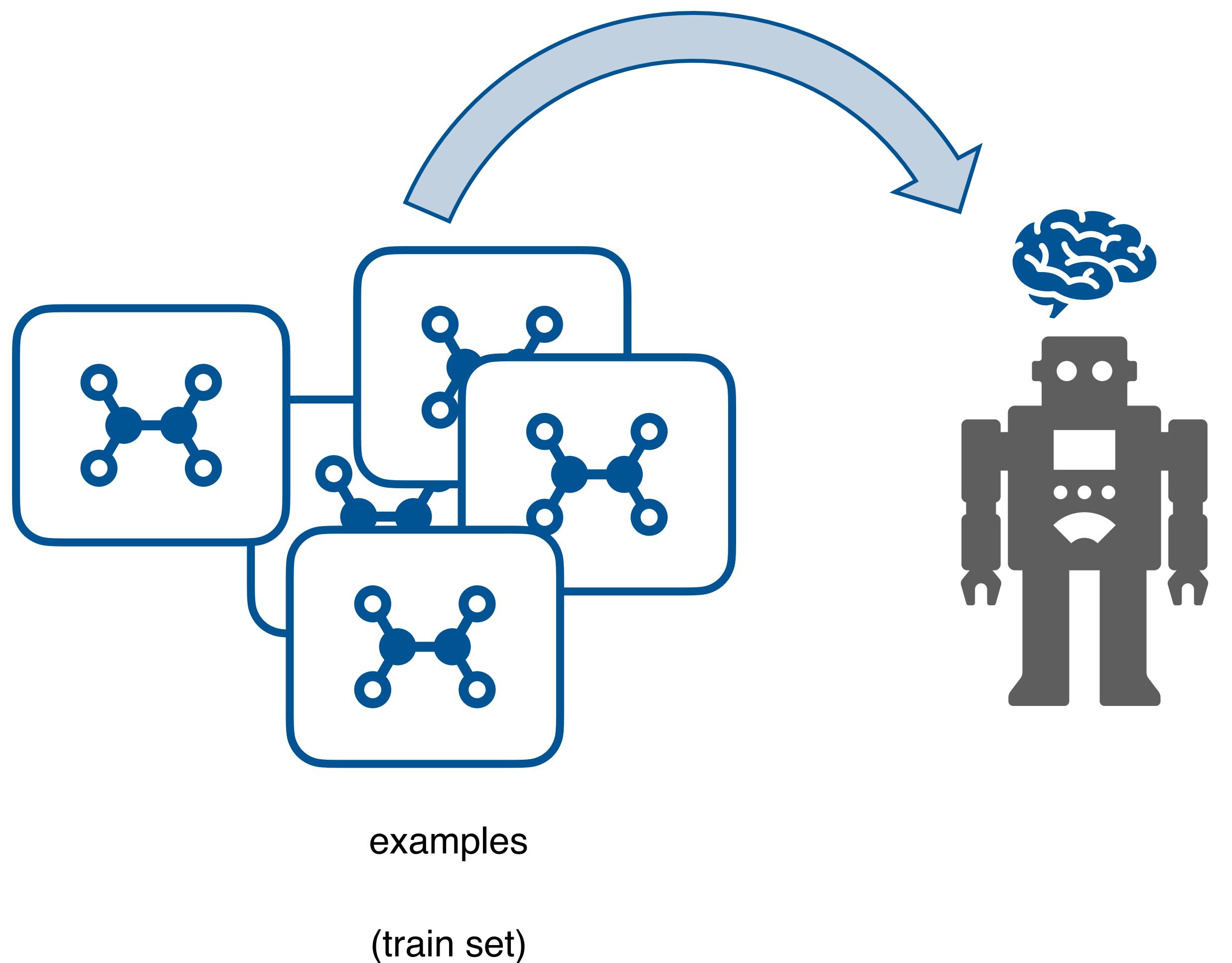


This is how we **should** represent ML:

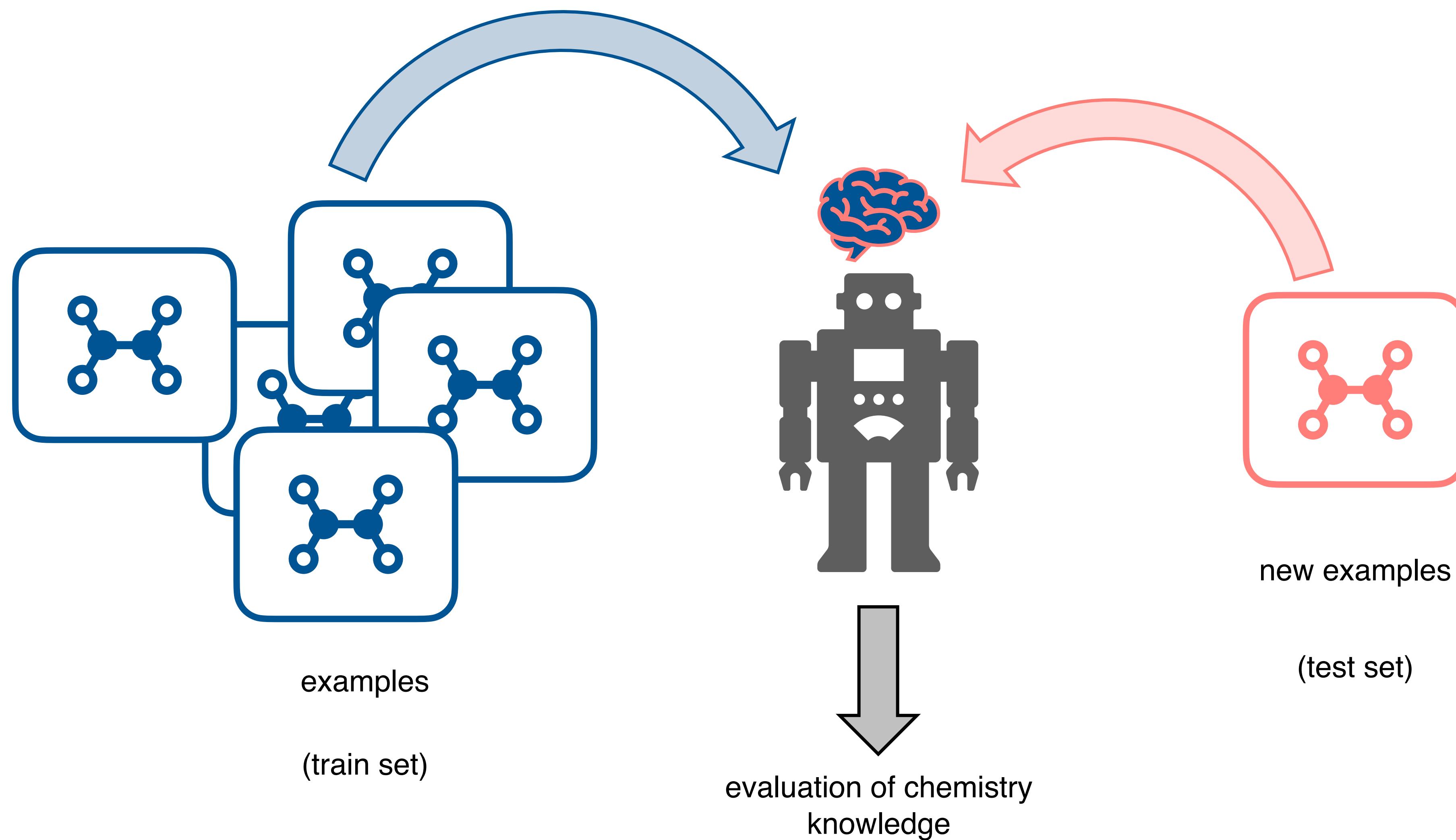


Very good at *one* task (e.g., good at understanding chemistry / molecules)

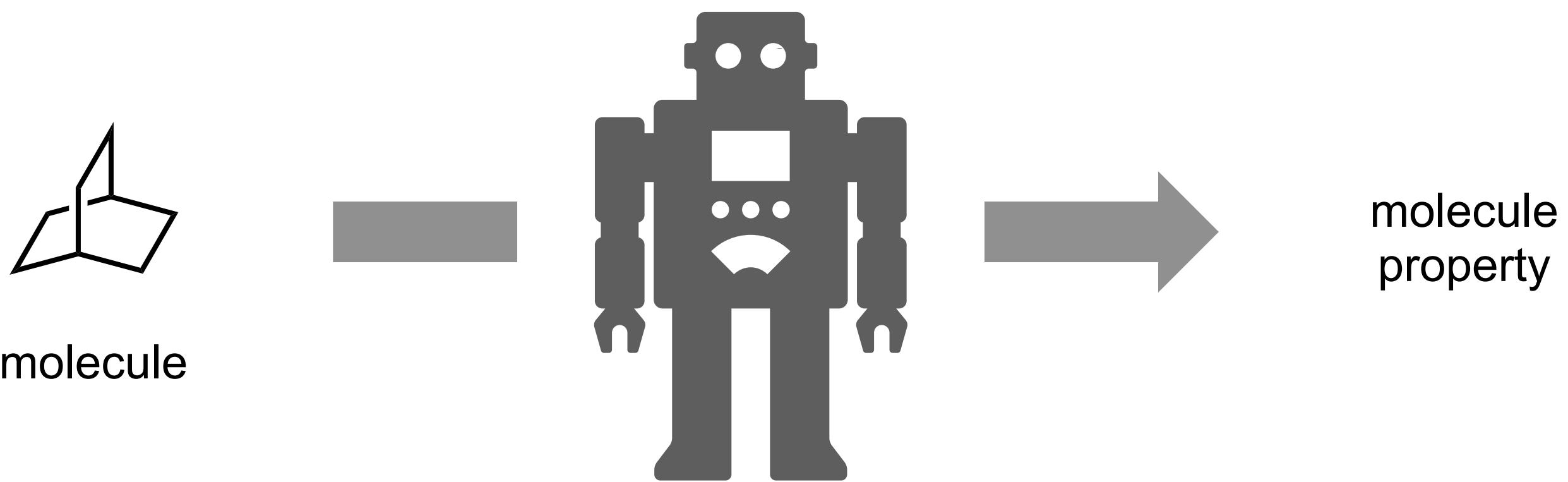
How do we do that?



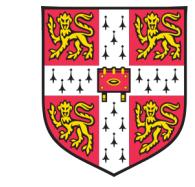
The Basic Idea Behind Machine Learning



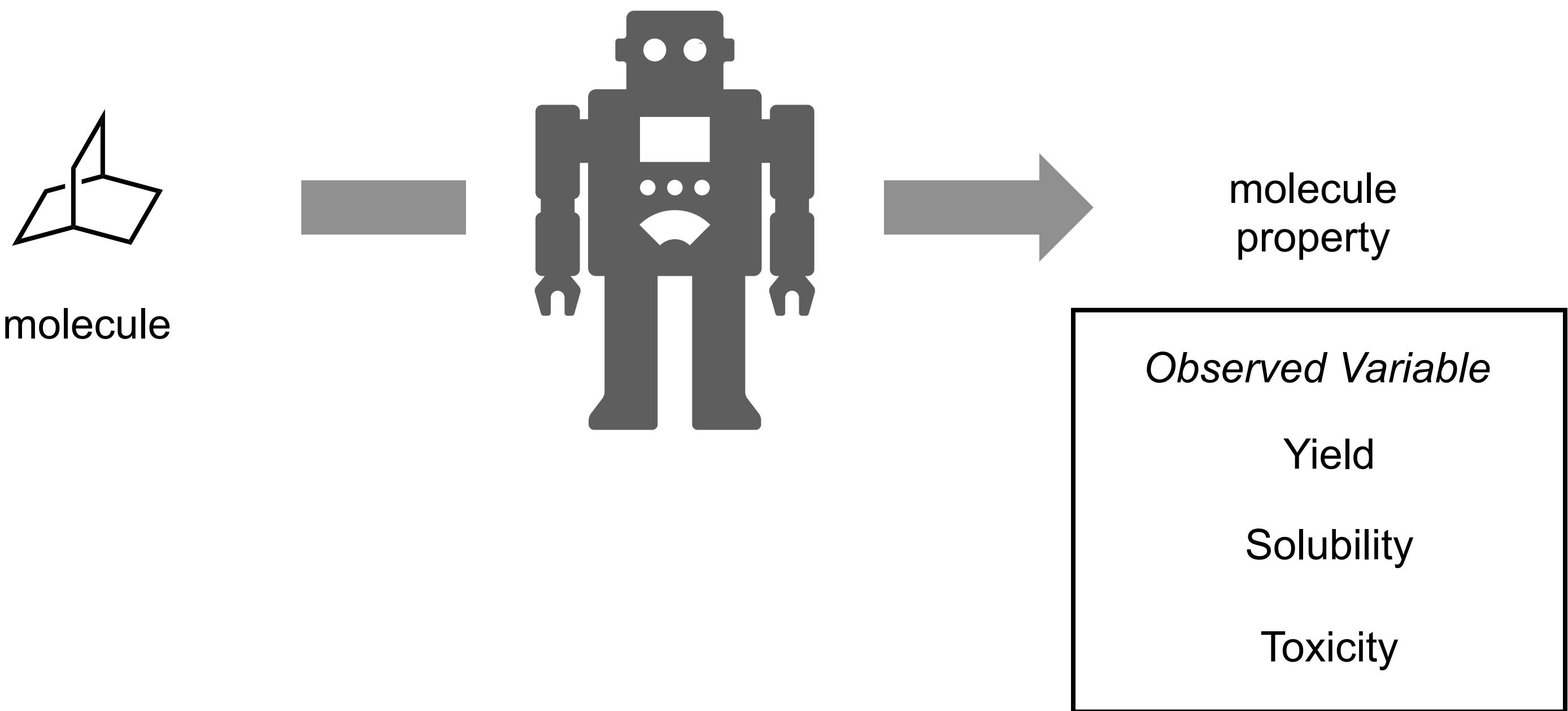
What is Machine Learning, Really?



