

## Science Use Case 3

# Data driven materials discovery for dye sensitized solar cells

ATPESC  
Aug 9, 2019

Álvaro V Mayagoitia  
Argonne CPS

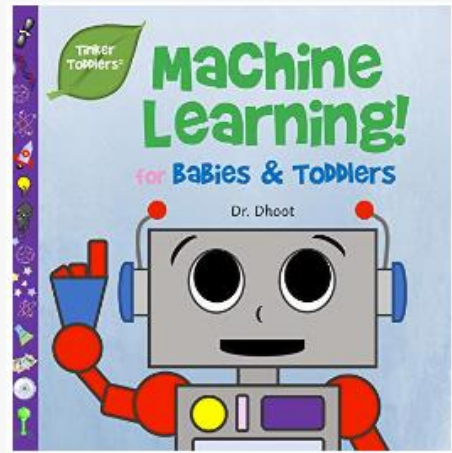
# Acknowledgement



## ALCF Acknowledgement

This research used resources of the Argonne Leadership Computing Facility, which is a DOE Office of Science User Facility supported under Contract DE-AC02-06CH11357.

# My Summer project



“...make any problem a machine learning problem..”

History of Science,  
summarized by Jim Cray

# 4<sup>th</sup> Paradigm Data-Intensive Scientific Discovery



## The FOURTH PARADIGM

DATA-INTENSIVE SCIENTIFIC DISCOVERY

EDITED BY TONY HEY, STEWART TANSLEY, AND KRISTIN TOLLE

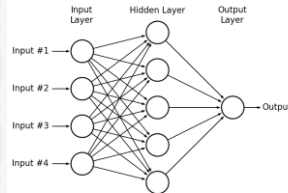
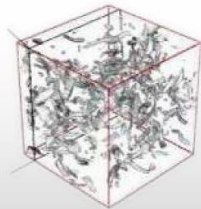
Experiments



Theory

$$\left(\frac{\dot{a}}{a}\right)^2 = \frac{4\pi G\rho}{3} - K \frac{c^2}{a^2}$$

Computation



Data science

1600

1950

2000



<https://www.microsoft.com/en-us/research/publication/fourth-paradigm-data-intensive-scientific-discovery/>

Computational Science Division

Argonne  
NATIONAL LABORATORY



# Chemical Compound Space

\*Estimated Energetically  
Possible Organic Molecules

$>10^{60}$

\*Nature, Insights, 2004



# Chemical Compound Space

\*Estimated Energetically  
Possible Organic Molecules  
>10<sup>60</sup>

\*Nature, Insights, 2004

Nature,  
Insights,  
2004

Total number of  
water molecules in  
Earth: 10<sup>40</sup>



Recent estimations say could exceed 10<sup>180</sup>

# Chemical Compound Space

\*Estimated Energetically  
Possible Organic Molecules  
>10<sup>60</sup>

\*Nature, Insights, 2004

Nature,  
Insights

Total number of  
water molecules in  
Earth: 10<sup>40</sup>



How much do we know of  
the chemical space?

Compiled from experiments  
since the early 1800s: 10<sup>8</sup>

Computationally:

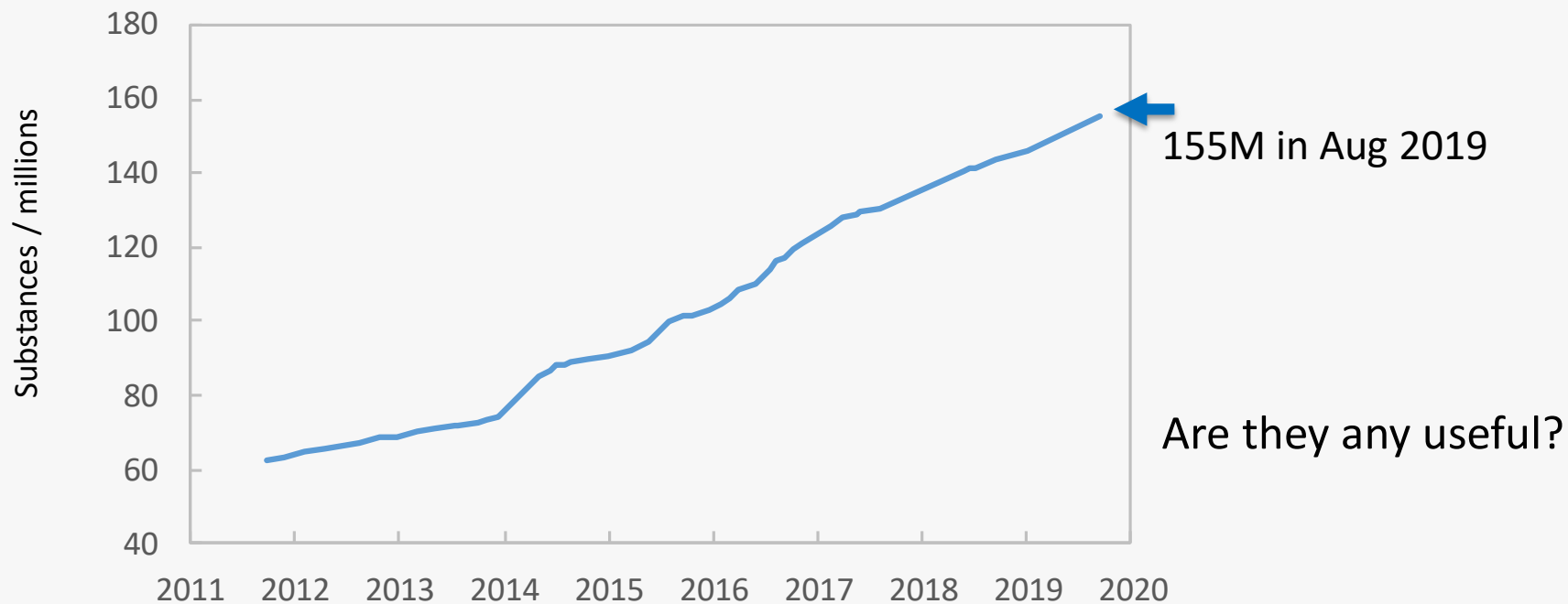
eg. Harvard Clean Energy  
project

10<sup>7</sup> Molecules

10<sup>7</sup> CPU hrs

10<sup>9</sup> Calculations

# Substances in CAS registry (ACS)

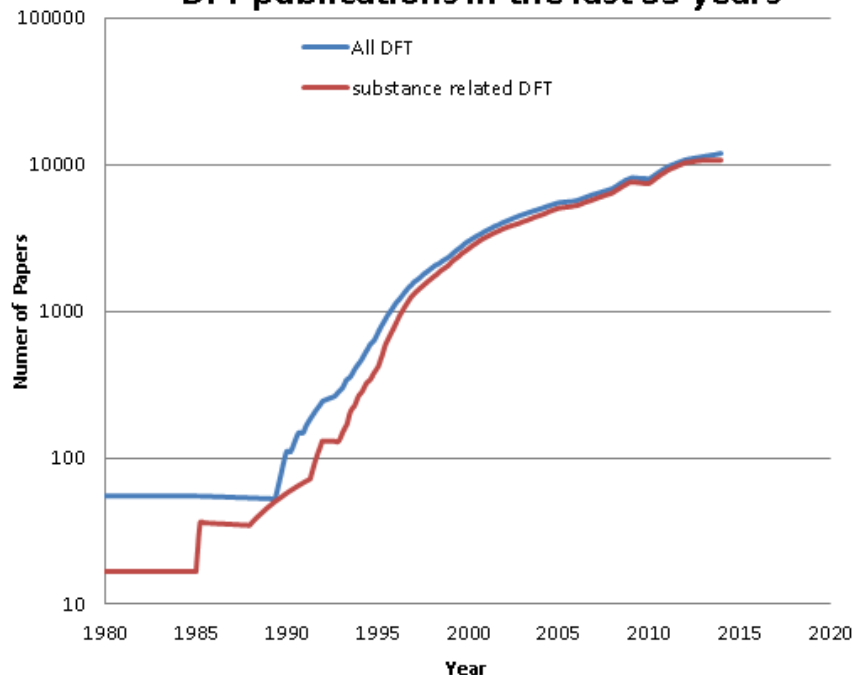


10M of new substances  
per year in average



# MATERIALS SCIENCE MODELING

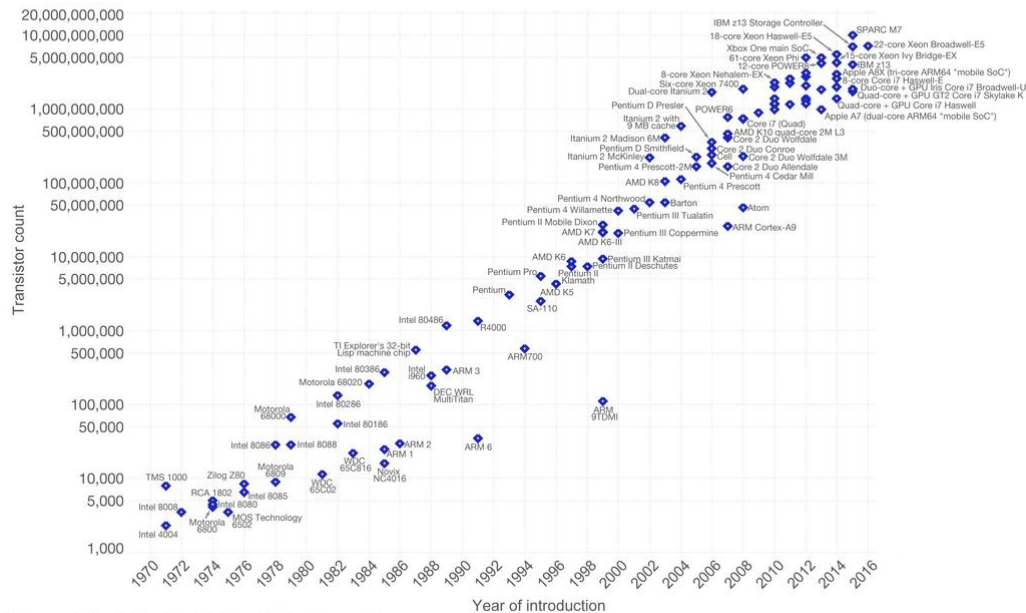
## DFT publications in the last 35 years



## Moore's Law – The number of transistors on integrated circuit chips (1971-2016)

Moore's law describes the empirical regularity that the number of transistors on integrated circuits doubles approximately every two years. This advancement is important as other aspects of technological progress – such as processing speed or the price of electronic products – are strongly linked to Moore's law.

