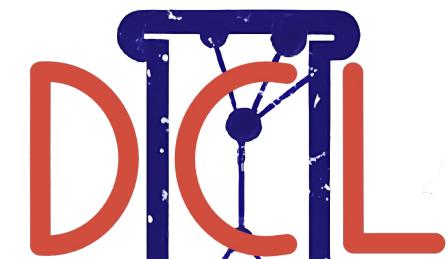


Simulations and machine learning for molecular design and reactivity

Kjell Jorner

Assistant Professor of Digital Chemistry @ ETH Zurich

Chalmers AI4Science Seminar, 2024-06-13



The Digital Chemistry Laboratory

Department of Chemistry and Applied Biosciences (D-CHAB)
Institute of Chemical and Bioengineering (ICB)



Principal investigator
Prof. Dr. Kjell Jorner

PhD students

Lauriane Jacot-Descombes
Giustino Sulpizio
Stefan Schmid
Luca Schaufelberger
Franziska Weißbach



Co-supervisor

Vignesh Ram Somnath (w. Krause)
Riccardo De Santi (w. Krause, He)



@DCL_ETHZ

Open position in Computer-aided Catalyst Design



- Closes **15.06** (on Saturday)
- Ideally computational science (chemistry) & machine learning profile

PhD position in computer-aided catalyst design

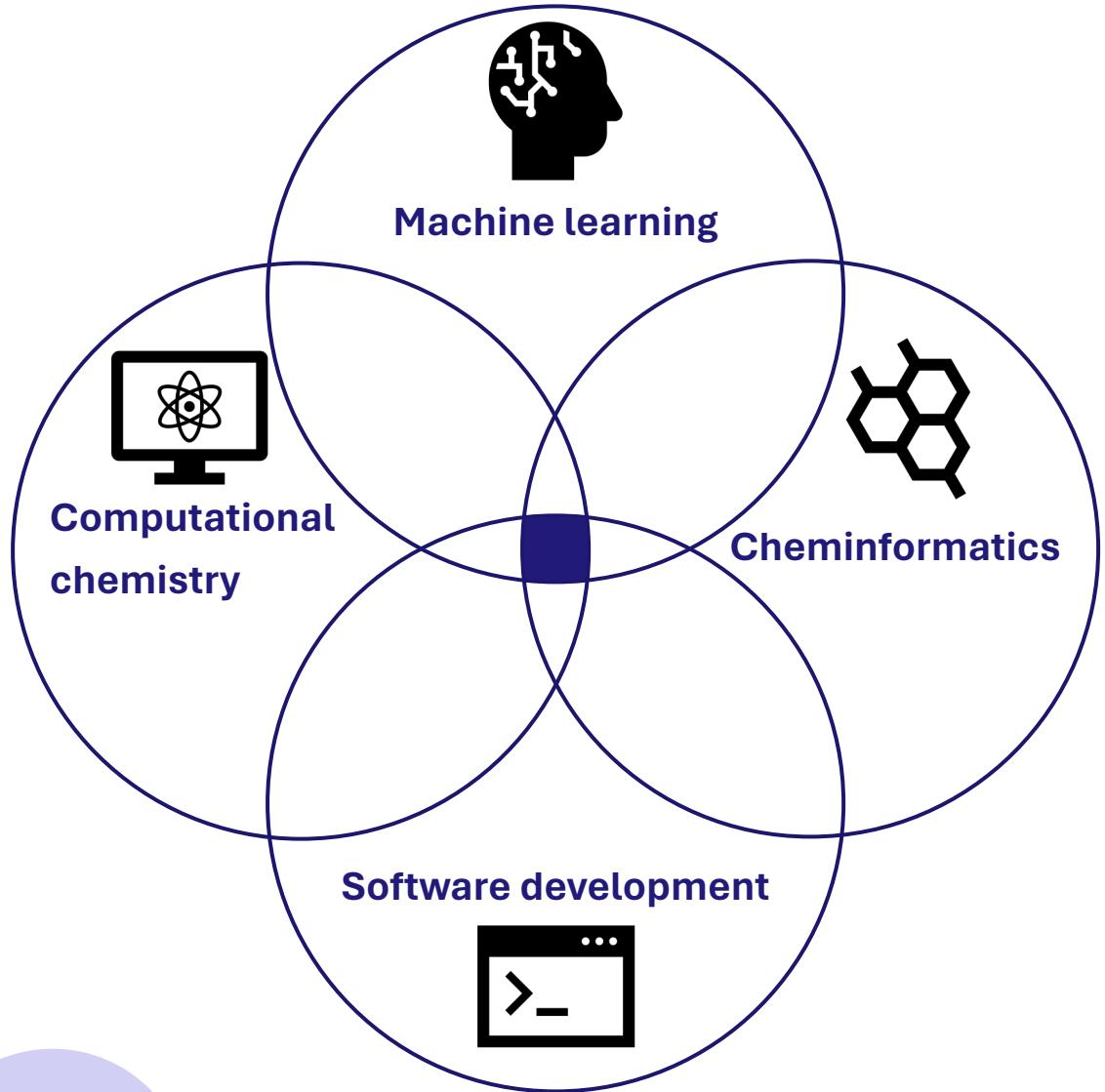
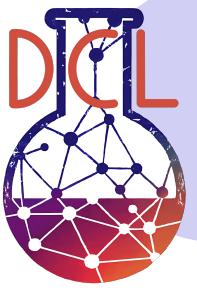
100%, Zurich, fixed-term