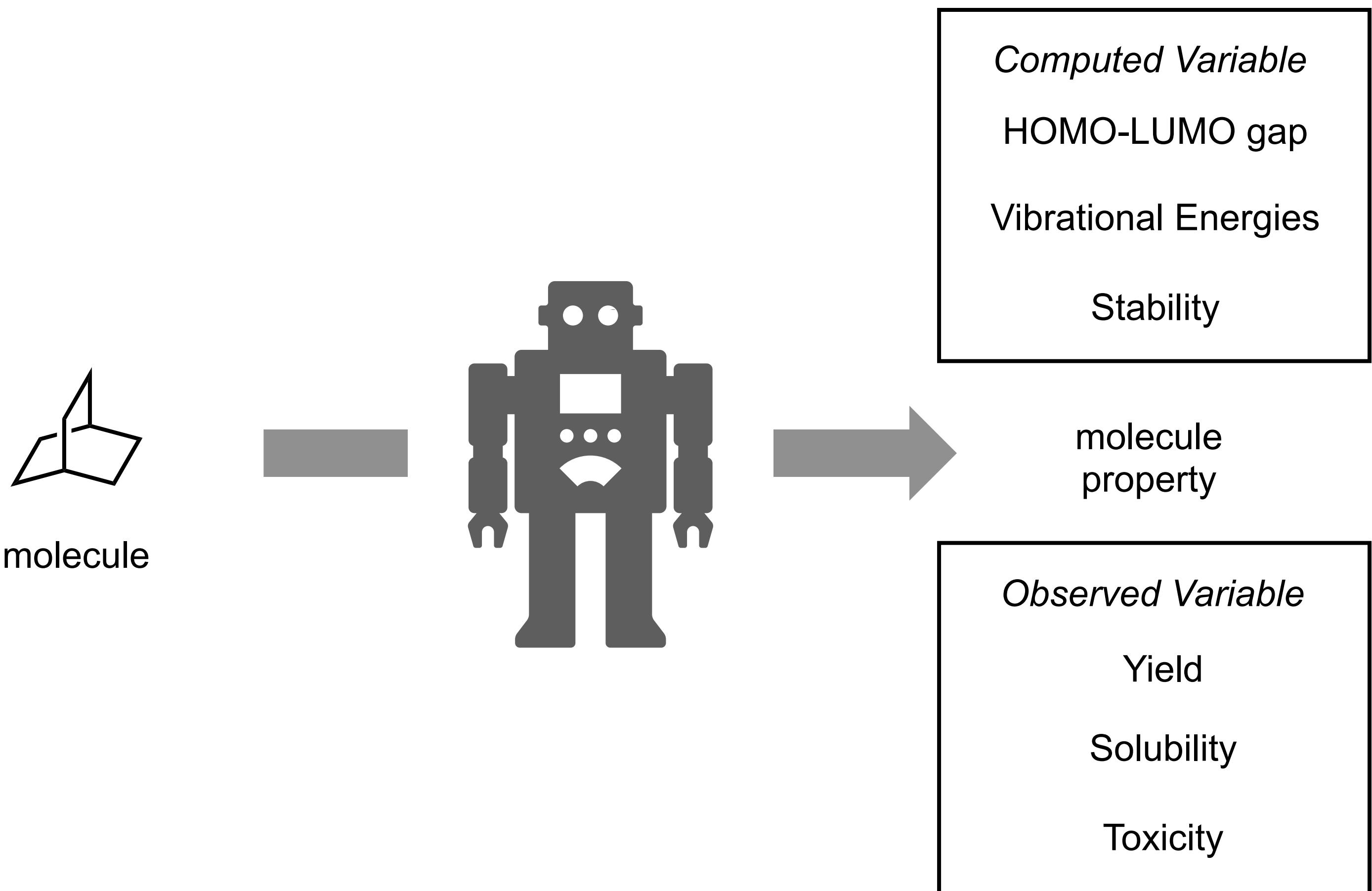
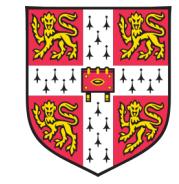
What is Machine Learning, Really?



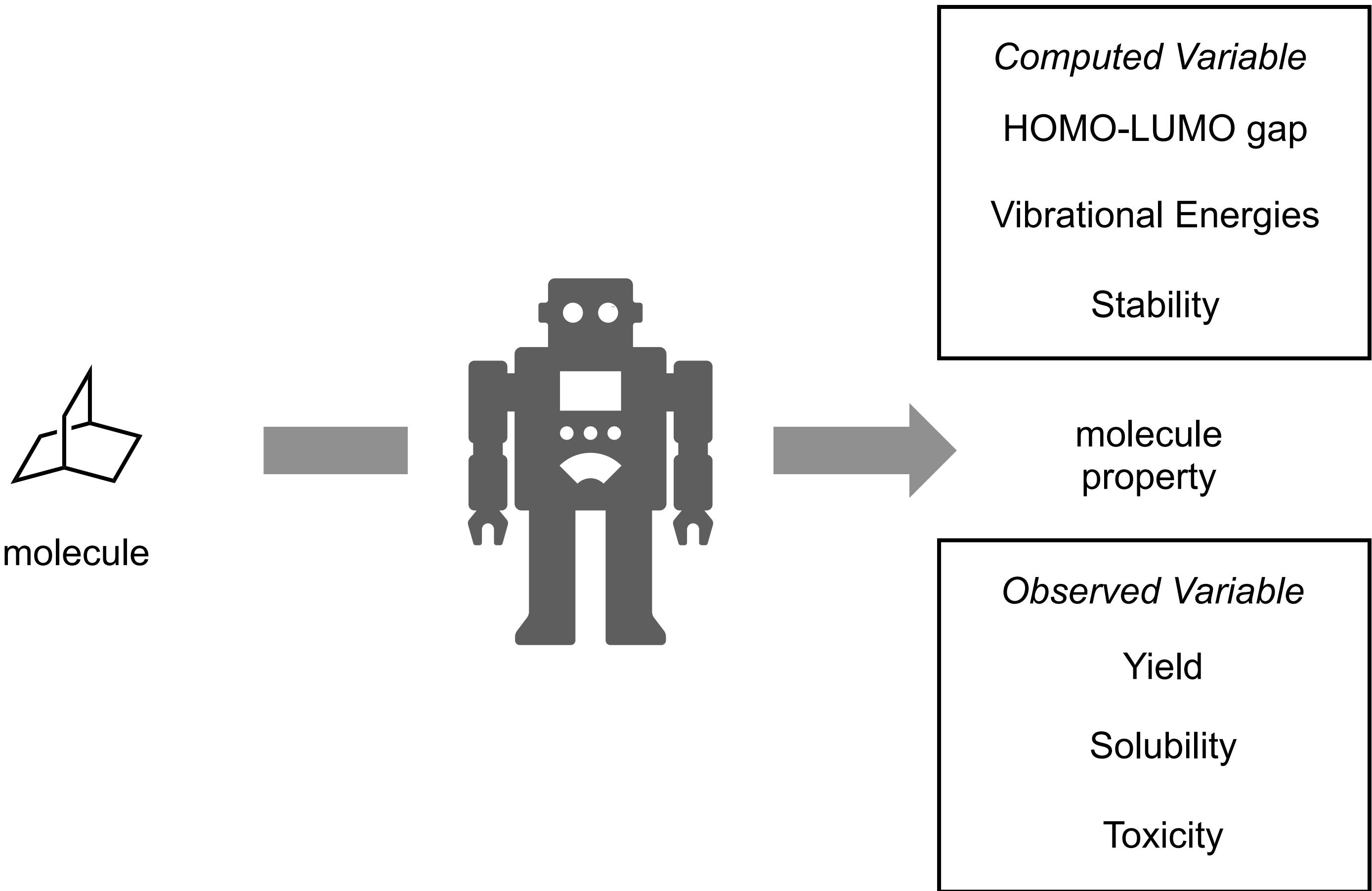
UNIVERSITY OF
CAMBRIDGE



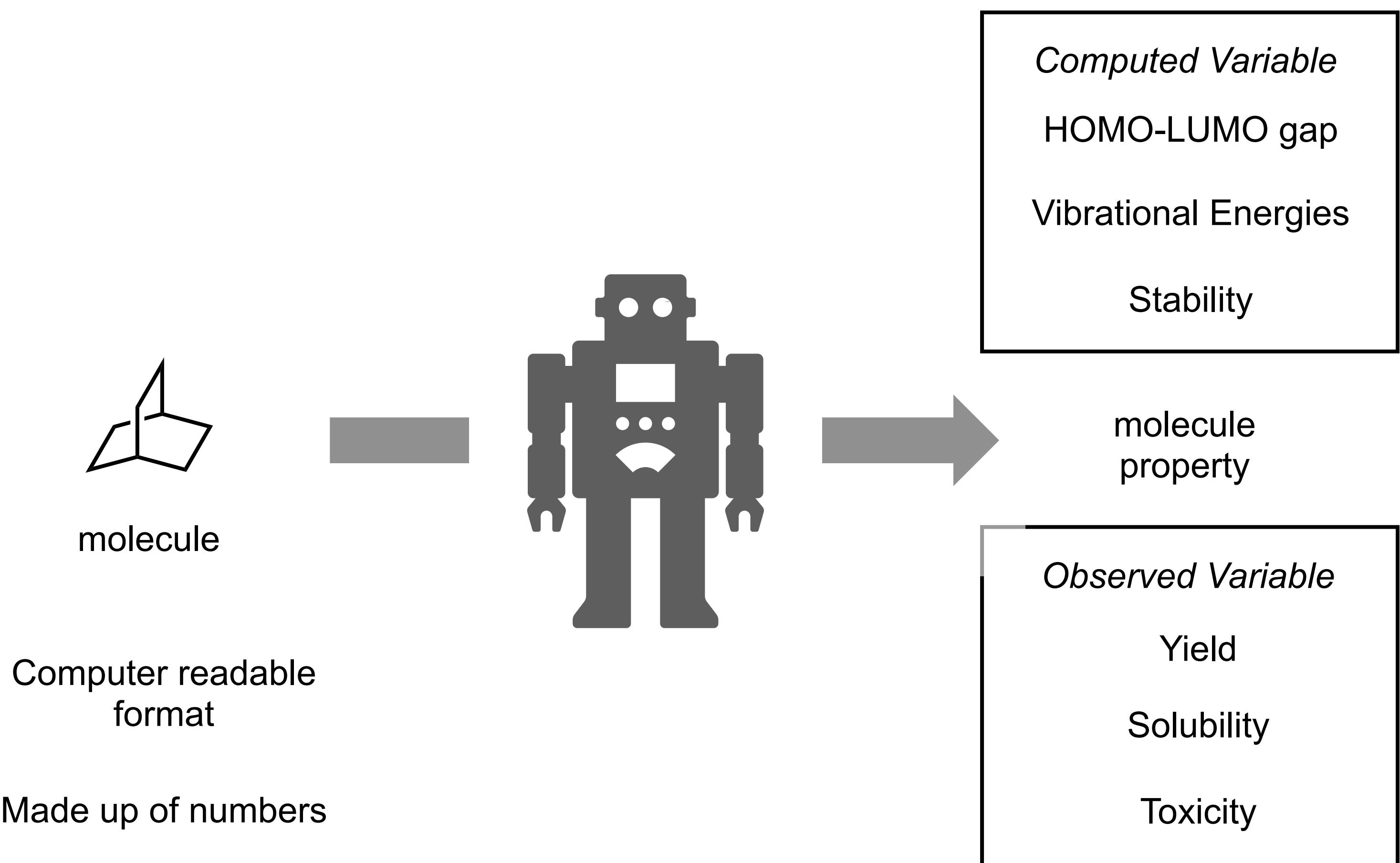
What is Machine Learning, Really?