Our World  
in Data



Data source: Wikipedia ([https://en.wikipedia.org/wiki/Transistor\\_count](https://en.wikipedia.org/wiki/Transistor_count))

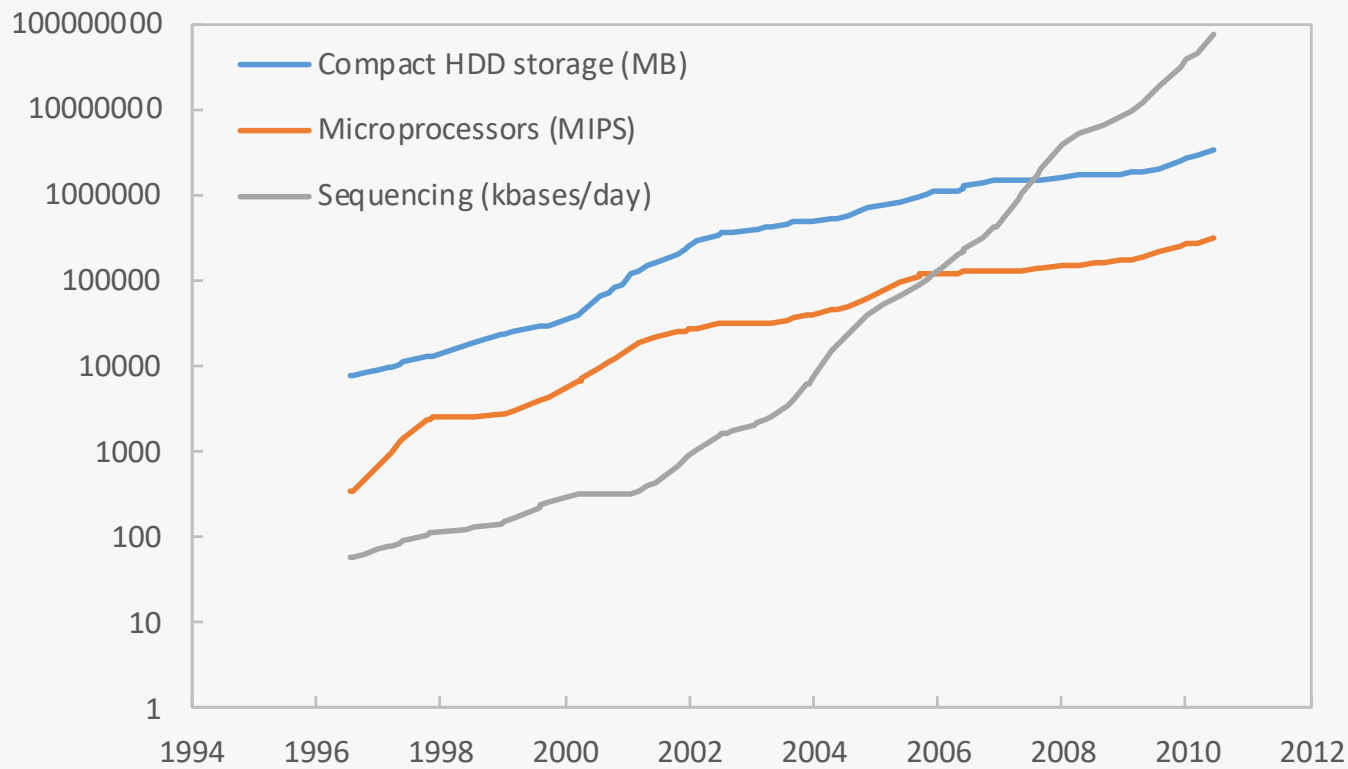
The data visualization is available at [OurWorldinData.org](https://ourworldindata.org). There you find more visualizations and research on this topic.

Licensed under CC-BY-SA by the author Max Roser.

\*Web of science queries

Computational Science Division

# Evolution of DNA sequencing

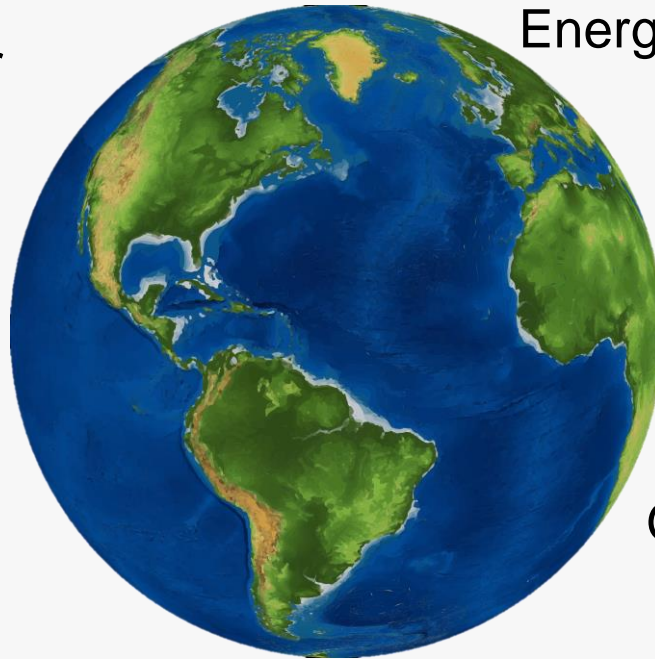


June 2019 : trillion bases a day

# Renewable sources of energy

Solar

$10^5$  TW at earth surface  
**10,000 TW** tech. value



Energy needed 2007: **15 TW**  
2017: **18 TW**  
2050: **~30 TW**  
2100: **~50 TW**

Wind

**14 TW**

Biomass **5-7 TW**

Geothermal **1.9 TW**

Tide/Ocean **0.7 TW**

*The Third Industrial Revolution* by Jeremy Rifkin, 2011

[IRENA report 2018](#)



