# An overview of the impact of machine learning in molecular sciences.

AI Training Series – Fall 2022
11/8/2022
Álvaro Vázquez Mayagoitia
Argonne Nat Lab – CPS Division

# ACKNOWLEDGMENTS

# DoE HPC Roadmap: Exascale computing project (2021-2025)



Frontier AMD CPU,
AMD GPU

Perlmutter AMD CPU,
Nvidia GPU

Aurora Intel CPU,
Intel GPU
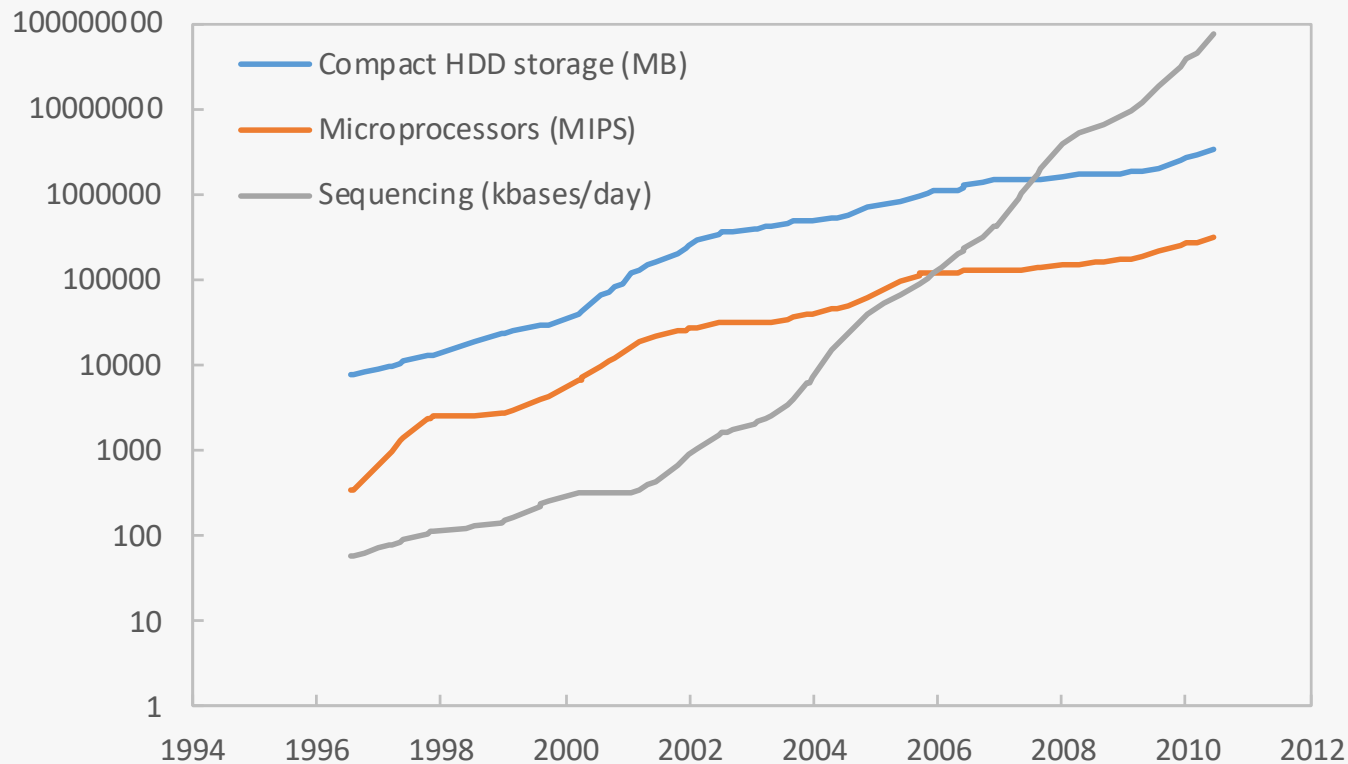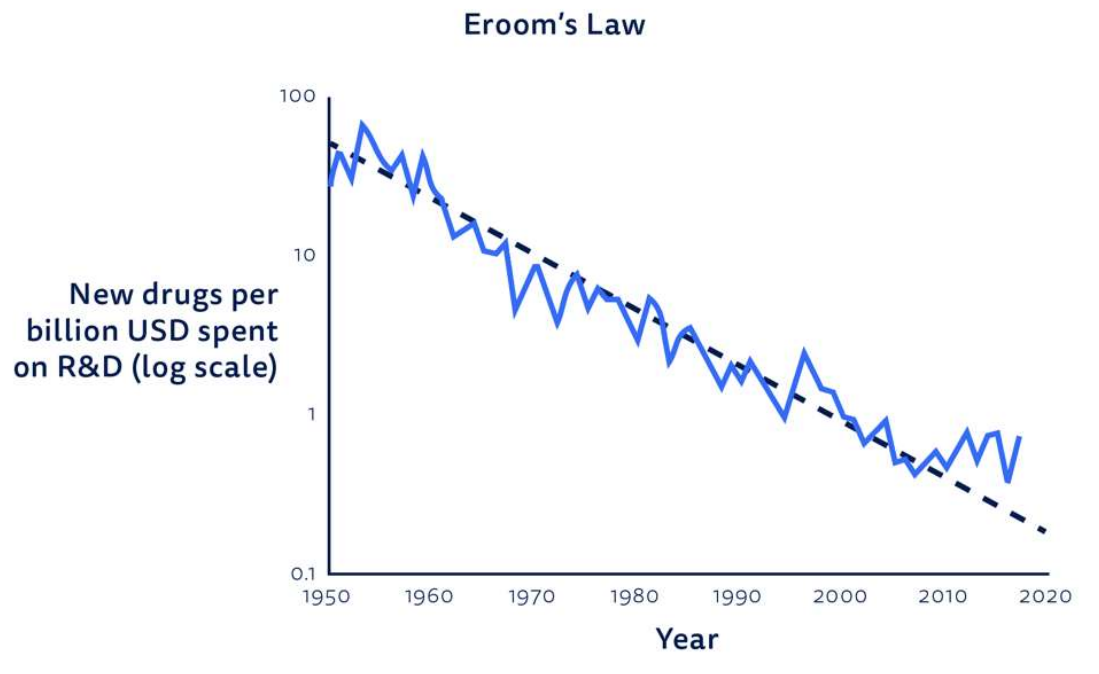
El Capitan AMD CPU,
AMD GPU