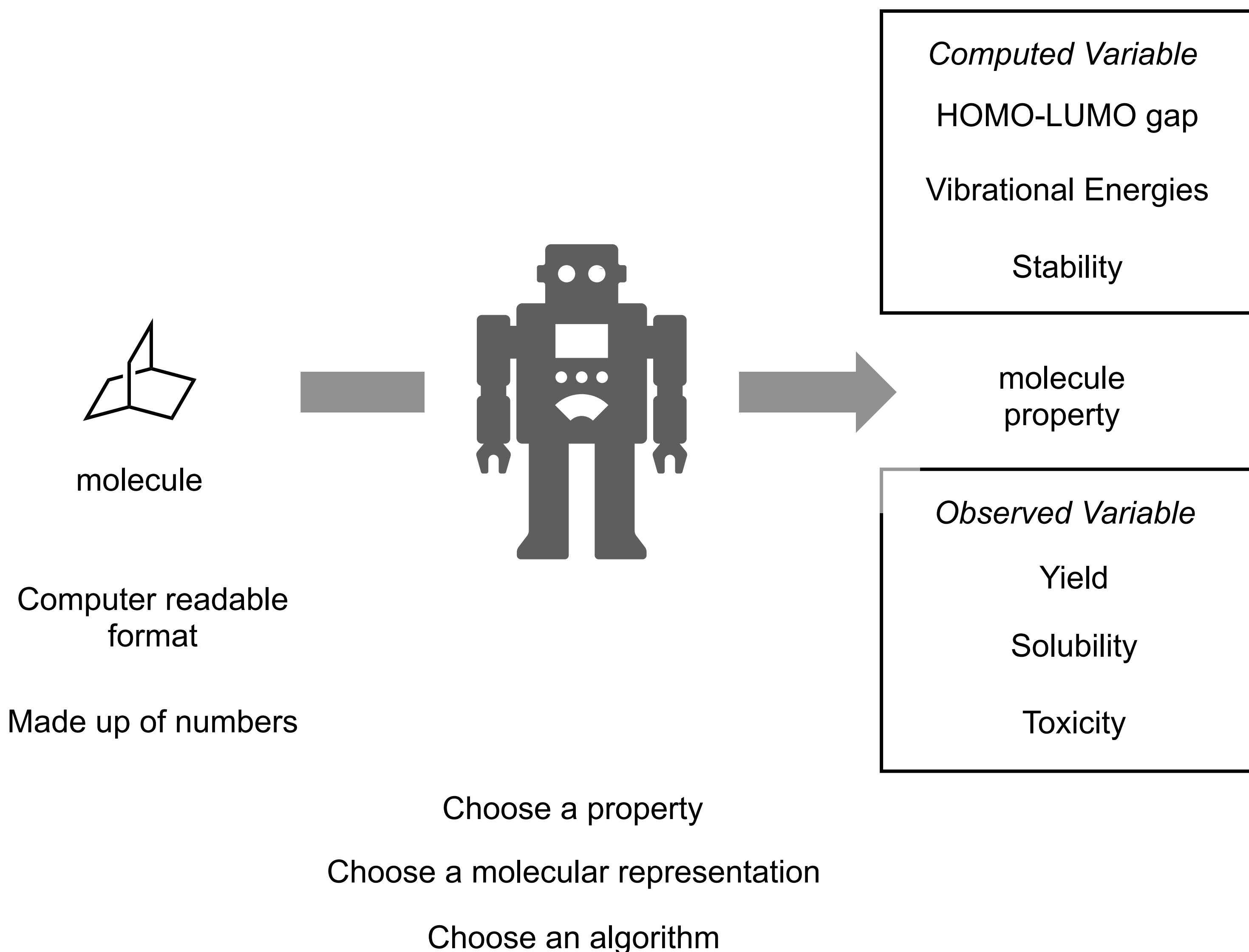
What is Machine Learning, Really?



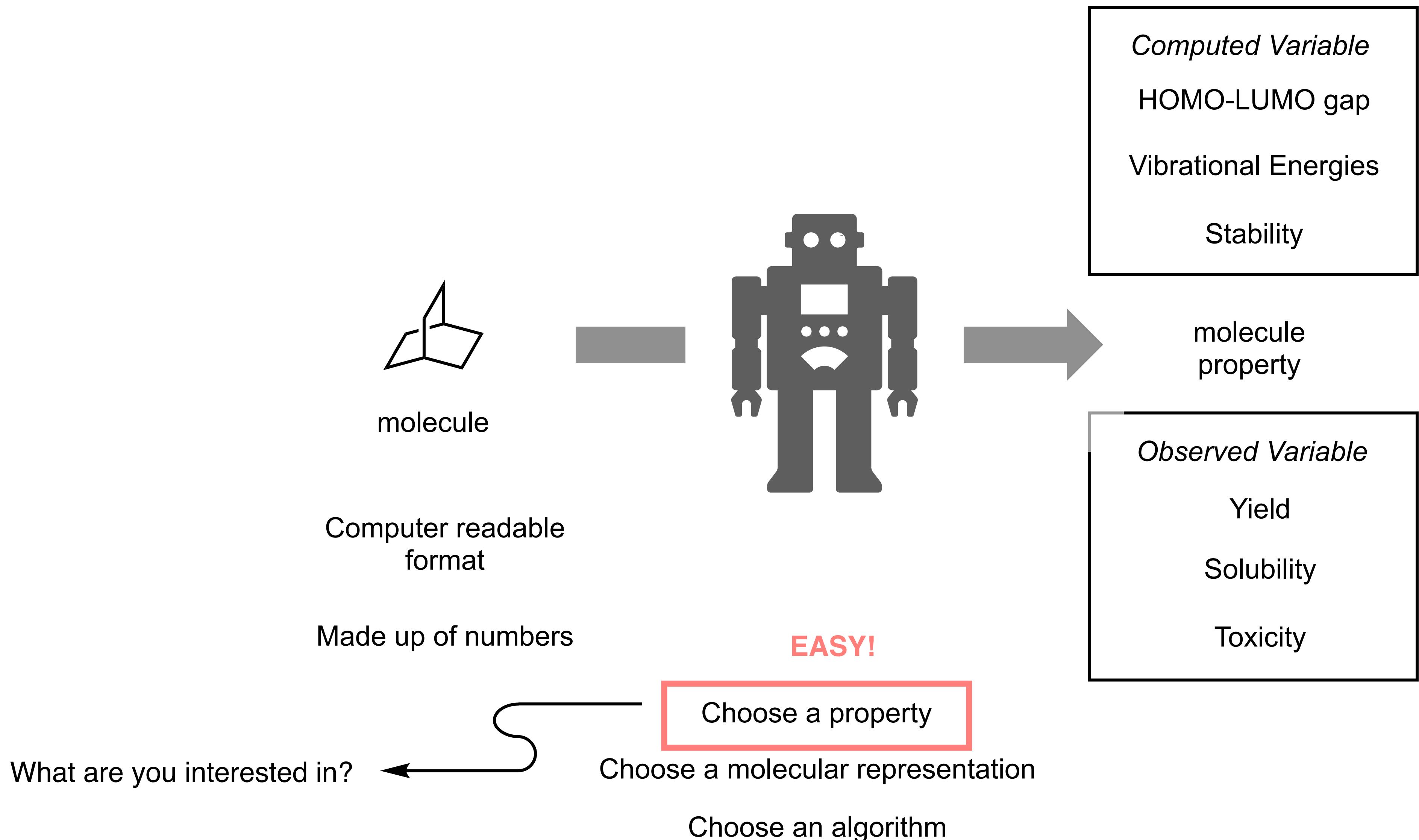
What is Machine Learning, Really?



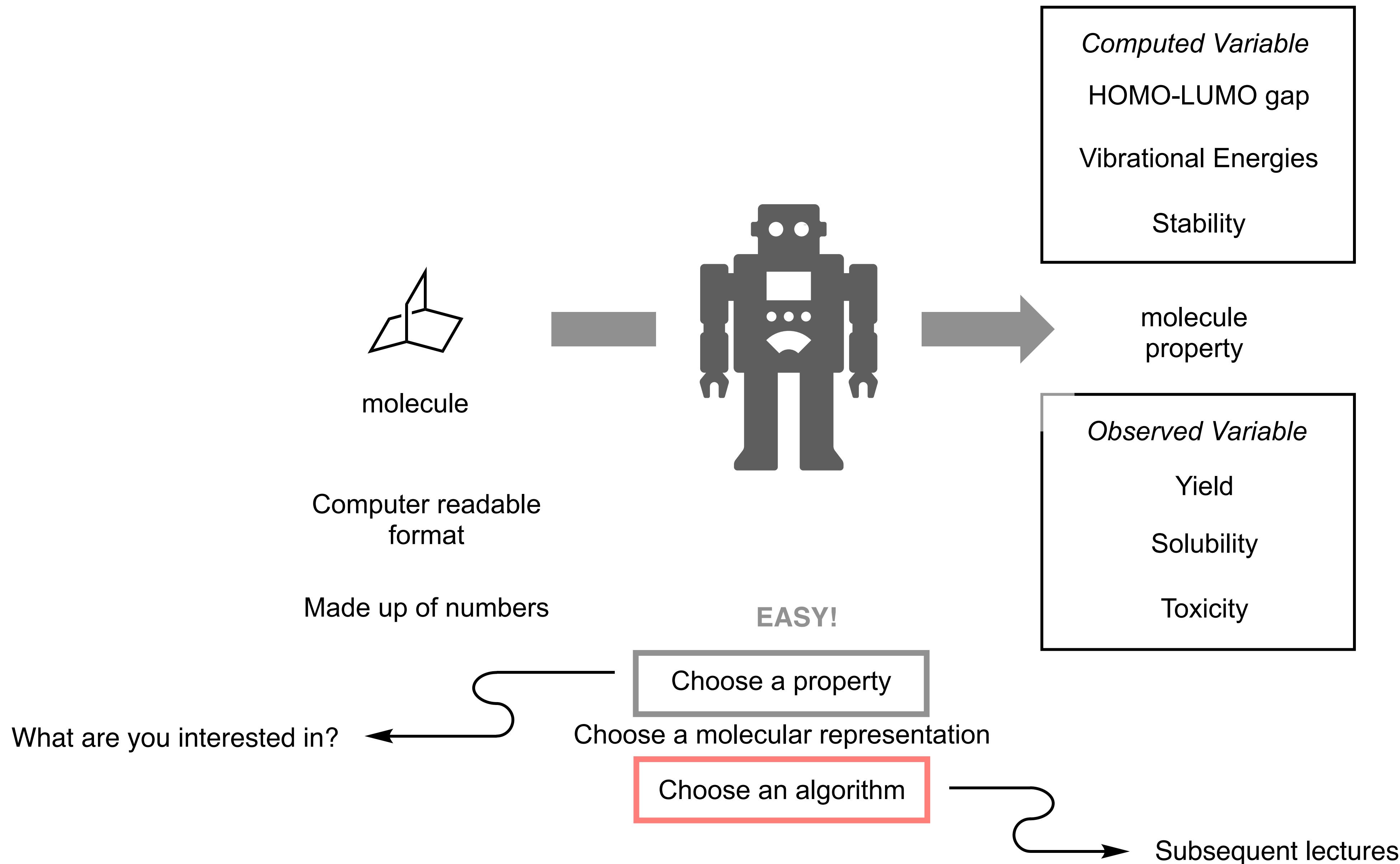
What is Machine Learning, Really?