# Cost per watt-hour of Solar energy

1977

\$76.67

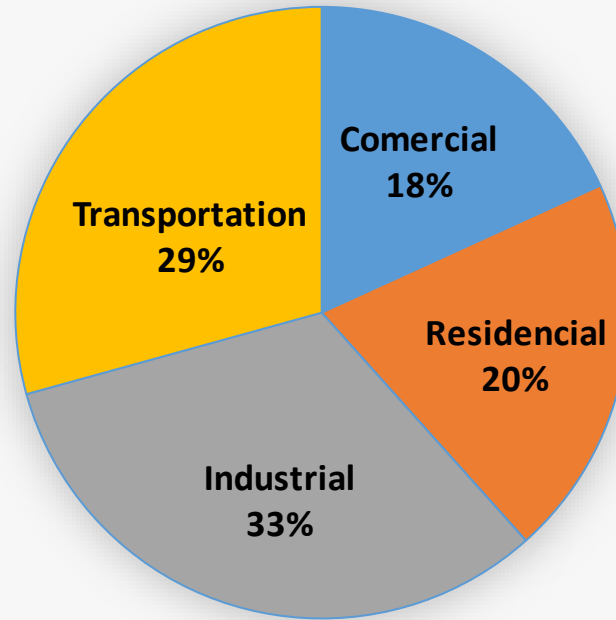


2019

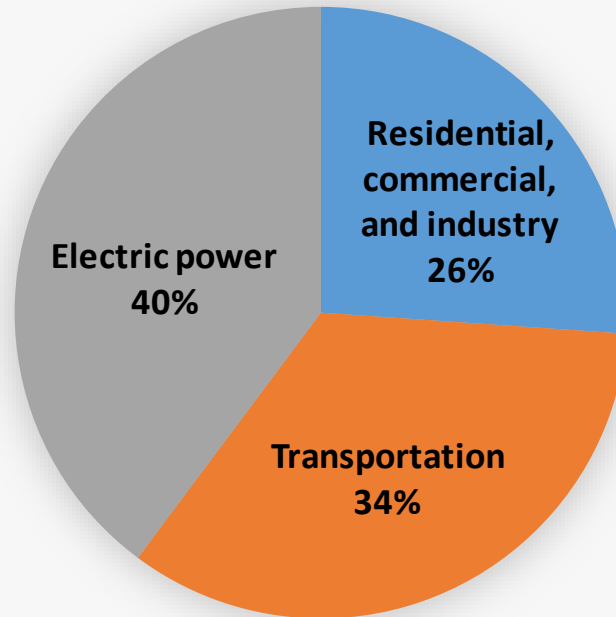
\$0.10



# Share of total US energy consumption by end-use sectors 2018



# Carbon footprint per sector 2018







Artificial Satellites



SF bus shelters



Food courts

Electric golf carts



# Solar windows

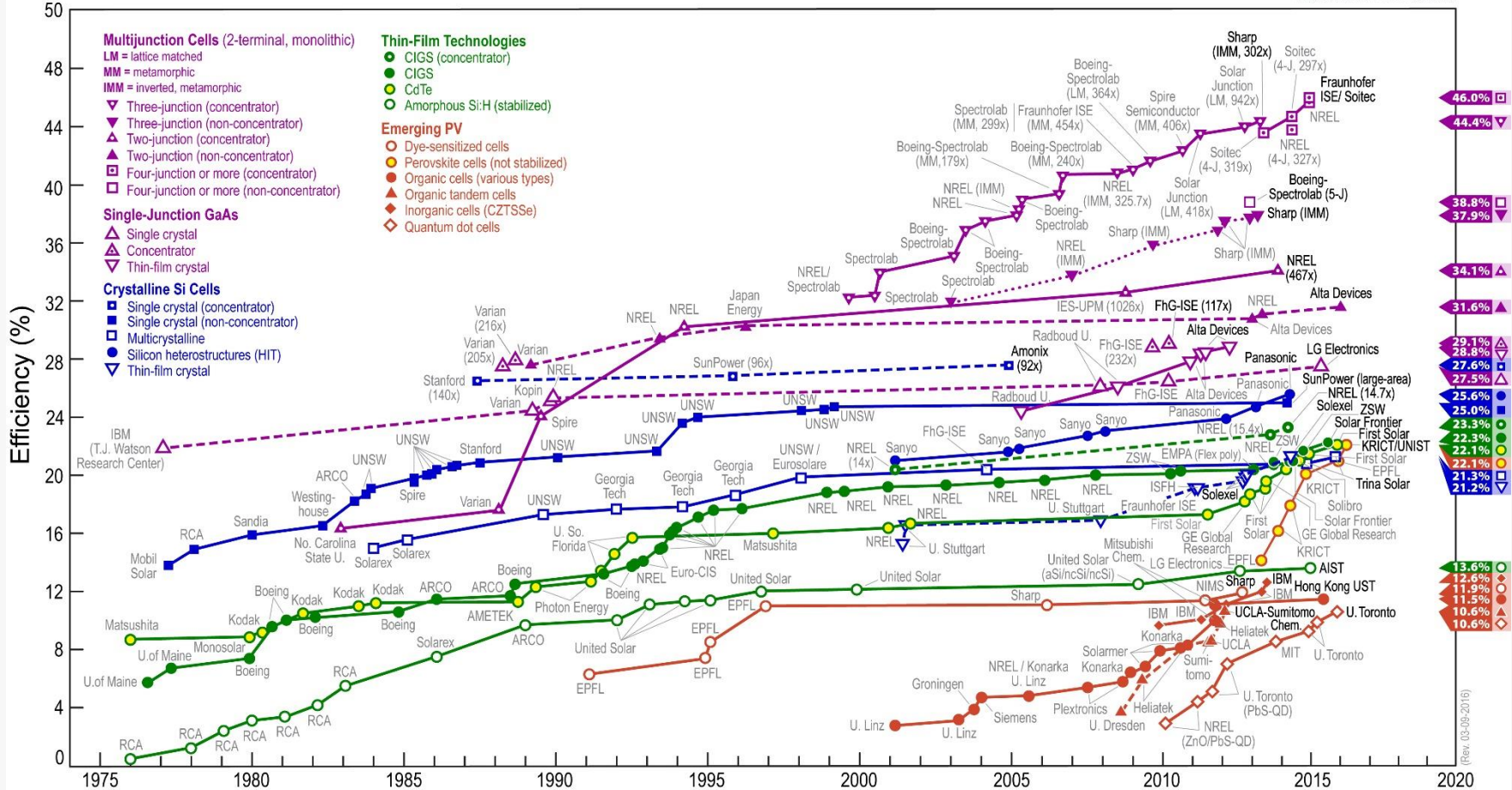


Harris Theatre – Chicago  
[north Millenium park]



# EVOLUTION OF SOLAR CELLS - ENERGY.GOV

## Best Research-Cell Efficiencies

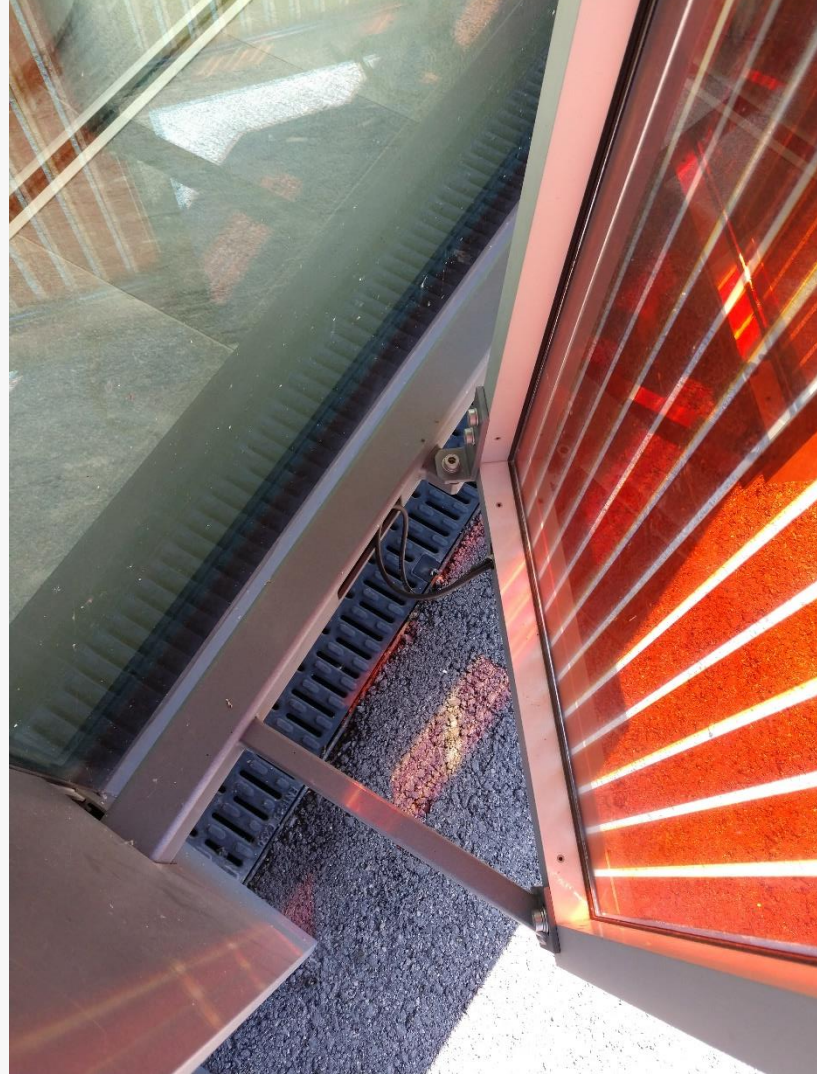
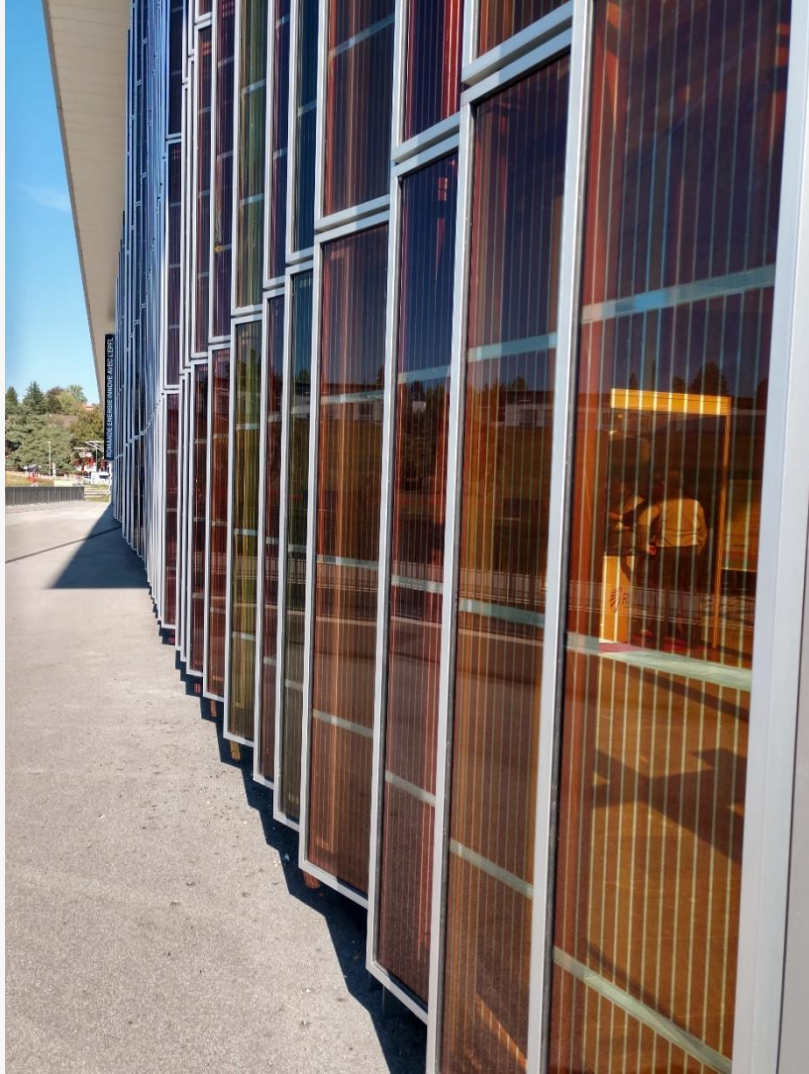