Argonne
NATIONAL LABORATORY

Moore's Law: The number of transistors on microchips has doubled every two years

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important for other aspects of technological progress in computing – such as processing speed or the price of computers.

# Evolution of DNA sequencing



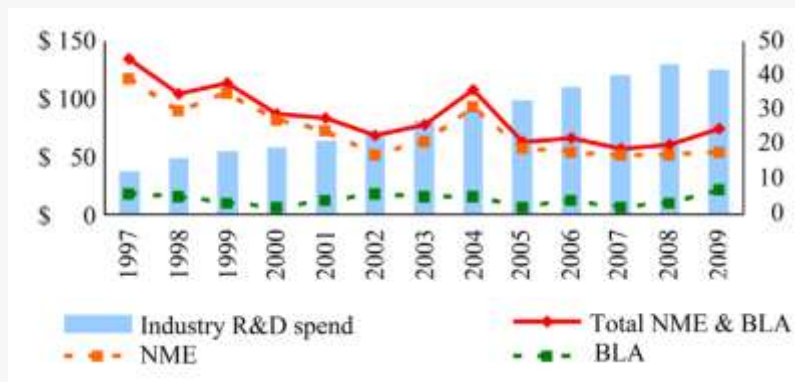In 2020's: trillion bases a day

Argonne
NATIONAL LABORATORY

# Decline in Pharmaceutical R&D efficiency

The **cost** of developing a new drug (~$2-3B) roughly **doubles every nine years**.