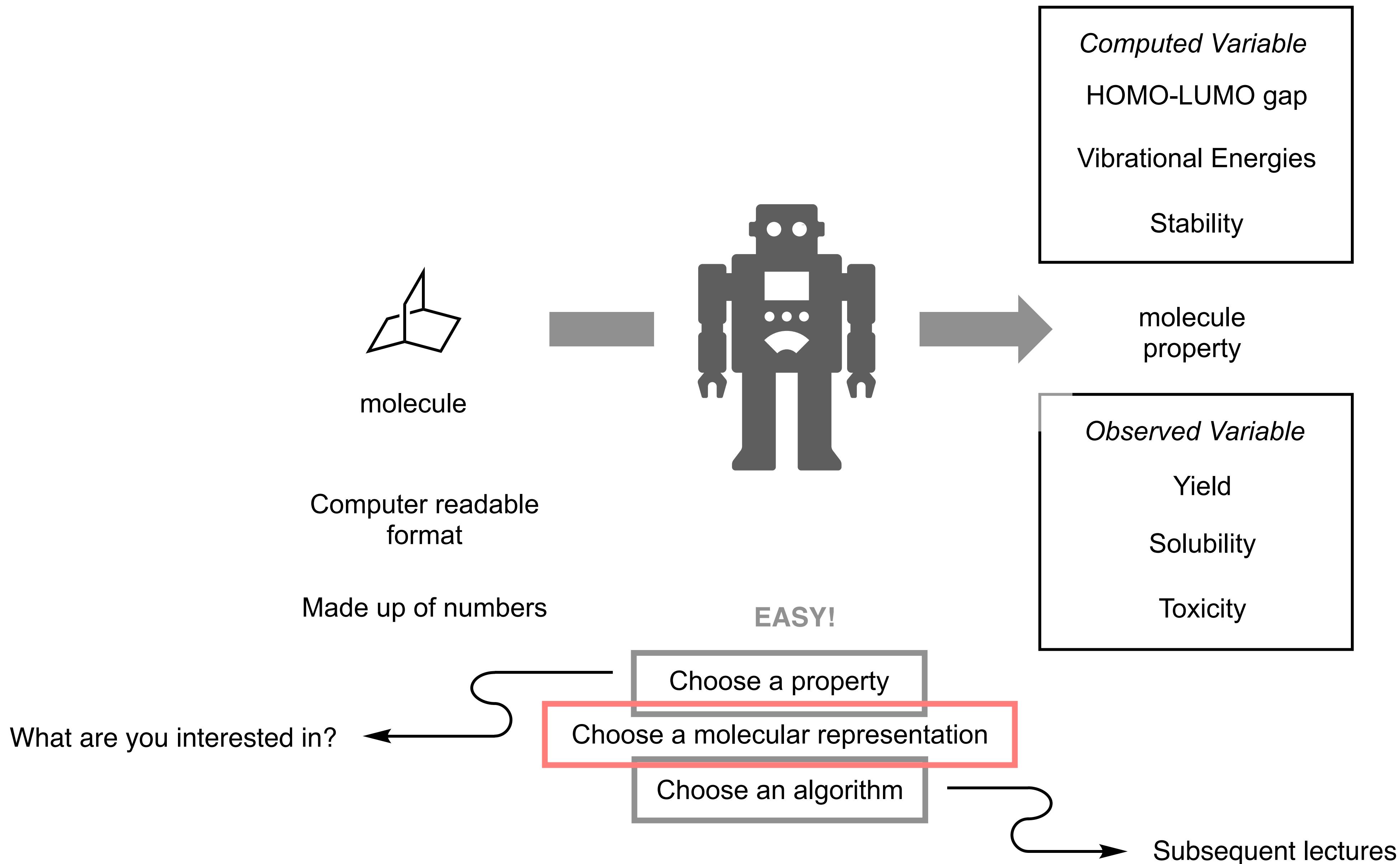
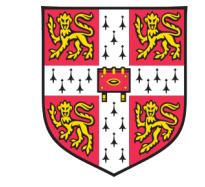
What is Machine Learning, Really?



What is Machine Learning, Really?



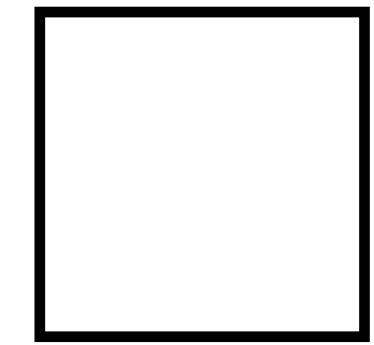
What is Machine Learning, Really?



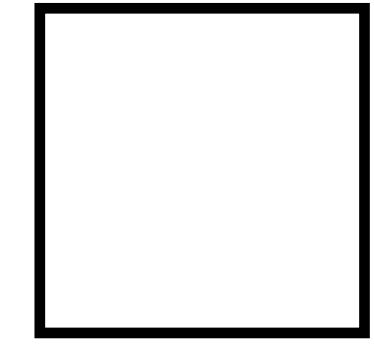
Something that stores information

Something that stores information

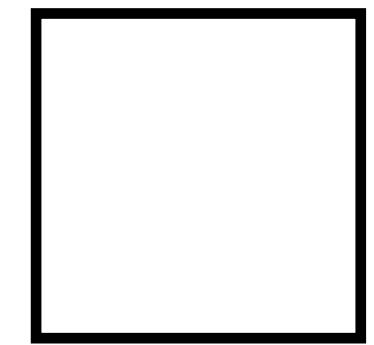
A



B



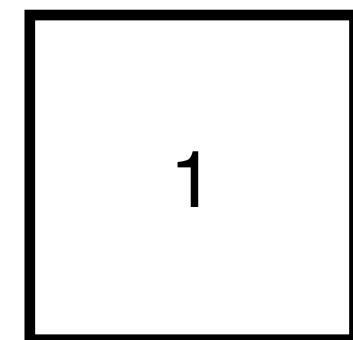
C



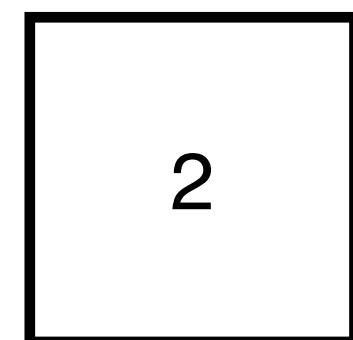
Something that stores information

Number
of Hats

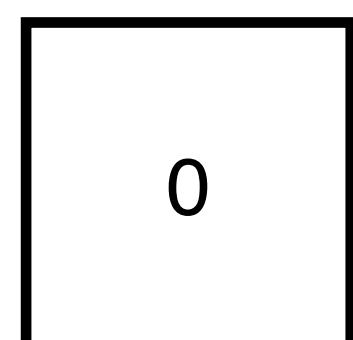
A



B



C



1D vectors

Something that stores information

Number Number
of Hats of Cats

A

1	2
---	---

B

2	5
---	---

C

0	3
---	---

2D vectors

Something that stores information

Number of Hats	Number of Cats	Favorite Number
----------------	----------------	-----------------

A

1	2	e
---	---	---

B

2	5	π
---	---	-------

C

0	3	$-\sqrt{2}$
---	---	-------------

3D vectors



Something that stores information

	Number of Hats	Number of Cats	Favorite Number	Birth-date	Birth Month	Birth Year	Owns A Bike?	Criminal Record?	Lives in England
--	----------------	----------------	-----------------	------------	-------------	------------	--------------	------------------	------------------

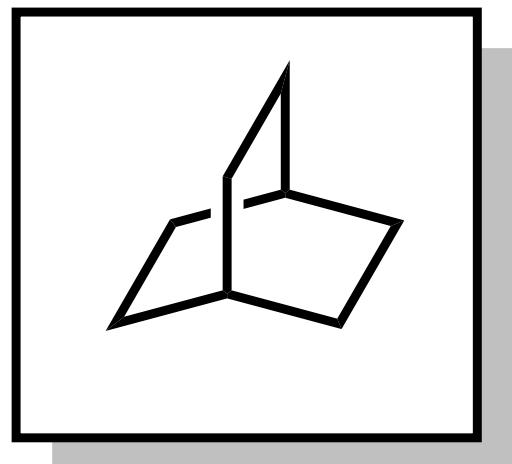
A	1	2	e	15	1	1967	0	1	0	...
---	---	---	---	----	---	------	---	---	---	-----

B	2	5	π	30	9	2001	0	0	0	...
---	---	---	-------	----	---	------	---	---	---	-----

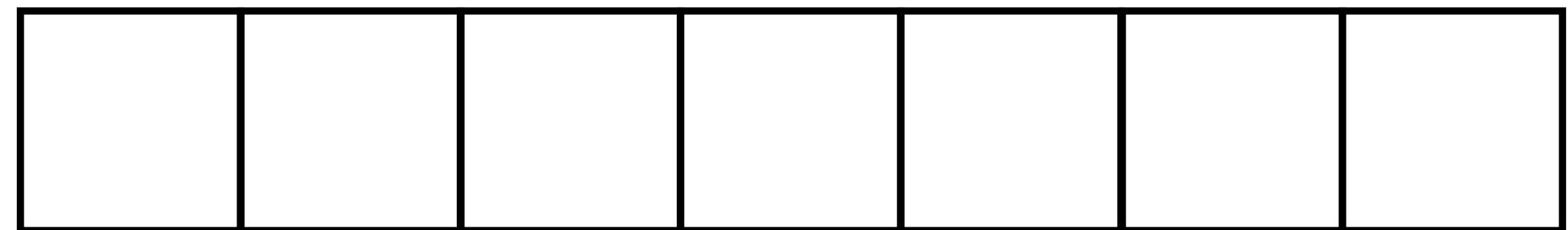
C	0	3	$-\sqrt{2}$	22	5	1991	1	0	1	...
---	---	---	-------------	----	---	------	---	---	---	-----

0 = No 1 = Yes

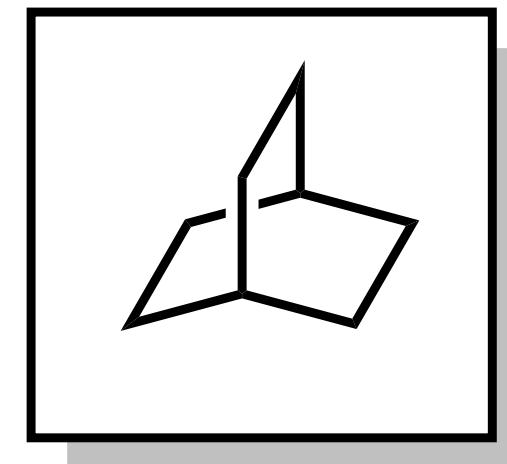
Nth dimensional vector



Presence / Absence “important” features
Molecular Fingerprints



NOTE: This is an oversimplified version of fingerprints!

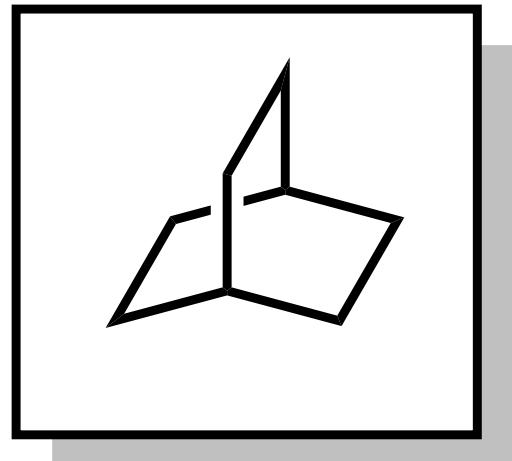


Presence / Absence “important” features

Molecular Fingerprints

OH	CH ₂	CH ₃	bridge-head C	CHO	Ph	CH ₂ -CH ₂
0	1	0	1	0	0	1

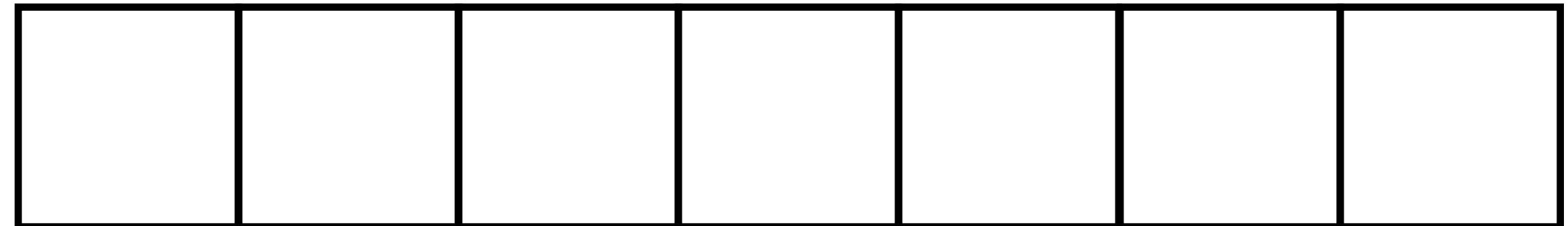
NOTE: This is an oversimplified version of fingerprints!