# Solar windows



SwissTech Convention Center EPFL



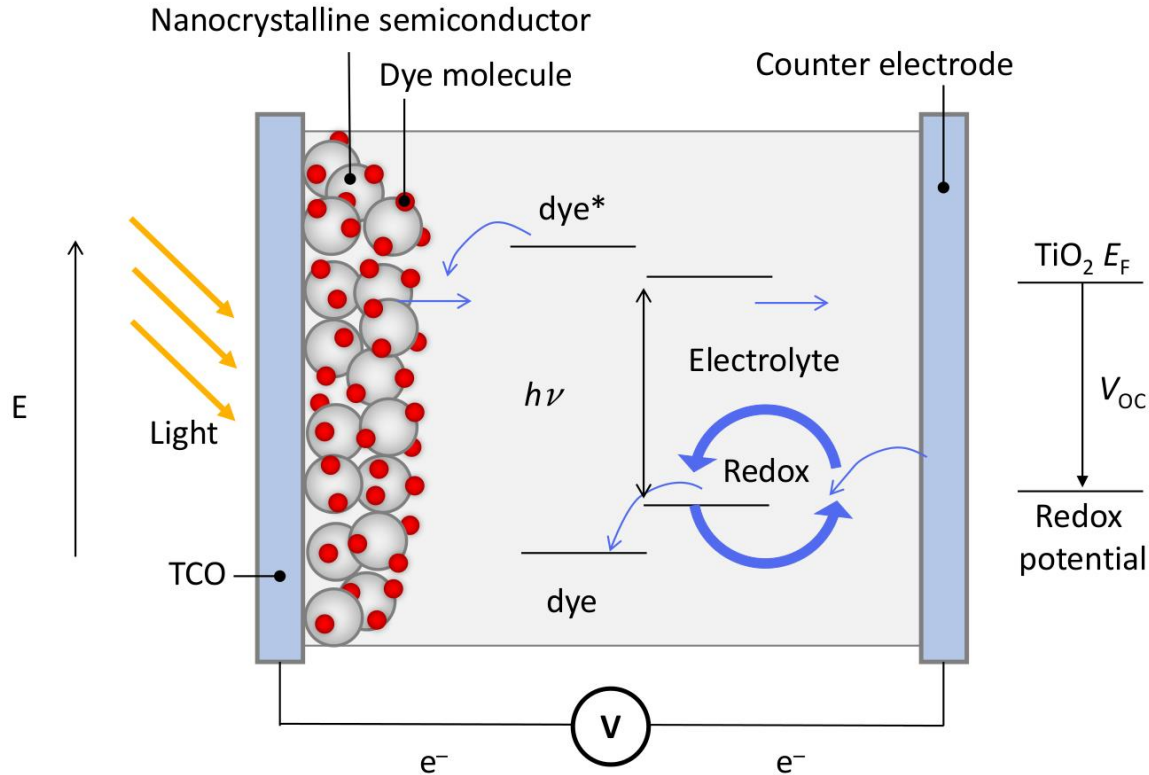






# dye-sensitized SOLar cell

## Operational Mechanism



## A low-cost, high-efficiency solar cell based on dye-sensitized colloidal TiO<sub>2</sub> films

Brian O'Regan\* & Michael Grätzel†

Institute of Physical Chemistry, Swiss Federal Institute of Technology, CH-1015 Lausanne, Switzerland

**1991 Conversion efficiency: 7.9%**

# DYE-SENSITIZED SOLAR CELL

- Emerging technology
- Cost effective (good price-to-performance ratio) 😊
- Less efficient than Si-base cells ☹️

## Organometallic cells

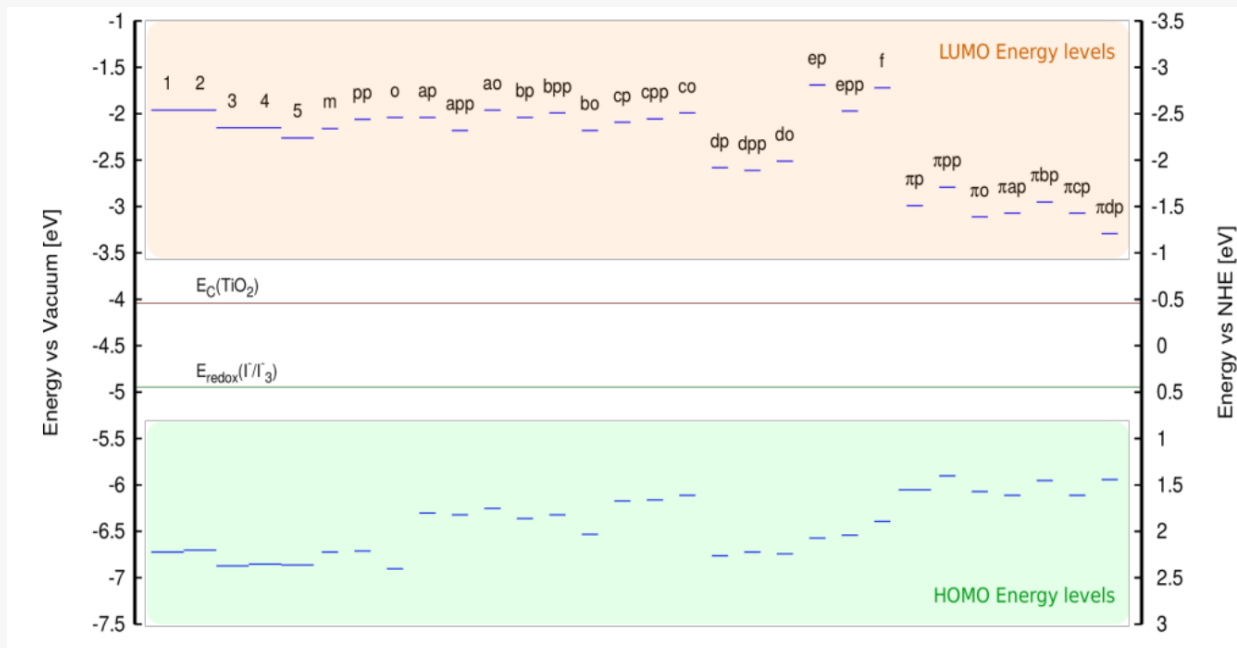
- Nazeeruddin et al JACS 1993 (ruthenium-based) N719: 10.4% eff.
- Yella et al Science 2011 (Zn-porphyrin-based) : 12.3 % eff.
- Burschka et al Nature 2013 (lead-iodide-based): 15% eff.

## Organic cells

- Daeneke et al Nat. chem. 2011 (carbazole-base): 7.5% eff.
- Zeng et al Chem. Mater. 2010 (thiophene-base): 10.3% eff.

# ENCODING STRUCTURE-FUNCTION

Screening with TDDFT is costly



Rules:

$$\epsilon_{LUMO} > E_{C\text{TiO}_2}$$

$$\epsilon_{HOMO} < E_{\text{Electrolyte}}$$

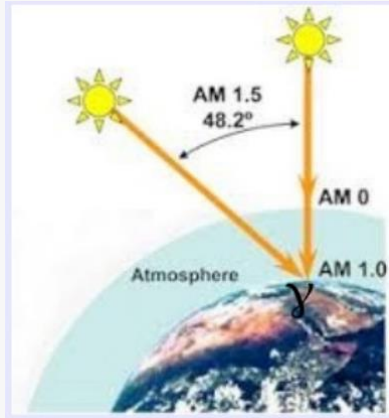
# MAXIMIZING LIGHT HARVESTING EFFICIENCY