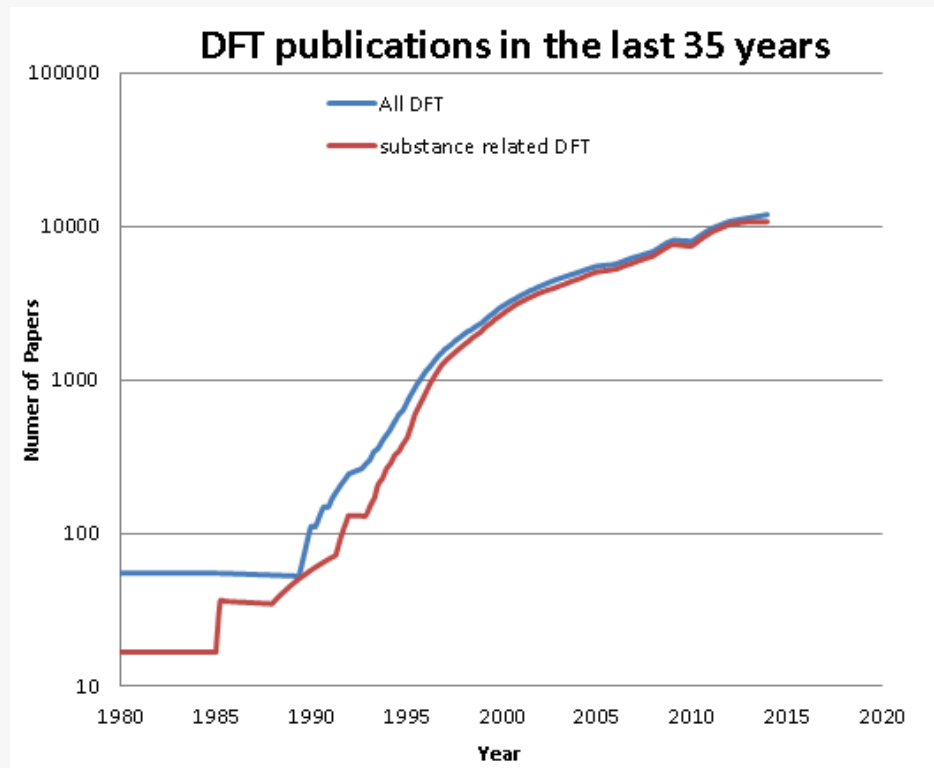Scannell et al. Nature Reviews Drug Discovery, 2012, 11, 191-200
Olexandr Isayev http://olexandrisayev.com

# R&D in Drug Discovery





NME: New Molecular Entities
BLA: Biological App.

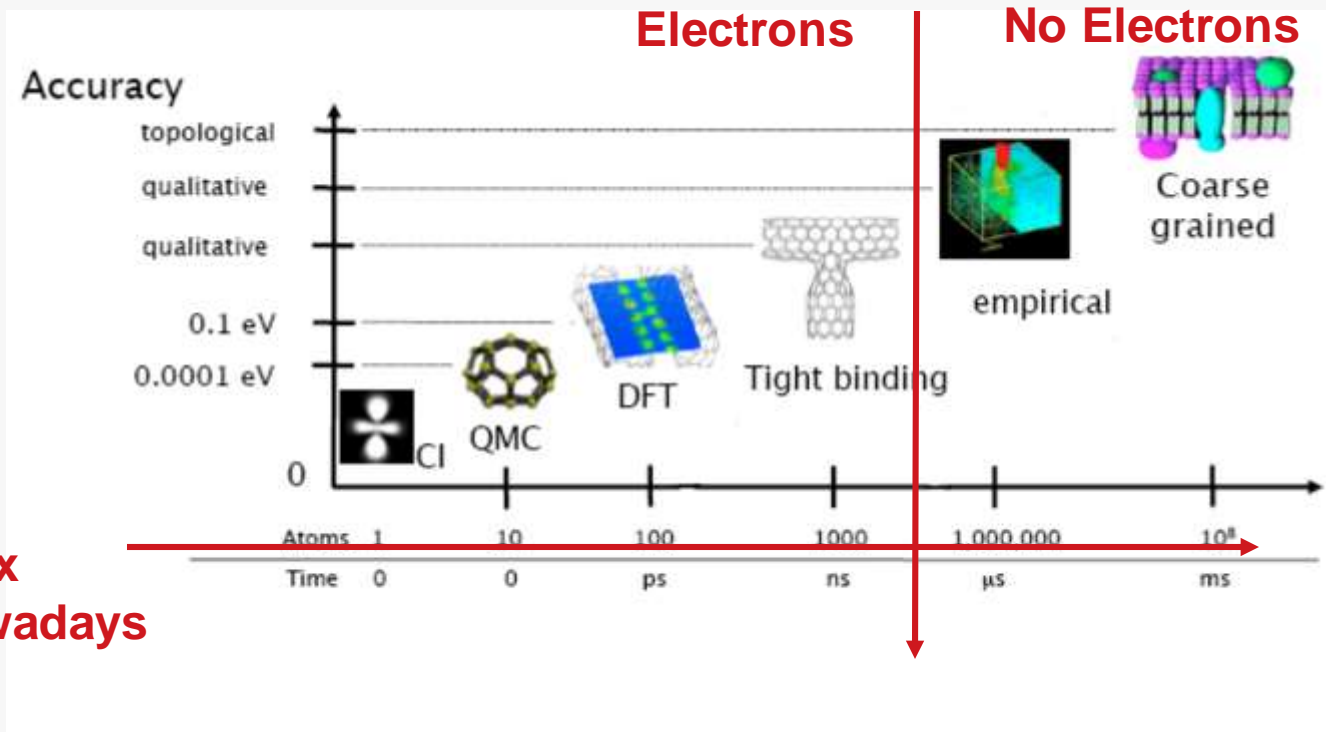Acta Pharmaceutica Sinica B 2014;4(2):112–119