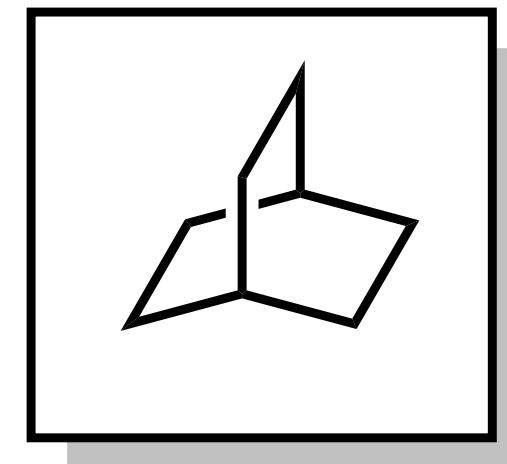
Chemical Properties

Typically User Defined

Often incorporates quantum chemical properties



NOTE: This is an oversimplified version of featurization!



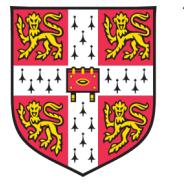
Chemical Properties

Typically User Defined

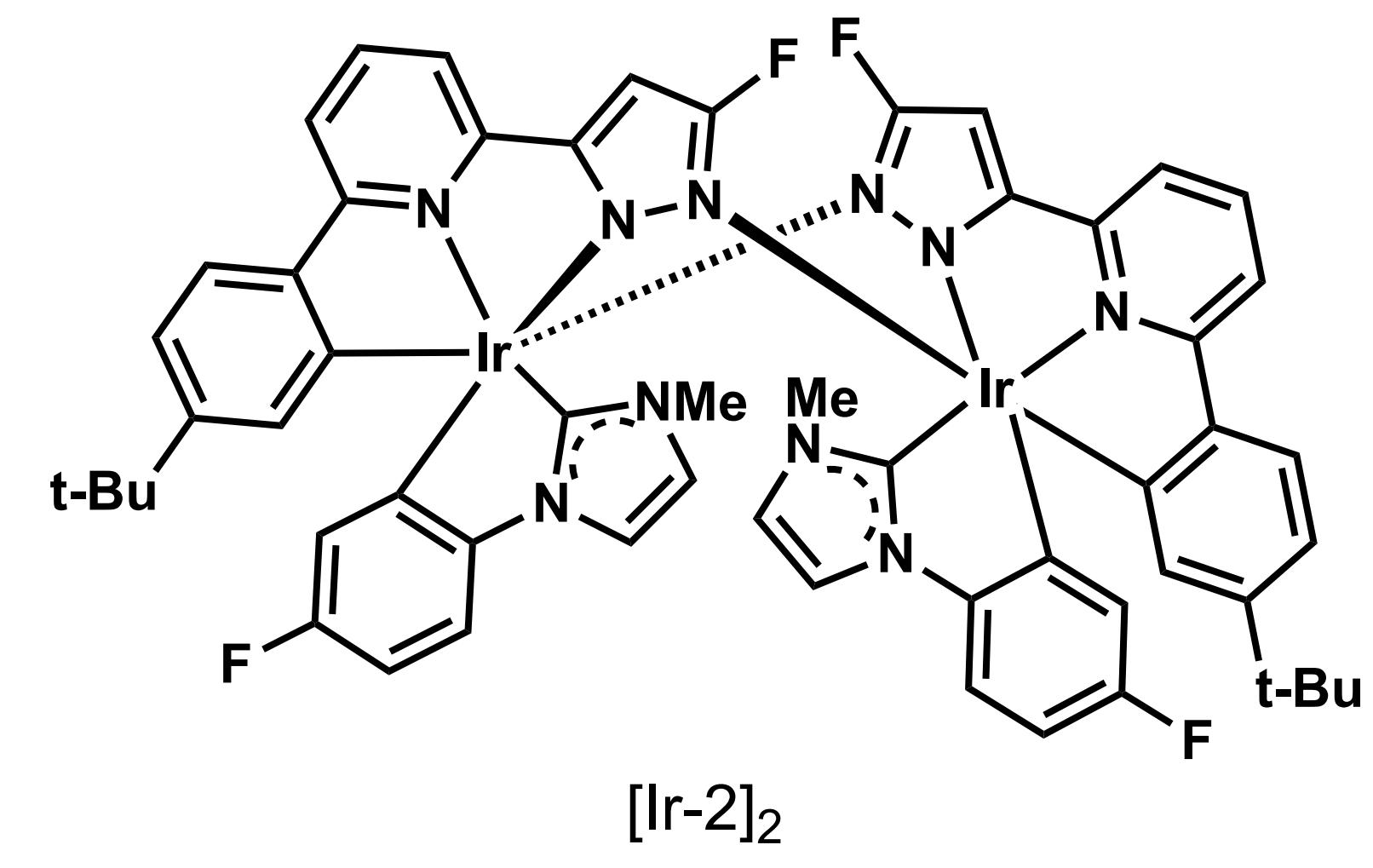
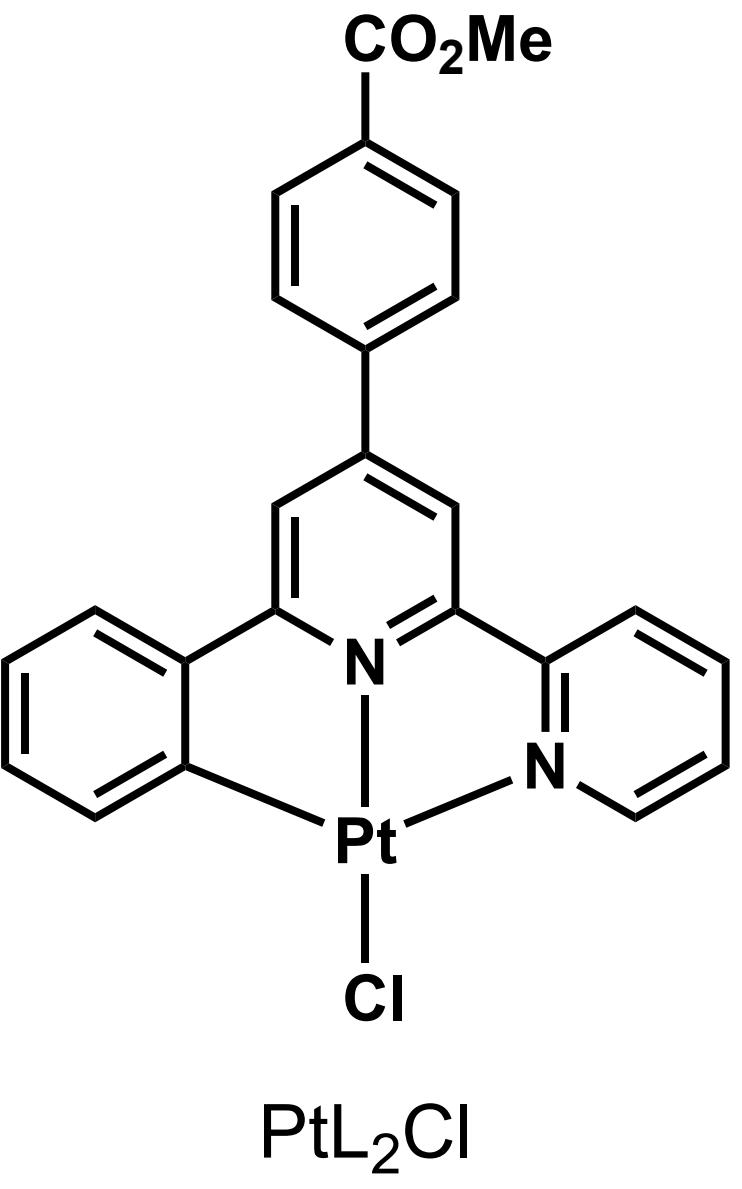
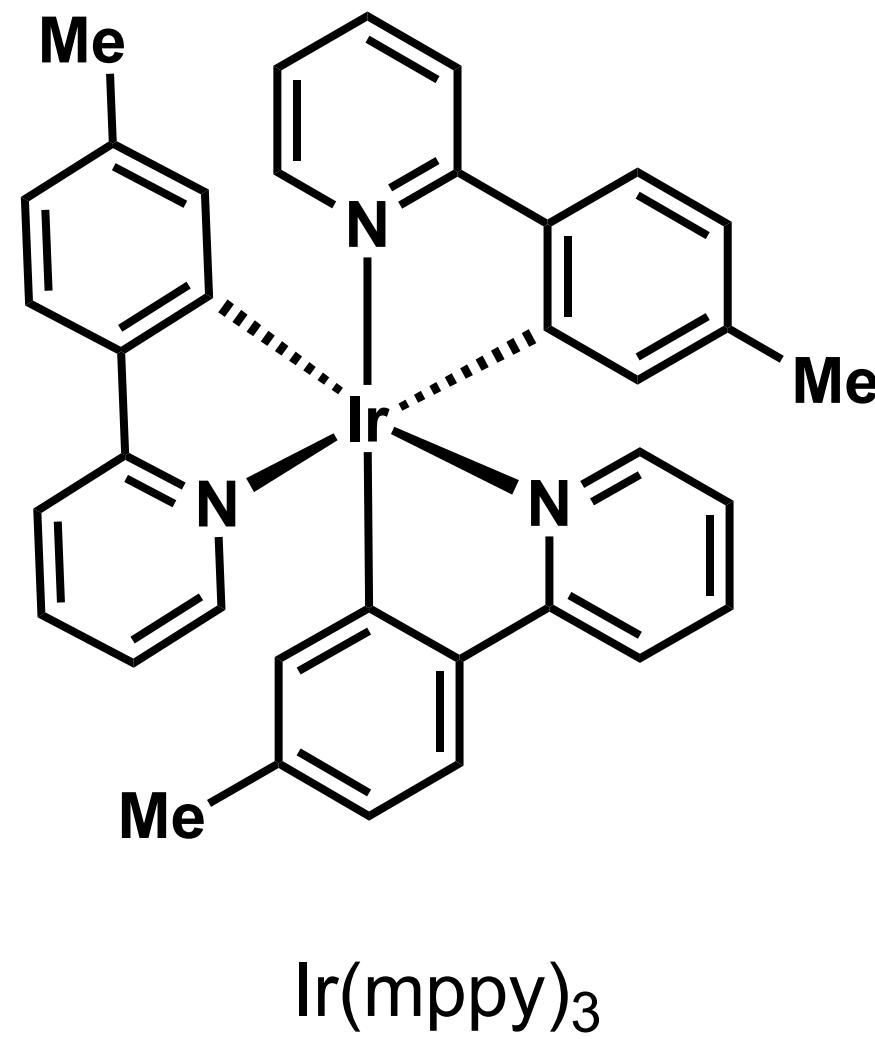
Often incorporates quantum chemical properties

logP	HOMO	LUMO	gap	Num Atoms	Num Hs	Dihedral Angle 1
2.2	-10	-5	5	22	14	120

NOTE: This is an oversimplified version of featurization!