$$AM = \frac{L}{L_0} \approx \frac{1}{\cos z}$$

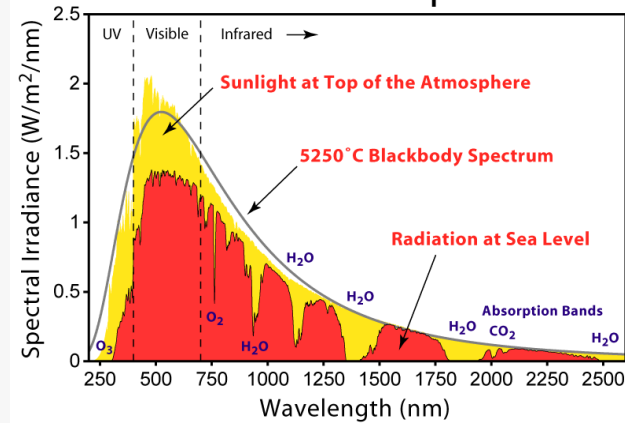
$L_0$  = zenith path length

$L$  = path length to the atmosphere

$z$  = zenith angle



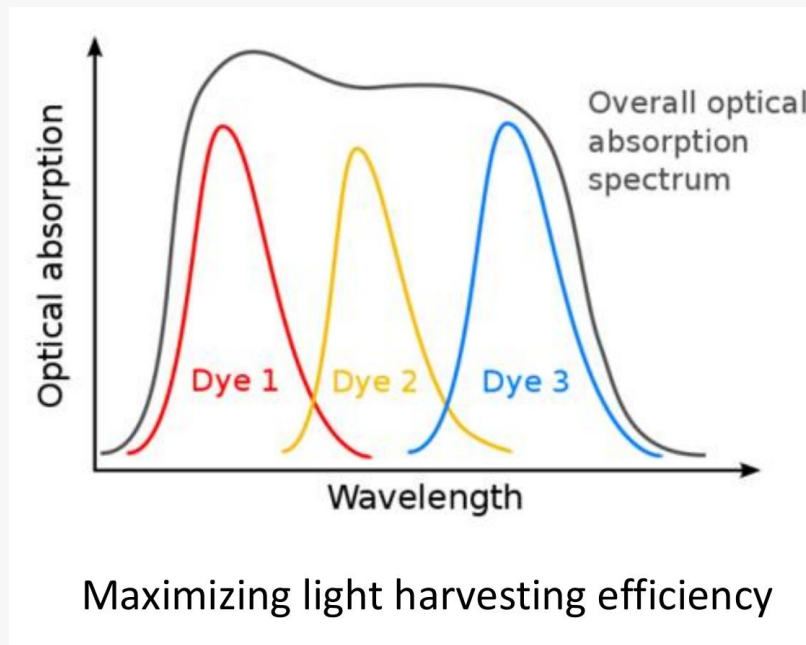
## Solar Radiation Spectrum



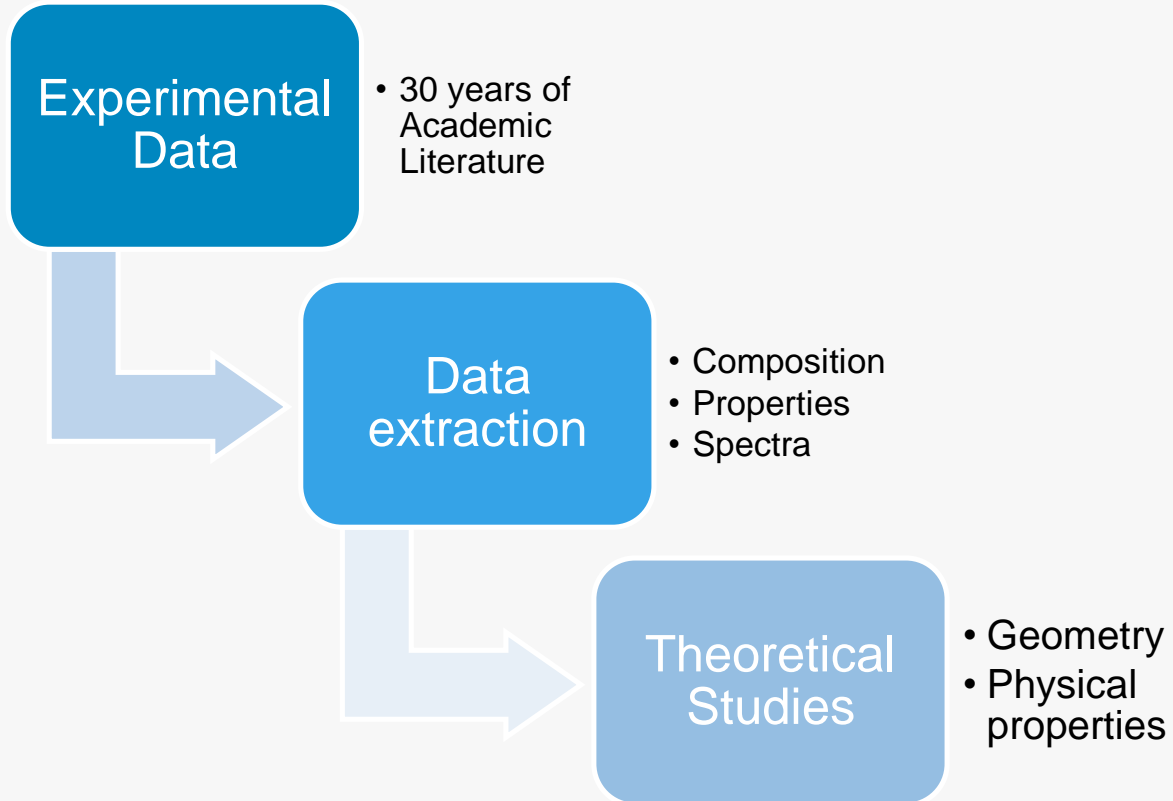
[https://en.wikipedia.org/wiki/Air\\_mass\\_\(solar\\_energy\)](https://en.wikipedia.org/wiki/Air_mass_(solar_energy))



# MAXIMIZING LIGHT HARVESTING EFFICIENCY



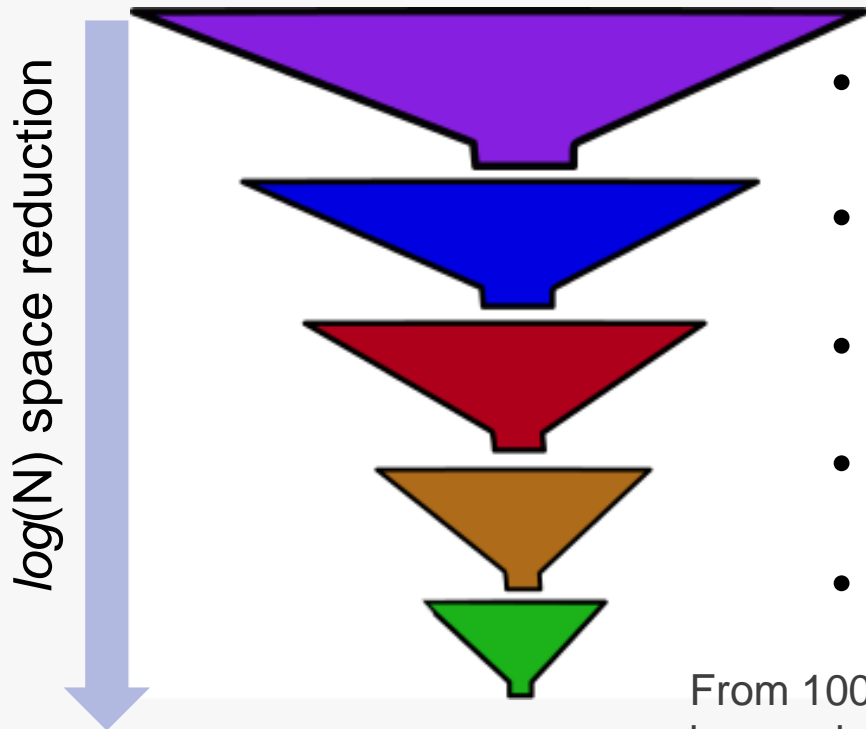
# Work flow



# Funnel approach

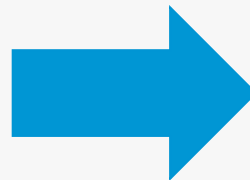
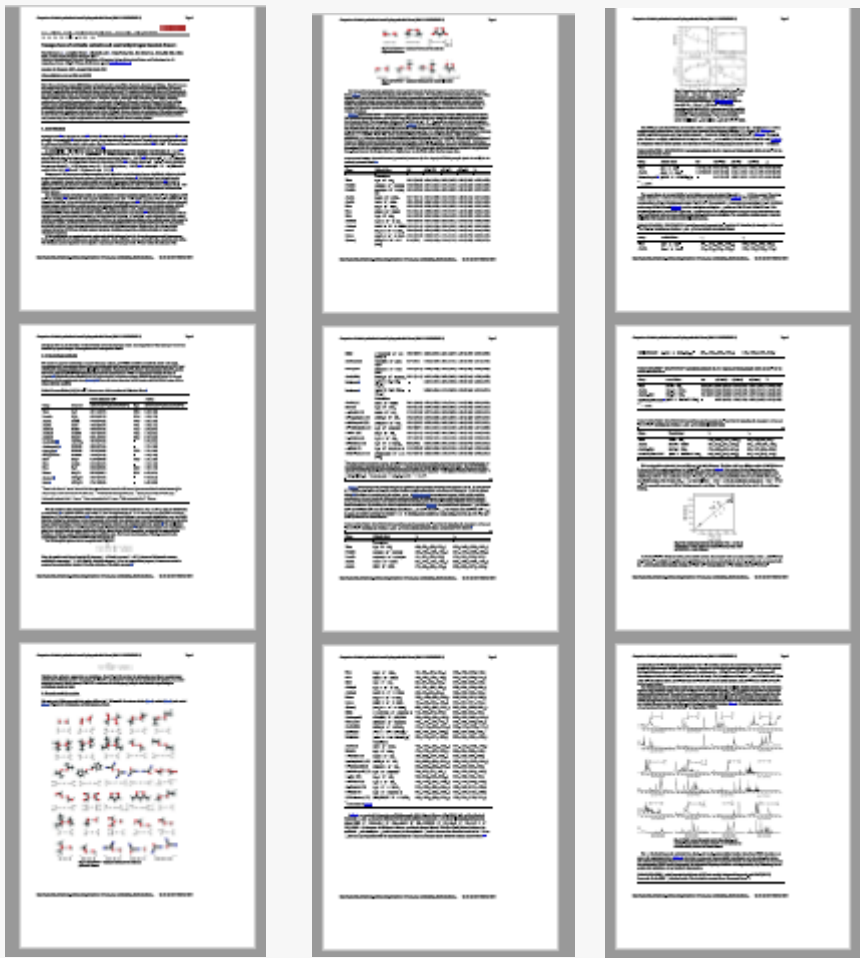
## Screening the chemical space

Funnel / Filters



- Size, spectra, charges
- Optoelectronics rules
- Semi-empirical Methods
- Density Functional Methods
- “Gold” standard methods

From 100 kilo molecules, which ones could be good dyes?

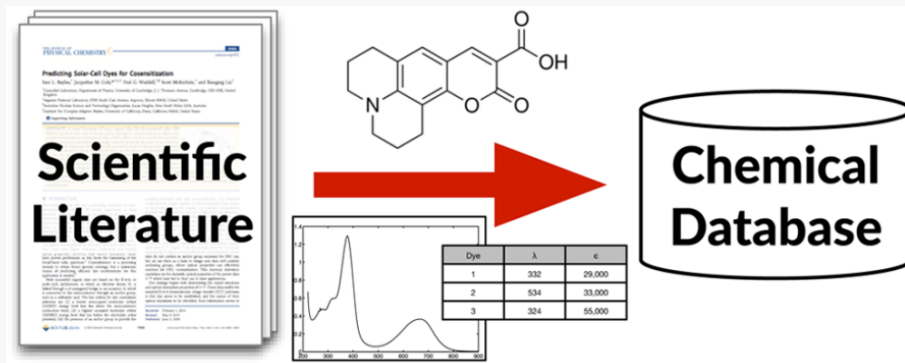


Compound	Lamda max	Ext. Coef.
Triphenylmethane	480	4320
Indigo	613	320
anthraquinone	320	-----

How long would it take a person to get this information?