# MATERIALS SCIENCE MODELING



*Web of science queries

Argonne
NATIONAL LABORATORY

# Substances in CAS registry (ACS)



155M in Aug 2019

Are they any useful?

~~10M~~ 50M of new substances per year in average

250M in Abril 2021

Argonne
NATIONAL LABORATORY

**Electrons**   **No Electrons**

Accuracy

topological — — — — — — — — — — — — — — — — — — — — — — — —

qualitative — — — — — — — — — — — — — — — — — — — — — — — Coarse grained

qualitative — — — — — — — — — — — — — — — — — — — — — —

0.1 eV — — — — — — — — — — — — — — — — empirical

0.0001 eV — — — — — — — Tight binding

DFT

CI   QMC

0

| Atoms | 1 | 10 | 100 | 1000 | 1.000.000 | $10^8$ |
|-------|---|----|-----|------|-----------|--------|
| Time | 0 | 0 | ps | ns | μs | ms |

**100x nowadays**

Argonne ▲
NATIONAL LABORATORY

$$\hat{H}(\mathbf{R})\Psi(\mathbf{R},\mathbf{r}) = E\Psi(\mathbf{R},\mathbf{r})$$

Equation that describes the properties of an atom-scale system (time independent)



This is a hyper surface that can be approximated with empirical models.
ML models can remove bias and deal with very complex functions.



**REVIEWS**

Structure prediction drives materials discovery

Computational Science Division

# Data driven AI/ML interatomic potentials for large scale multi physics simulations

- Currently developing AI/ML models, to predict with QM accuracy, energies and forces for in- and out- equilibrium, comparable to experimental observations.
- Our models will ultra fast infer ( using GPU and FPGA accelerators) properties for MD and PIMD.

X-ray data (densities, crystal structure, temperature changes etc) → Fit of initial parameters to reproduce experimental setup

Simulations in ALCF to support observations, and experiments with different compositions ← Training data generation, and AI models production

Feedback to experiments

Water



Hafnia



Graphene



Perovskites

Molten salts

Silica

Single atom catalysis

...

Argonne
NATIONAL LABORATORY

Project: ML inter atomic potential for all hafnia phases.