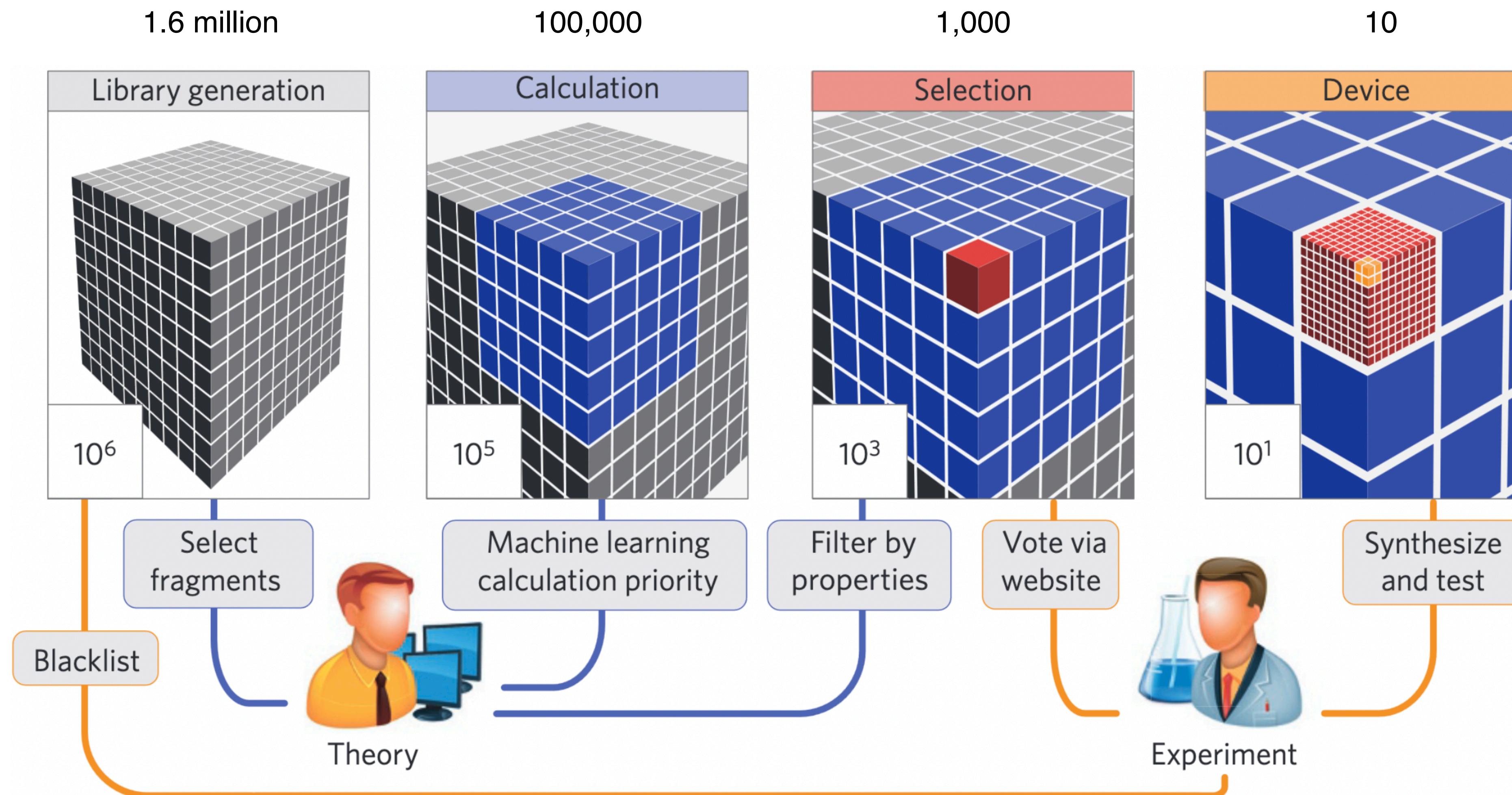
Recent OLED Molecules:



Expensive heavy metals

*Can we identify new non-heavy metal containing
OLED molecules?*

ML-assisted library prioritization of new OLED materials

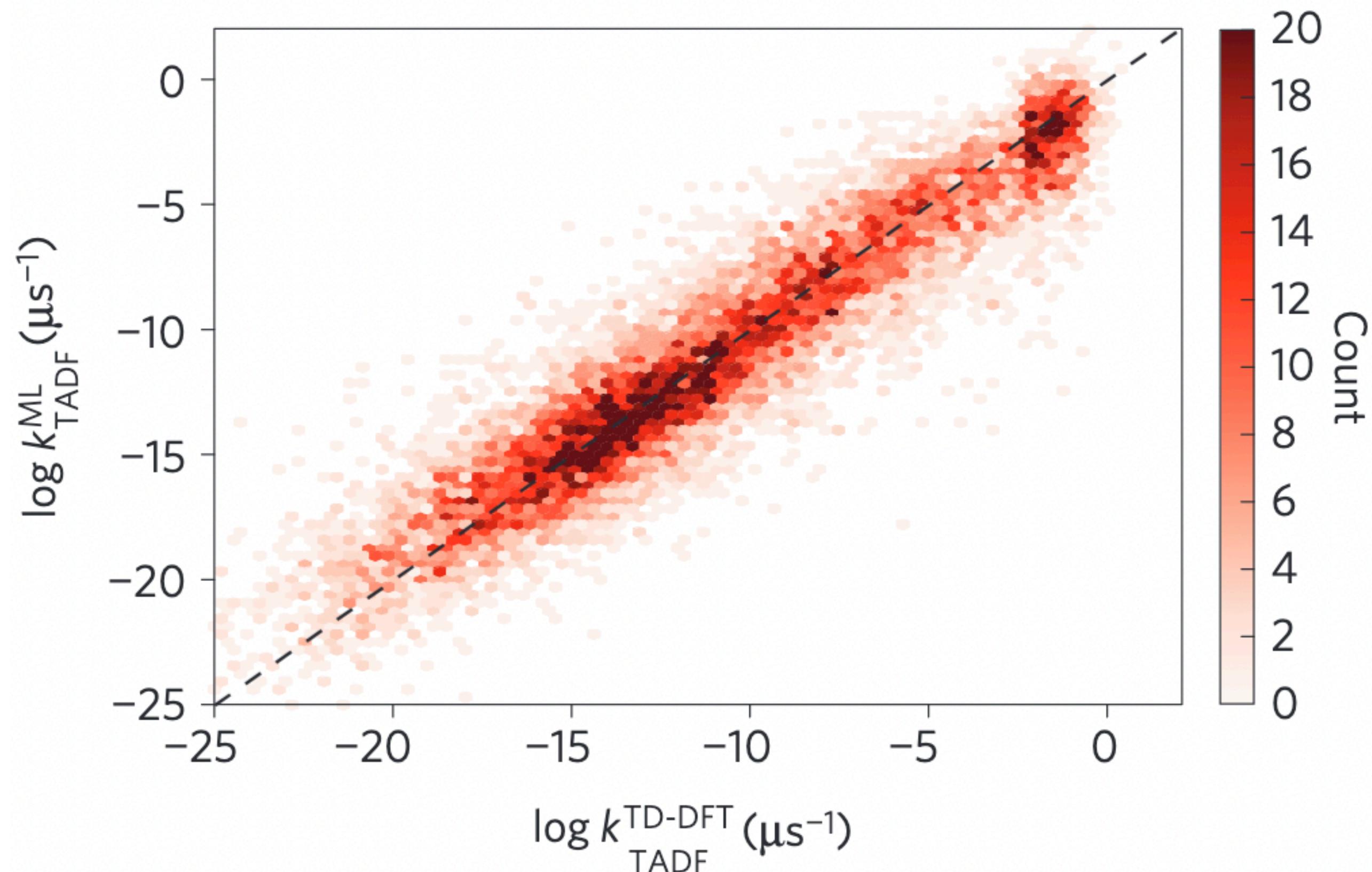


1.6 million compounds is too many to experimentally validate **AND** to perform quantum simulations upon.

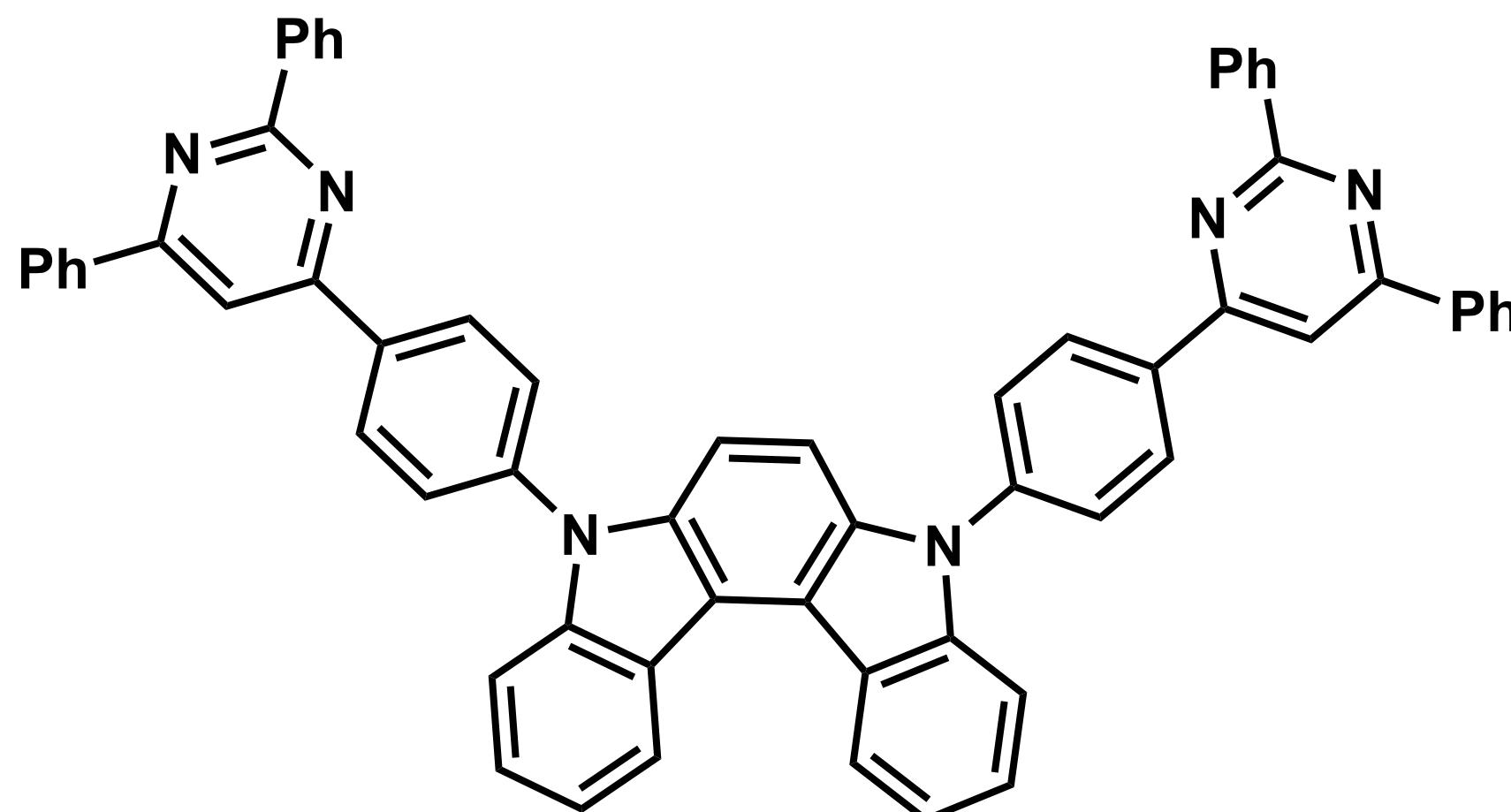
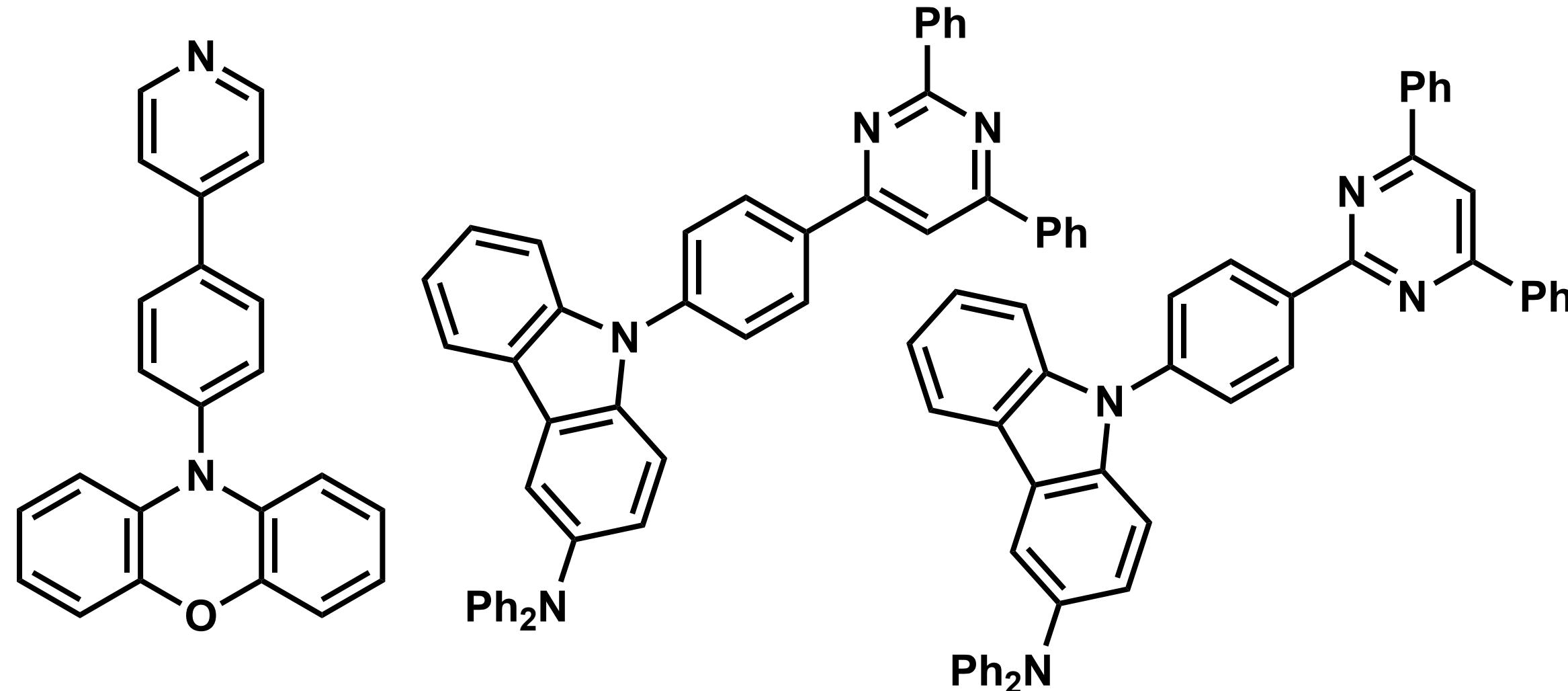


Randomly chose 40,000 library compounds for DFT computation of the rate of thermally activated delayed fluorescence (k_{TDAF}).