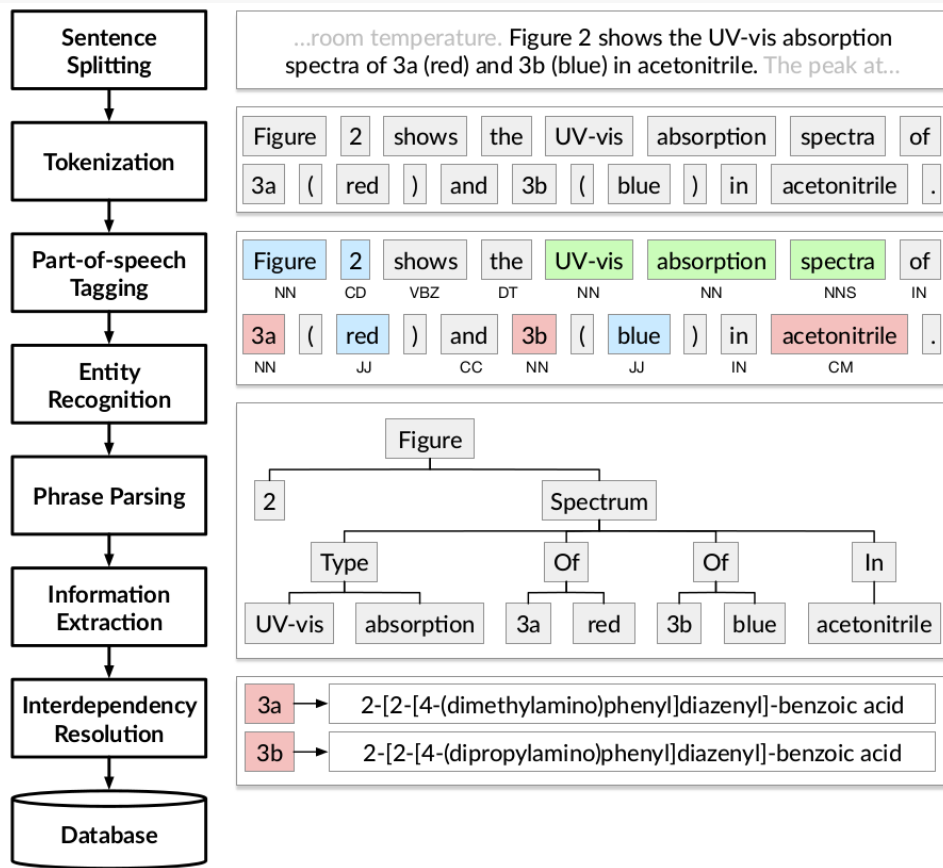
# A Toolkit for Automated Extraction of Chemical Information from the Scientific Literature



<http://chemdataextractor.org>



# NATURAL LANGUAGE PROCESSING PIPELINE



# DATABASE

The dye 2-[2-[4-(dimethylamino)phenyl]diazenyl]-benzoic acid (**3a**) was added...

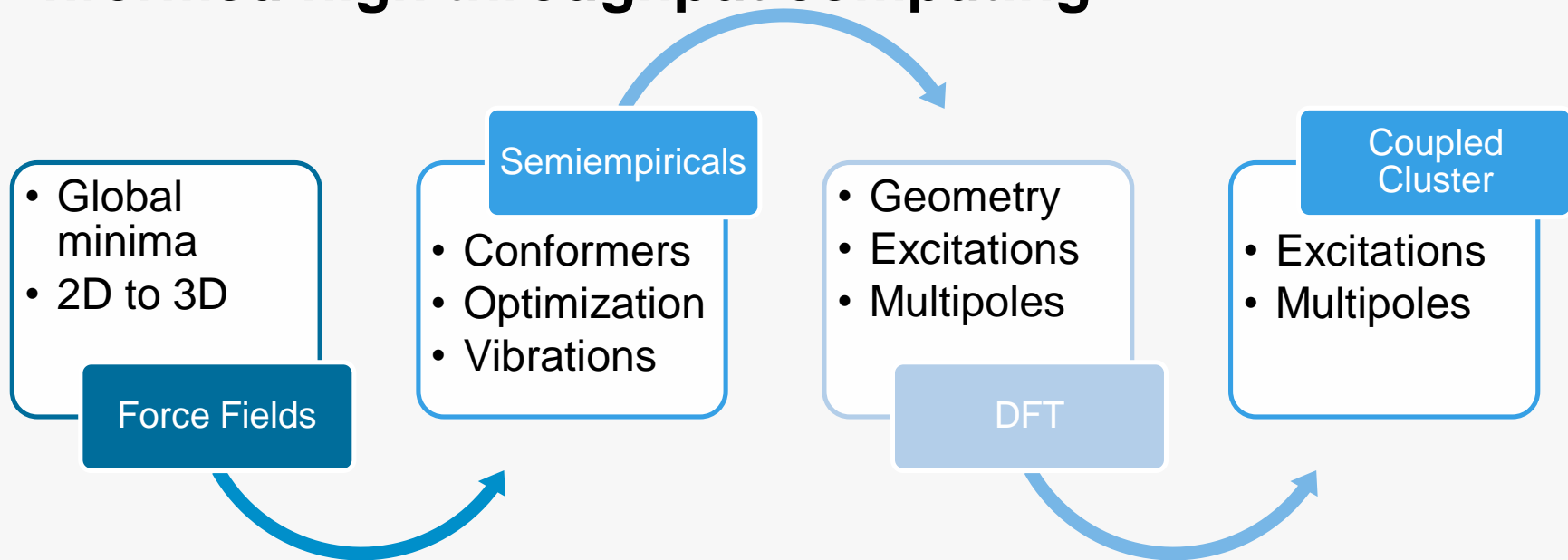
UV-vis spectra were recorded using an Agilent8453 diode array spectrophotometer.

```
{  
  "name": "2-[2-[4-(dimethylamino)phenyl]diazenyl]-benzoic acid",  
  "label": "3a",  
  "uvvis": [ {  
    "solvent": "acetonitrile",  
    "apparatus": "Agilent8453 diode array spectrophotometer",  
    "peaks": [ { "wavelength": "448", "extinction": "29,000" } ],  
  } ],  
}
```

**Figure 2:** UV-vis absorption spectra of 3a in acetonitrile.

Dye	$\lambda_{\max}/\text{nm}$ ( $\epsilon/\text{M}^{-1} \text{cm}^{-1}$ )
3a	448 (29,000)
3b	415 (48,000)

# Informed high throughput computing

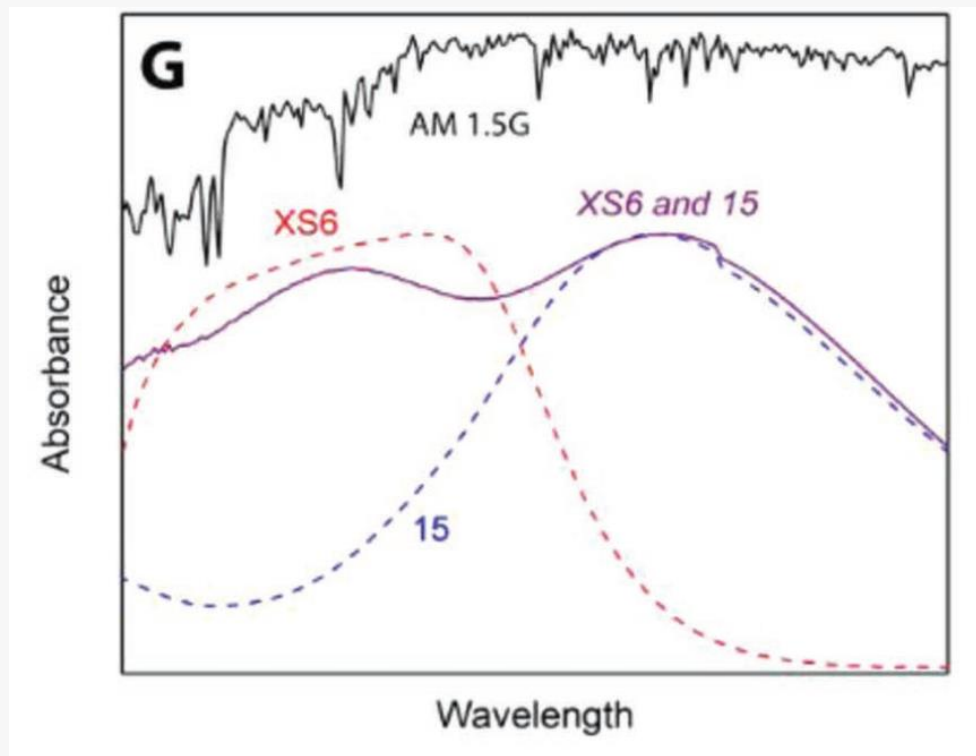


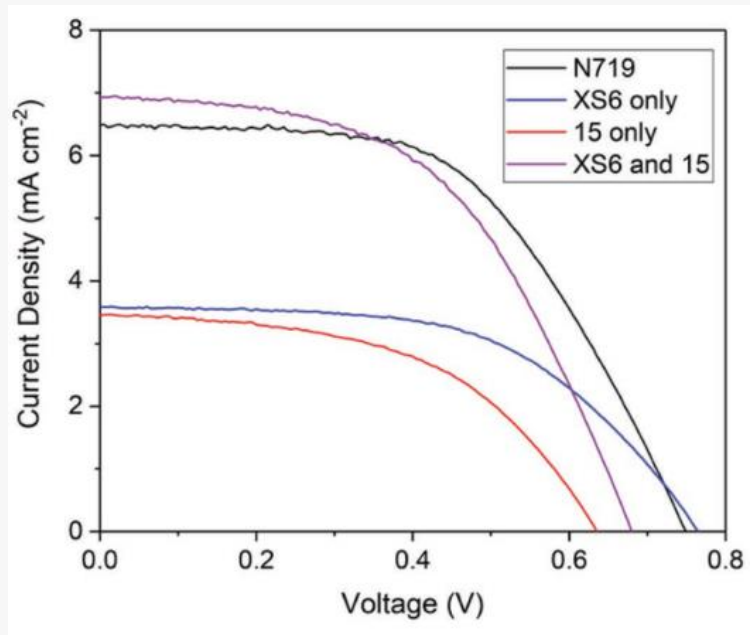
Composite of codes:

- Babel
- Rdkit
- MOPAC
- ORCA
- NWChem

Sample name	AFM parameters				XRR parameters			
	Mean height [nm]	Max height [nm]	Aggregate coverage [%]	Number of aggregates [ $\mu\text{m}^{-2}$ ]	Dye layer thickness [ $\text{\AA}$ ]	$\text{SLD}_{\text{dye}}$ [ $\times 10^{-6} \text{\AA}^{-2}$ ]	Surface roughness [ $\text{\AA}$ ]	Surface coverage [%]
Singly sensitized working electrodes								
C1 only	$5 \pm 1$	$7 \pm 2$	$3 \pm 6$	$2 \pm 3$	$43.5 \pm 0.9$	$6.6 \pm 0.5$	$5.6 \pm 0.7$	$55 \pm 4$
8c only	$5 \pm 1$	$6 \pm 2$	$3 \pm 2$	$3 \pm 2$	$26.6 \pm 0.9$	$5.1 \pm 0.9$	$3.3 \pm 0.8$	$39 \pm 7$
XS6 only	$4.9 \pm 0.4$	$6.0 \pm 0.7$	$1.0 \pm 0.1$	$0.3 \pm 0.2$	$23.6 \pm 0.5$	$8.7 \pm 0.4$	$3.7 \pm 0.5$	$73 \pm 3$
H3 only	$9 \pm 1$	$15 \pm 3$	$0.3 \pm 0.1$	$0.18 \pm 0.05$	$27 \pm 1$	$6.7 \pm 0.5$	$3.7 \pm 0.5$	$55 \pm 4$
15 only	$8 \pm 2$	$15 \pm 3$	$7 \pm 2$	$1.1 \pm 0.4$	$24.3 \pm 0.3$	$7.8 \pm 0.4$	$2.7 \pm 0.3$	$62 \pm 3$
Co-sensitized working electrodes								
<i>C1 then 15</i>	$6 \pm 2$	$10 \pm 3$	$1.3 \pm 0.5$	$0.7 \pm 0.2$	$33.7 \pm 0.5$	$5.9 \pm 0.7$	$3.1 \pm 0.6$	$49 \pm 6$
<i>C1 and 15</i>	$7 \pm 2$	$12 \pm 4$	$2.0 \pm 0.5$	$0.9 \pm 0.6$	$21.5 \pm 0.8$	$6.3 \pm 0.9$	$3.8 \pm 0.7$	$52 \pm 7$
<i>H3 then C1</i>	$8 \pm 2$	$16 \pm 4$	$3 \pm 3$	$0.4 \pm 0.2$	$42 \pm 1$	$6.0 \pm 0.6$	$5.2 \pm 0.7$	$49 \pm 5$
<i>C1 and H3</i>	$5 \pm 1$	$8 \pm 3$	$2 \pm 1$	$2 \pm 2$	$25.4 \pm 0.4$	$8.5 \pm 0.4$	$3.0 \pm 0.5$	$69 \pm 3$
<i>8c then 15</i>	$6 \pm 1$	$9 \pm 2$	$1.1 \pm 0.2$	$0.7 \pm 0.6$	$30.9 \pm 0.4$	$6.9 \pm 0.4$	$3.9 \pm 0.6$	$54 \pm 3$
<i>8c and 15</i>	$4.6 \pm 0.3$	$5.8 \pm 0.4$	$12 \pm 9$	$16 \pm 5$	$31 \pm 2$	$5.7 \pm 0.5$	$7 \pm 2$	$45 \pm 4$
<i>H3 then 8c</i>	$5.5 \pm 0.7$	$8 \pm 1$	$3 \pm 2$	$1 \pm 1$	$37.2 \pm 0.2$	$9.0 \pm 0.7$	$2.9 \pm 0.4$	$70 \pm 5$
<i>8c and H3</i>	$5.2 \pm 0.7$	$7 \pm 2$	$2 \pm 2$	$1 \pm 1$	$27.5 \pm 0.4$	$8.0 \pm 0.4$	$3.3 \pm 0.6$	$63 \pm 3$
<b><i>XS6 then 15</i></b>	<b><math>6 \pm 1</math></b>	<b><math>8 \pm 2</math></b>	<b><math>0.7 \pm 0.3</math></b>	<b><math>0.8 \pm 0.3</math></b>	<b><math>18.8 \pm 0.3</math></b>	<b><math>8.7 \pm 0.5</math></b>	<b><math>3.6 \pm 0.4</math></b>	<b><math>72 \pm 4</math></b>
<b><i>XS6 and 15</i></b>	<b><math>7.8 \pm 0.7</math></b>	<b><math>11 \pm 1</math></b>	<b><math>0.3 \pm 0.1</math></b>	<b><math>0.24 \pm 0.09</math></b>	<b><math>18.6 \pm 0.3</math></b>	<b><math>8.8 \pm 0.5</math></b>	<b><math>3.4 \pm 0.4</math></b>	<b><math>73 \pm 4</math></b>
<i>XS6 then H3</i>	$5.5 \pm 0.7$	$7.6 \pm 0.8$	$0.3 \pm 0.1$	$0.25 \pm 0.04$	$21.0 \pm 0.3$	$9.6 \pm 0.5$	$4.1 \pm 0.4$	$79 \pm 4$
<i>XS6 and H3</i>	$5.3 \pm 0.8$	$7 \pm 1$	$0.3 \pm 0.1$	$0.2 \pm 0.1$	$21.6 \pm 0.6$	$8.7 \pm 0.6$	$4.0 \pm 0.5$	$71 \pm 5$

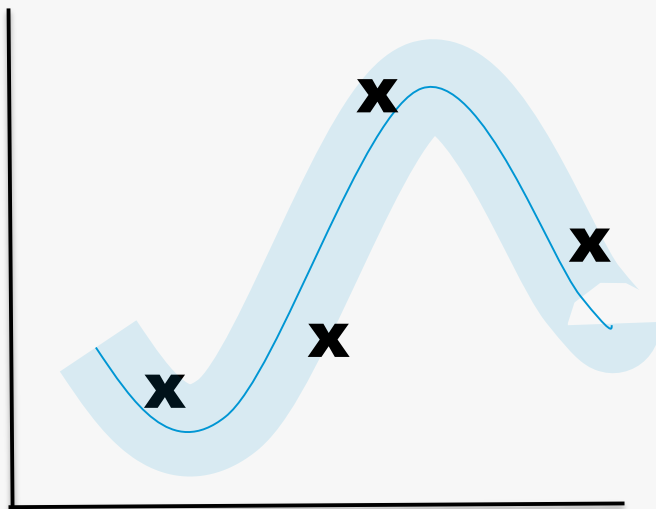
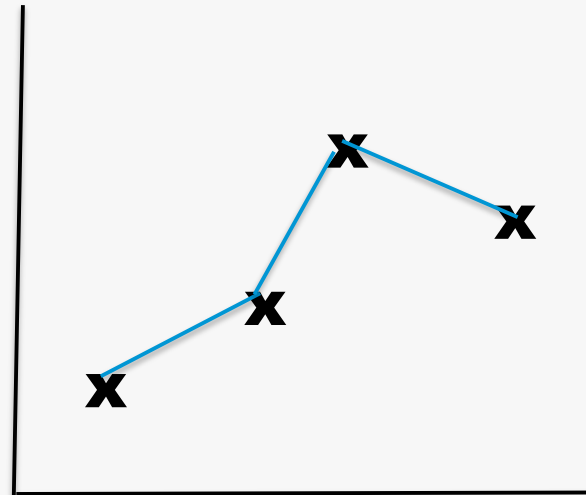
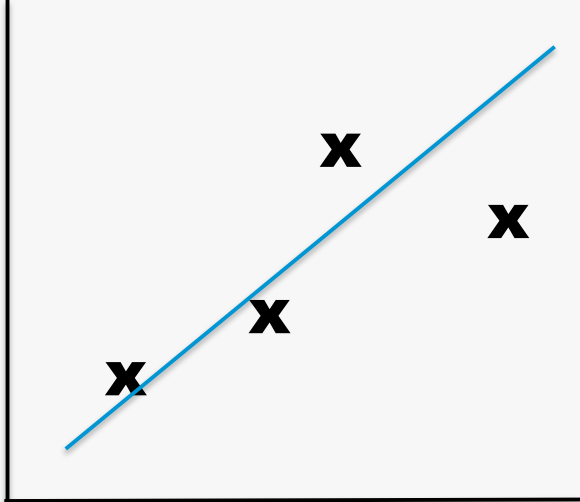