(1) Diffraction Experiment

Crystal structure or Pair Distribution Function

(2) Phase space exploration

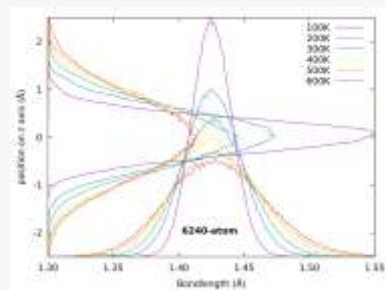Clustering based active learning

KS-DFT AIMD Single point

ML-IP driven MD

• Train data    • Test data

(3) Iterative training
ML-IP fitting

PRL,  2021, Editor's Pick

Argonne
NATIONAL LABORATORY

# VAE-MNIST

## CVAE-MNIST

Argonne
NATIONAL LABORATORY

# VAE jointly trained with a Regressor



Bombarelli *et al.,ACS Cent. Sci*., 2018, 4 (2), pp 268–276

Computational Science Division

# Funnel approach

Virtual screening of the
chemical space



Computational
cost

Molecules most likely
to be of interest

Argonne
NATIONAL LABORATORY

# Informed high throughput computing

**Force Fields**
- Global minima
- 2D to 3D

**Semiempiricals**
- Conformers
- Optimization
- Vibrations

**DFT**
- Geometry
- Excitations
- Multipoles

**Coupled Cluster**
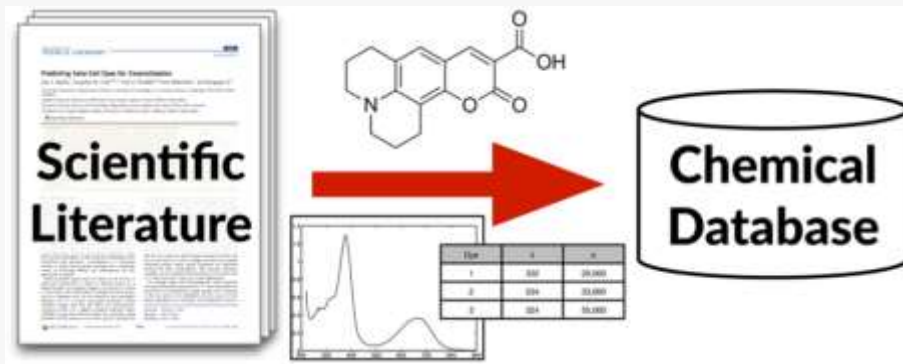- Excitations
- Multipoles

Composite of codes:
- Babel
- Rdkit
- MOPAC
- ORCA
- NWChem

A Toolkit for Automated Extraction of Chemical Information from the Scientific Literature
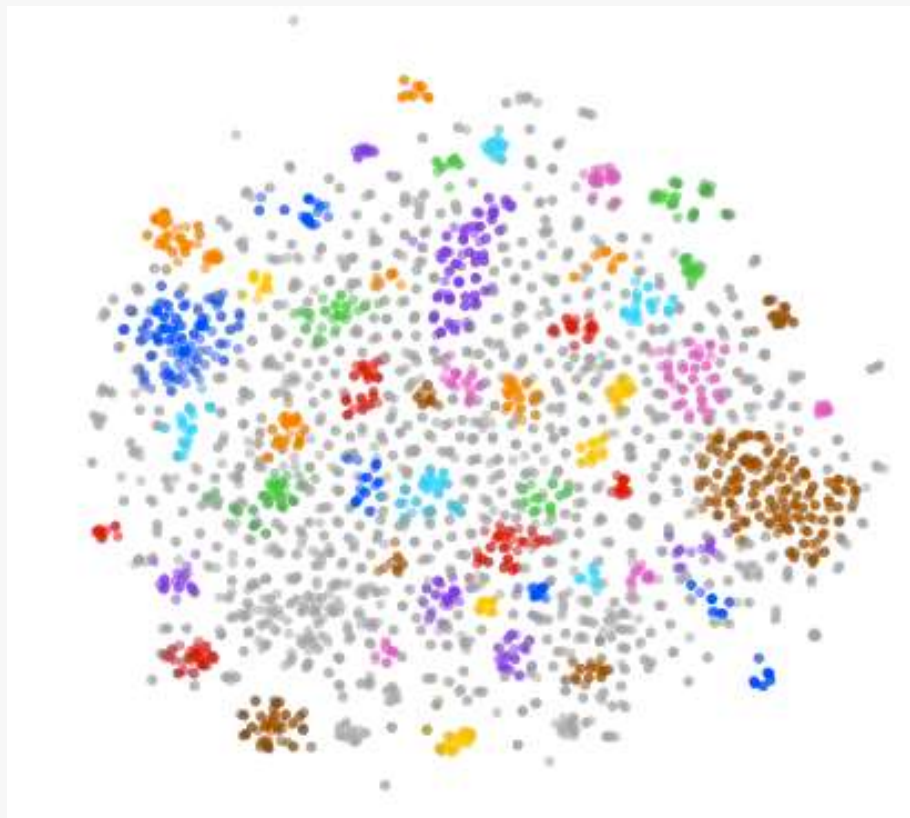
http://chemdataextractor.org

Swain, M. C., & Cole, J. M. J. Chem. Inf. Model. 2016

# NATURAL LANGUAGE PROCESSING PIPELINE

Swain, M. C., & Cole, J. M. J. Chem. Inf. Model. 2016

# Molecular cartography - clusters



Dimension reduction with t-SNE
Clustering with HDBSCAN

Argonne
NATIONAL LABORATORY

# **Learn from data and feedback to experiments**

## **Transition prediction**

TDDFT gap prediction – We used Gaussian Process and Circular Morgan Fingerprints to predict the first transition of the a reduce scale TDDFT (sTDA//wB97X-D3/TZVP), we found that this value is predictable. Similar result found for HOMO-LUMO DFT gap.

Argonne
NATIONAL LABORATORY

# Q&A