Developed a ML model that could accurate predict k_{TDAF} at fraction of time but with same level of accuracy.

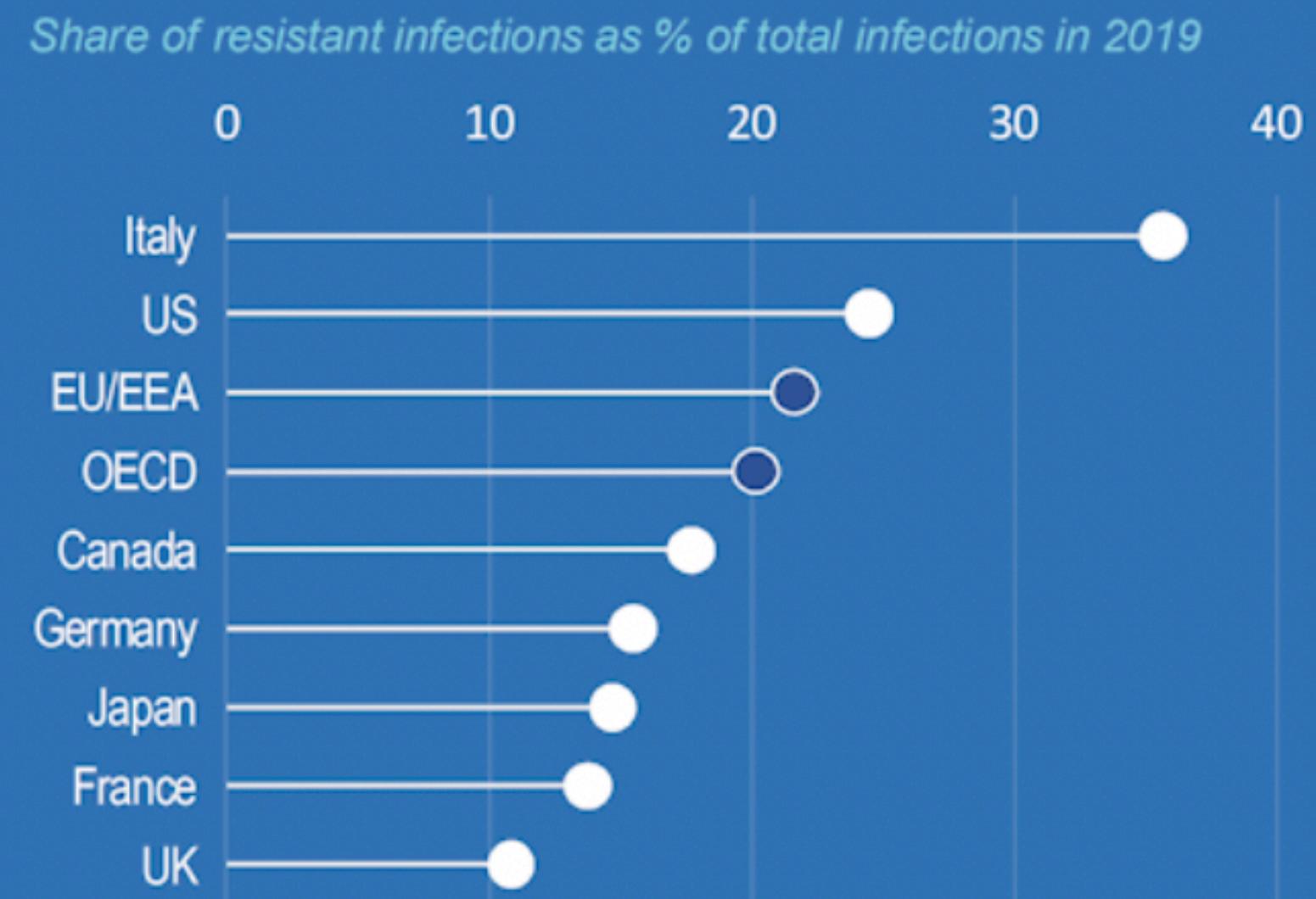


Experimentally validated promising compounds.



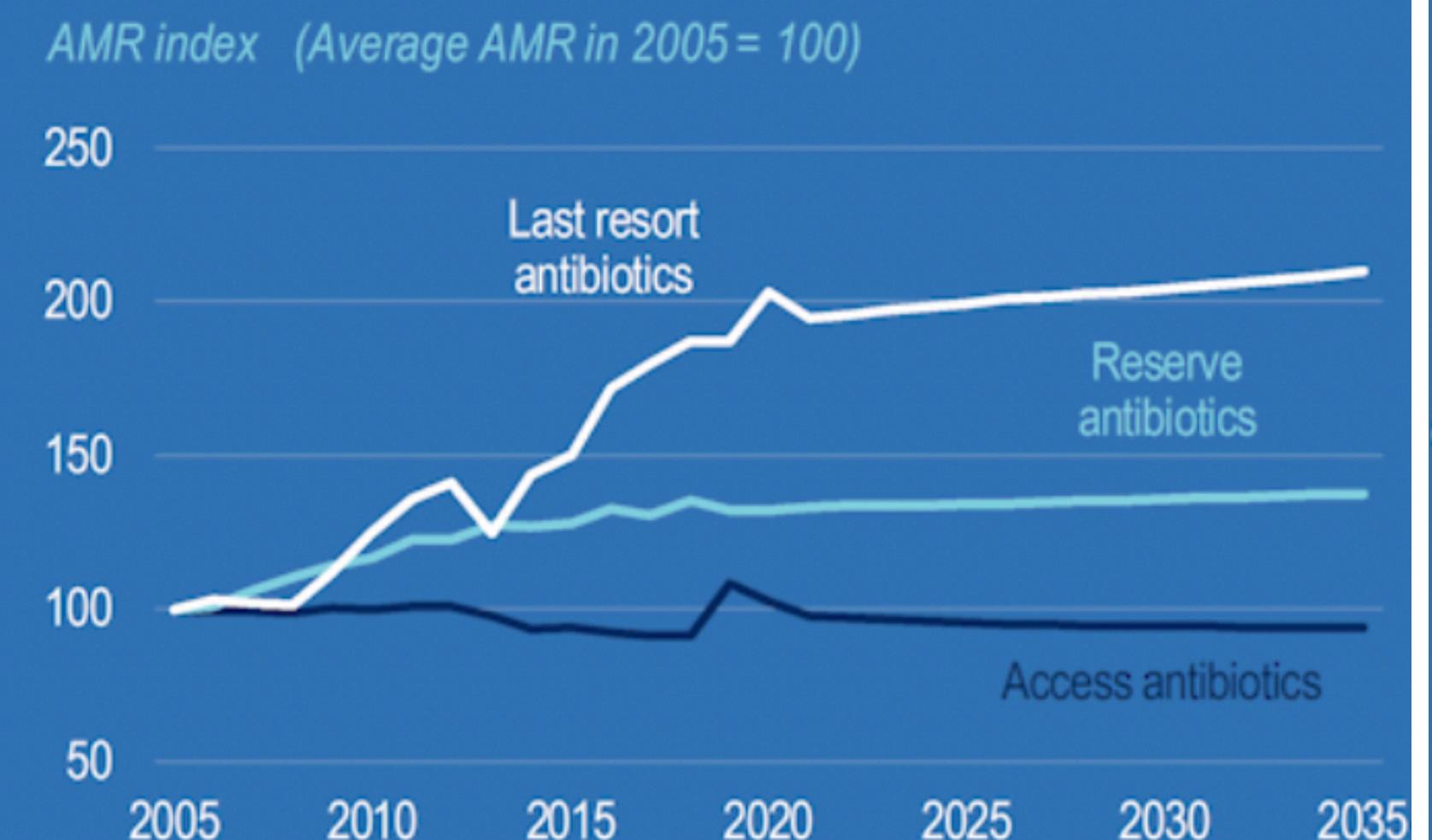
The AMR pandemic is here

One in five infections in OECD are resistant to antibiotic treatment. This will not improve without policy action.



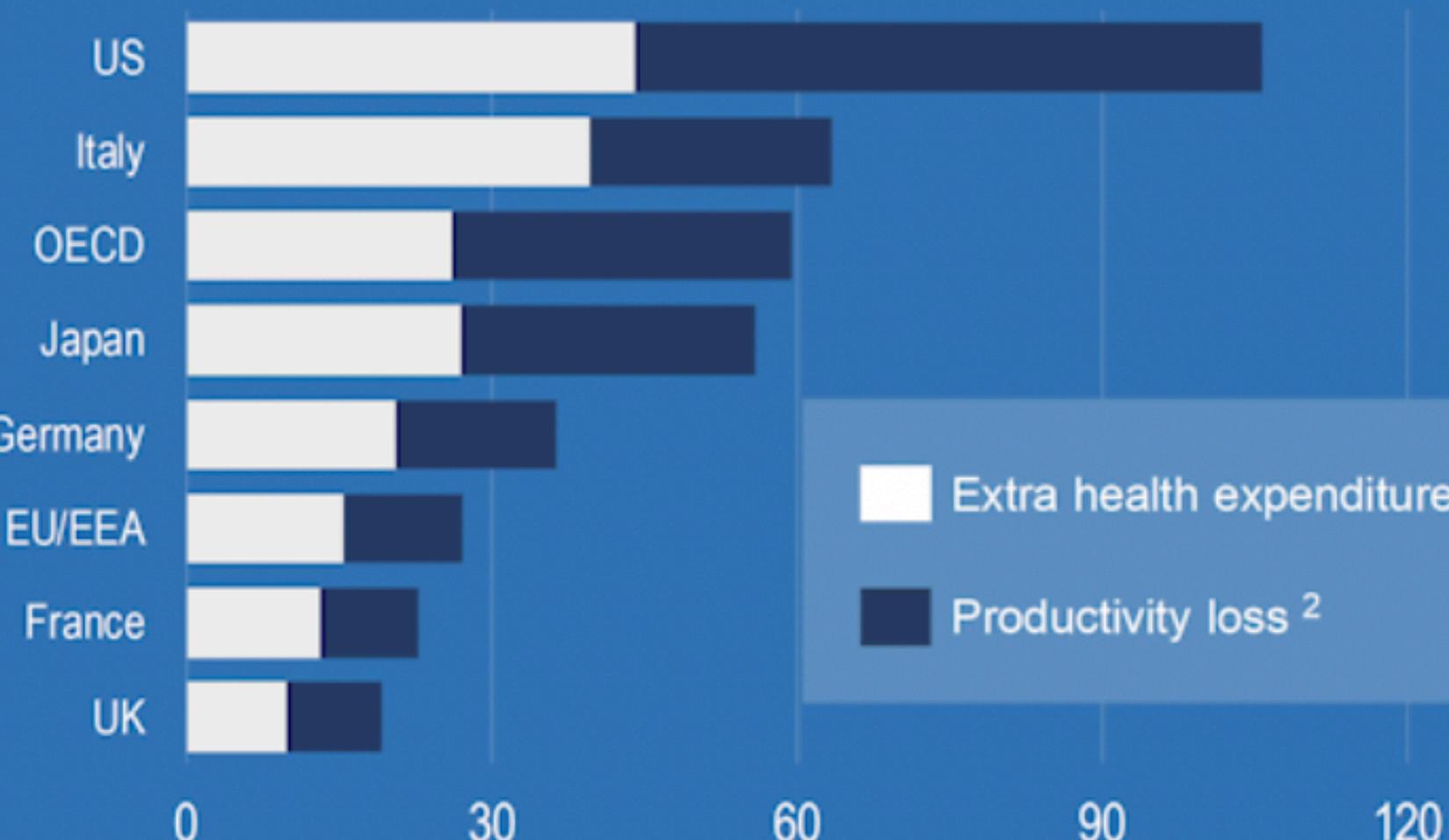
We are exhausting our antibiotic arsenal

Resistance to last resort drugs in OECD countries could more than double by 2035 compared to 2005.