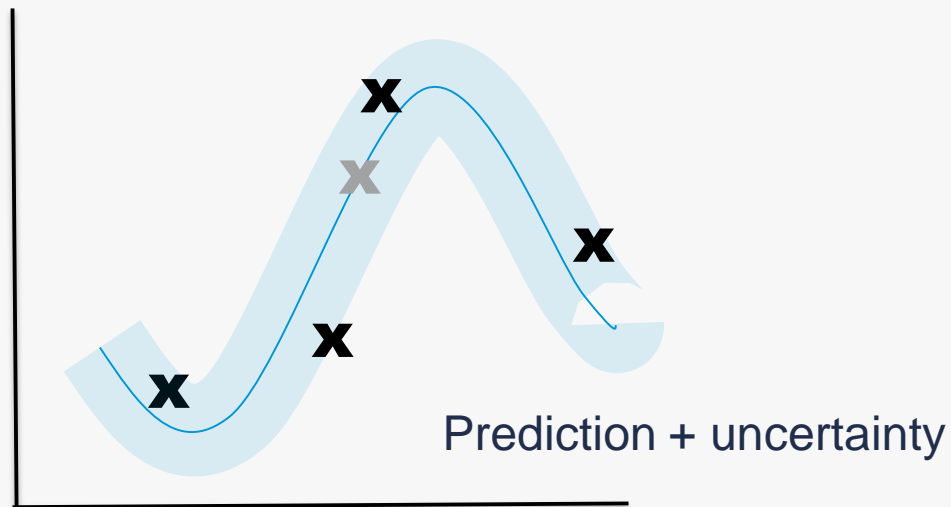
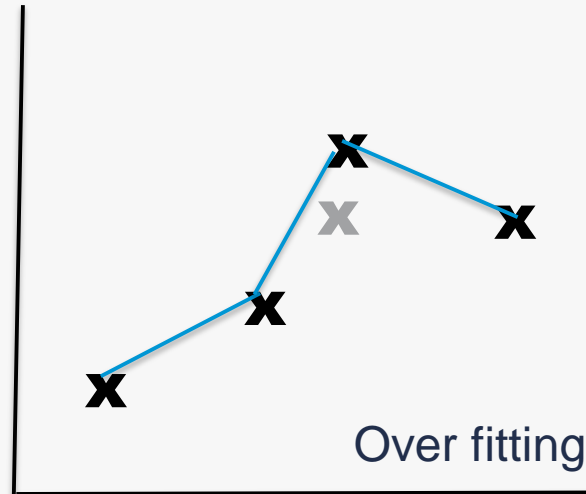
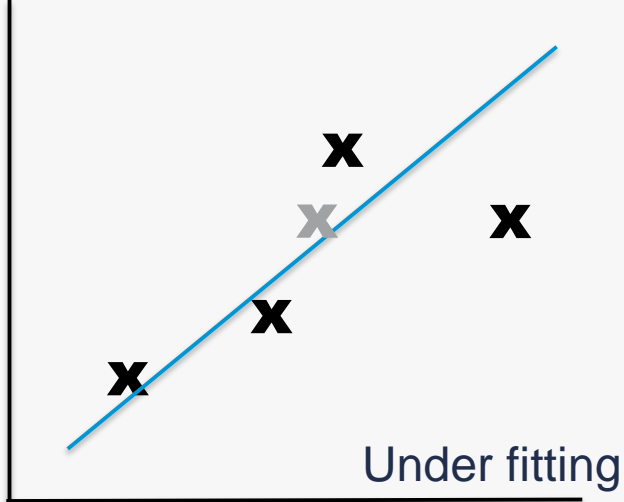




[pvinsights.com](http://pvinsights.com)

Computational Science Division





# Gaussian Process

## Quick introduction

Target function

$$y_i = f(x_i) + \epsilon_i$$


Noise function

$$\epsilon_i = \begin{bmatrix} \mathbf{y} \\ y_* \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K & K_*^T \\ K_* & K_{**} \end{bmatrix}\right)$$

Here\* Means a point  
that we want to predict

Given  $\mathbf{y}$ , the probability of  $y_*$  is:

$$y_* | \mathbf{y} \sim \mathcal{N}(K_* K^{-1} \mathbf{y}, K_{**} - K_* K^{-1} K_*^T)$$

Prediction (or kriging) 

$$\bar{y}_* = K_* K^{-1} \mathbf{y}$$

Variation 

$$\text{var}(y_*) = K_{**} - K_* K^{-1} K_*^T$$



# Gaussian process

## Covariance matrix

### Covariance function

$$k(x, x') = \sigma_f^2 \exp \left[ \frac{-(x - x')^2}{2l^2} \right]$$

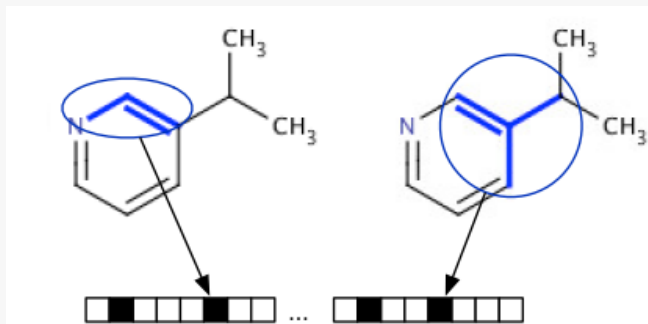
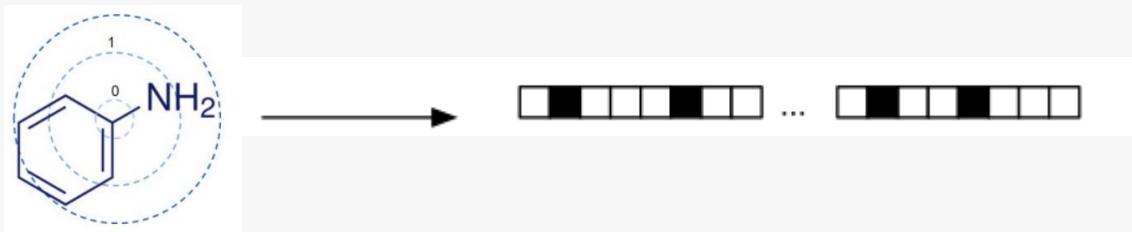
Example: Square exponential

### Covariance matrix

$$K = \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_n) \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_n) \\ \vdots & \vdots & \ddots & \vdots \\ k(x_n, x_1) & k(x_n, x_2) & \cdots & k(x_n, x_n) \end{bmatrix}$$

# Molecular Fingerprint examples

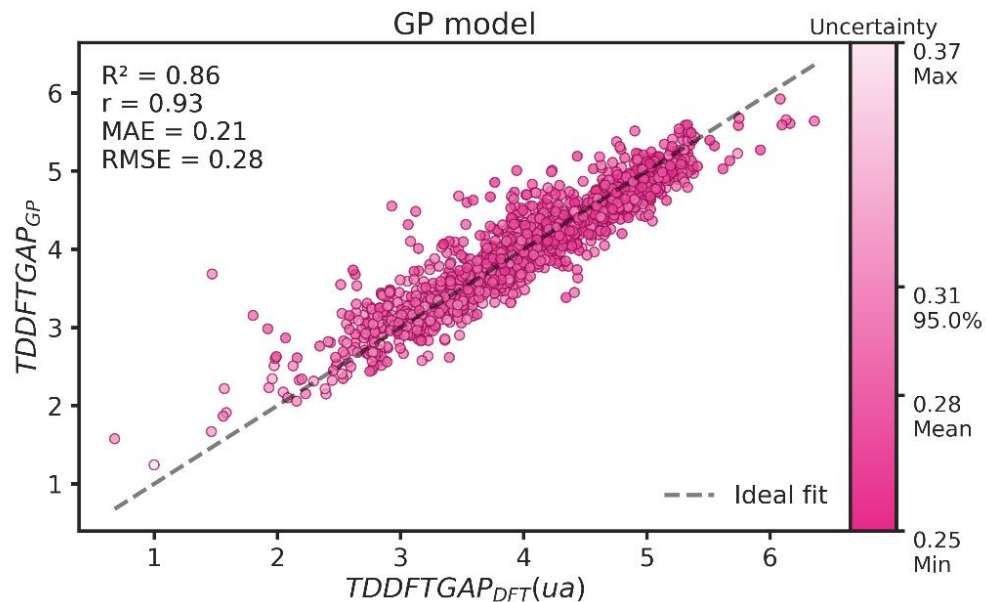
## Morgan Circular Fingerprint



# 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.

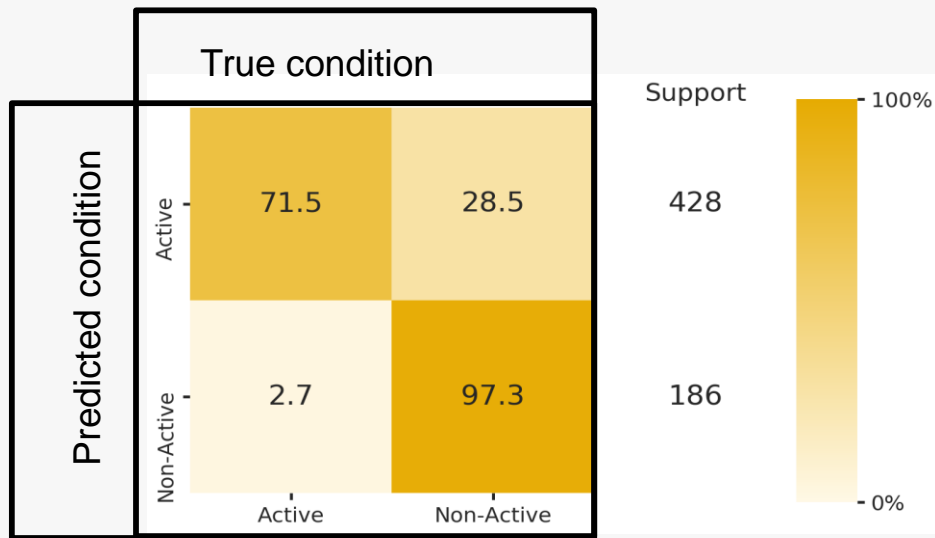


# Learn from data and feedback to experiments

## Is this optically active?

Oscillator Strength prediction –  
Transitions could not be optically  
active. We can predict which of  
electronic transitions have an  
oscillator strength  $< 0.8$  a.u. with an  
error 3%.

### Confusion matrix



# Future work

- **Extinction coefficients prediction could be improved adding extra information that could be slightly costly to get, such as orbital dipole moments or results from lower scale methods.**
- **Variational autoencoders (VAE) could help us to discover the most important molecular features of good dyes in the dataset. This is work in progress.**
- **Using Generative models to produce new molecules that optimize the physical chemical properties we want, and that are likely to exist (chemically stable) and can be synthesized in the lab.**
- **We will release a suit of code to simplify data driven materials research, with building blocks to tailor workflows for similar problems.**



# Q&A

[pvinsights.com](http://pvinsights.com)

Computational Science Division