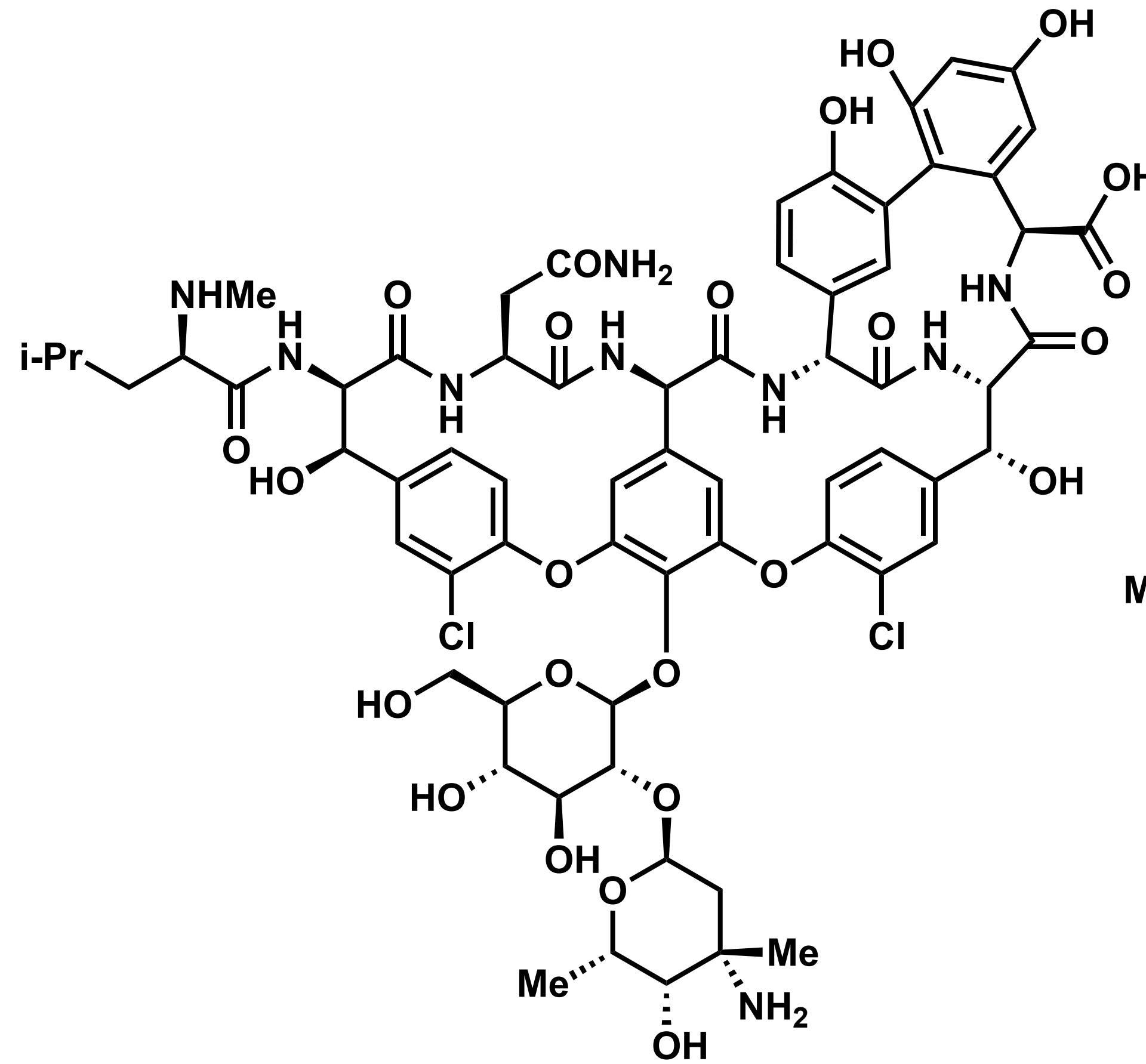
We pay a high price for inaction

Health and labour cost of resistant infections per year up to 2050, Per capita (USD PPP)

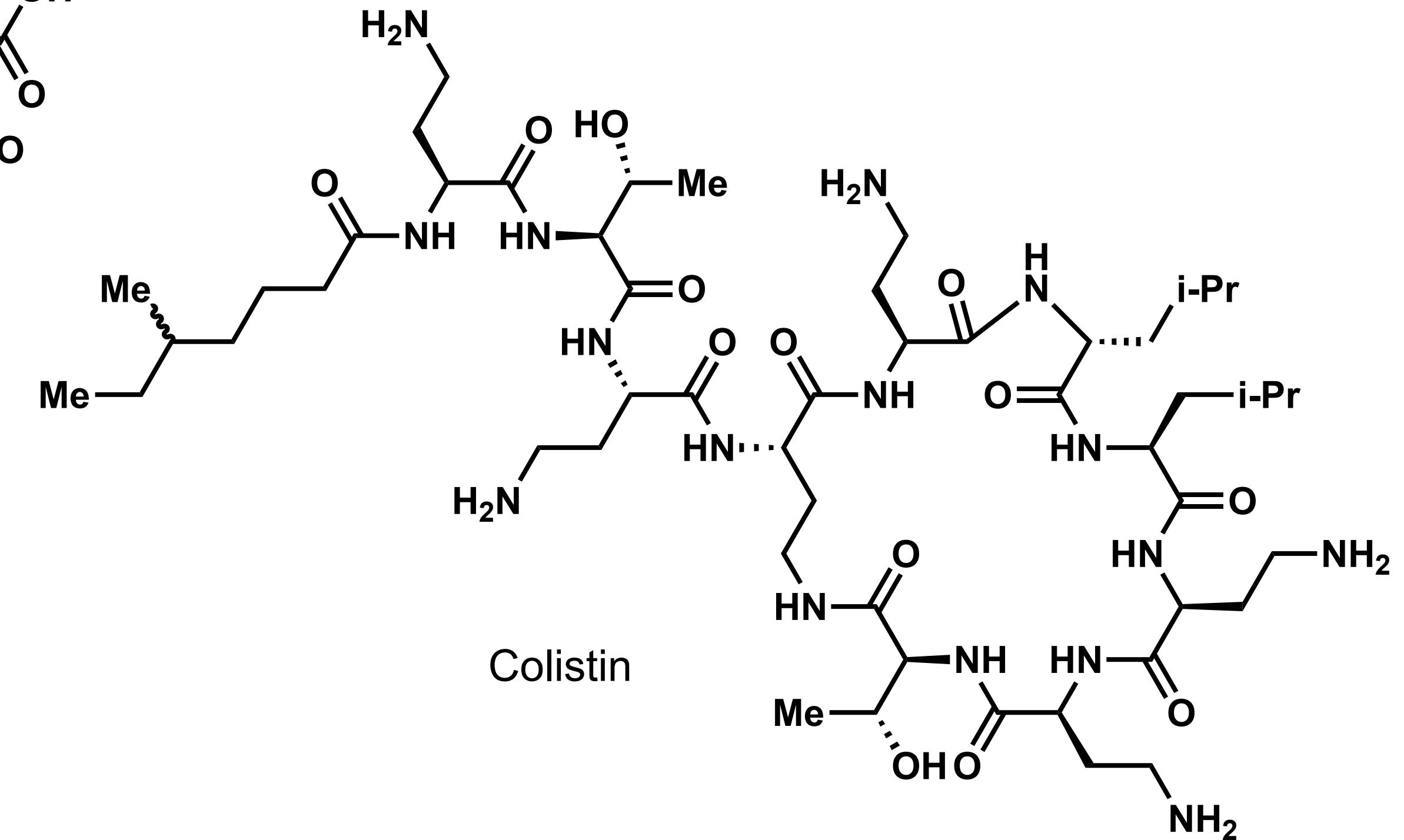


Need to discovery new antibiotics.

Antimicrobial Peptides:



Vancomycin



Colistin

Can we identify new antimicrobial peptides?



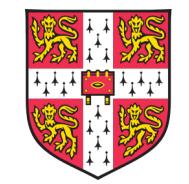
Test Set (Size)	Avg Number of Correct Predictions*	Avg ML Accuracy*
Unseen sequences from train data (~97) ^a	N/A	90.2%
Sequences annotated as antibacterial peptides in Swiss-Prot (39) ^b	26.0	66.7%
Swiss-Prot sequences that excluded peptides containing signal + antibacterial sequences (24) ^b	22.0	91.7%

* = Average over 3 different ML techniques

^a = Accuracy = correct predictions / total predictions

^b = Accuracy = correct predictions / (correct prediction + false positives)

Best average accuracy shown.



REPORT

ORGANIC CHEMISTRY

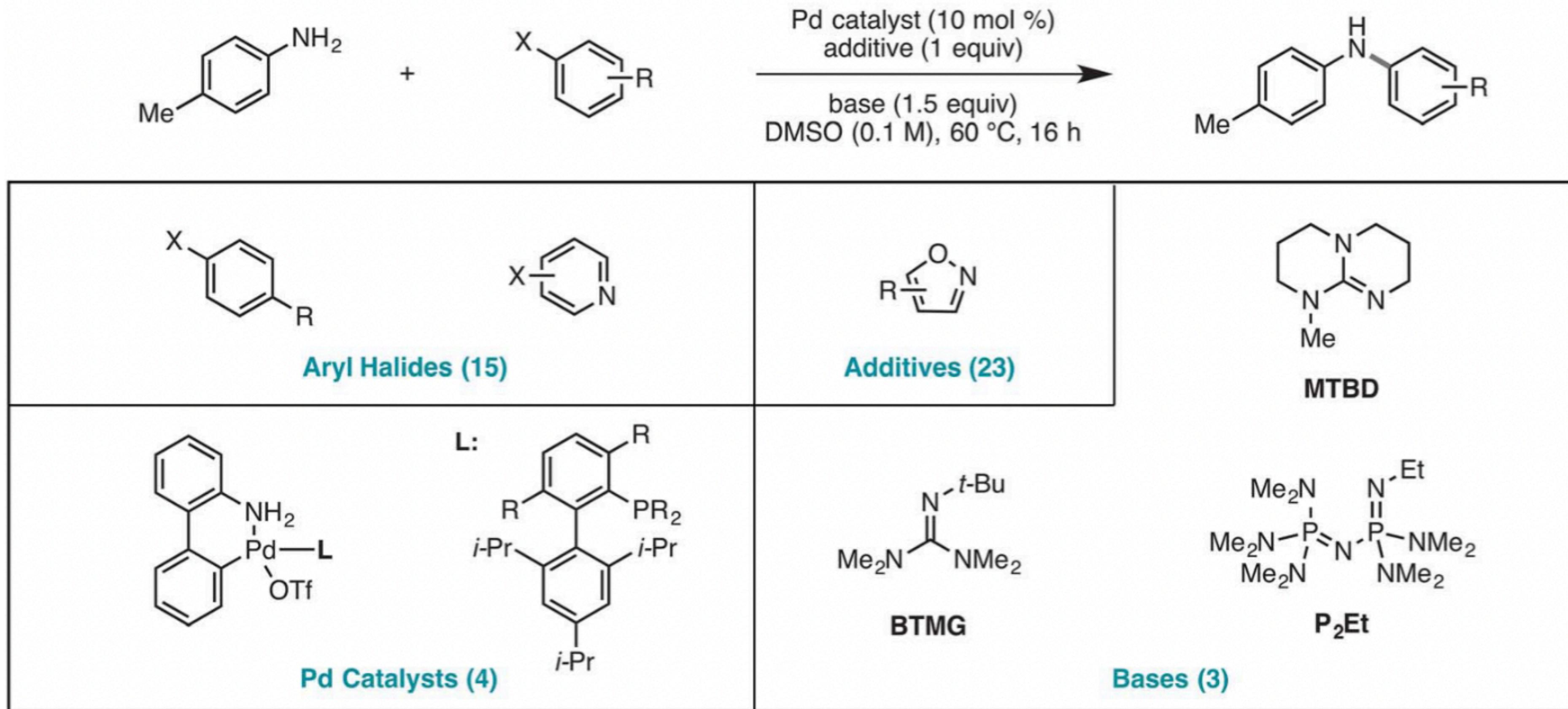
Predicting reaction performance in C-N cross-coupling using machine learning

Derek T. Ahneman,¹ Jesús G. Estrada,¹ Shishi Lin,²
Spencer D. Dreher,^{2*} Abigail G. Doyle^{1*}

Let's Meet The Dataset



UNIVERSITY OF
CAMBRIDGE



https://github.com/doylelab/rxnpredict/blob/master/data_table.csv