Print

The Digital Chemistry Laboratory is led by Prof. Dr. Kjell Jorner at the Institute of Chemical and Bioengineering, within the Department of Chemistry and Applied Biosciences at ETH Zurich. Our mission is to accelerate chemical discovery using digital tools. We predict chemical reactivity and molecular properties using the tools of machine learning, computational chemistry, and cheminformatics. Our ultimate goal is the computer-aided design of molecules and especially catalysts.



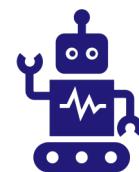
Digital Chemistry



Lab chemist

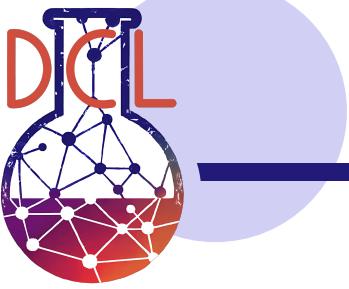


High-throughput
experimentation



Self-driving labs

Mission: Accelerate chemical discovery with digital tools

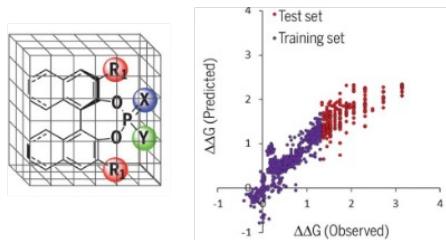


Reactivity and catalysis

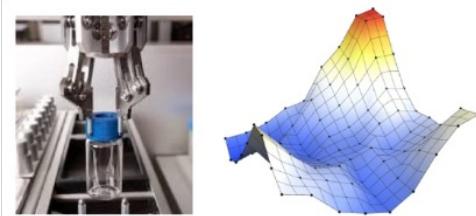


Strategic core research

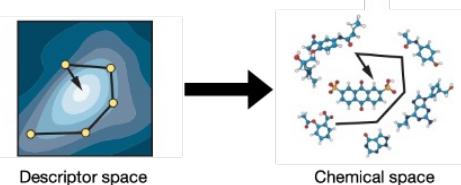
Reaction outcome prediction



Reaction optimization



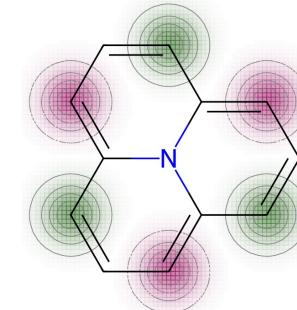
Catalyst & reaction design



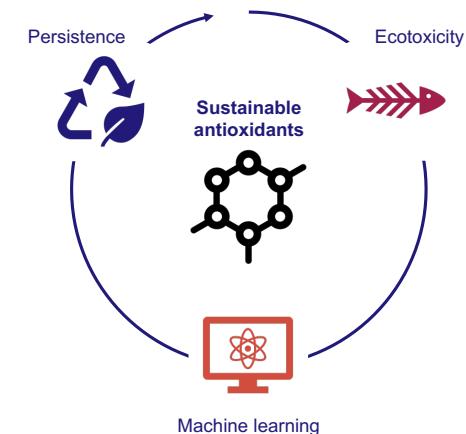
Opportunity-driven

Computer-aided molecular design

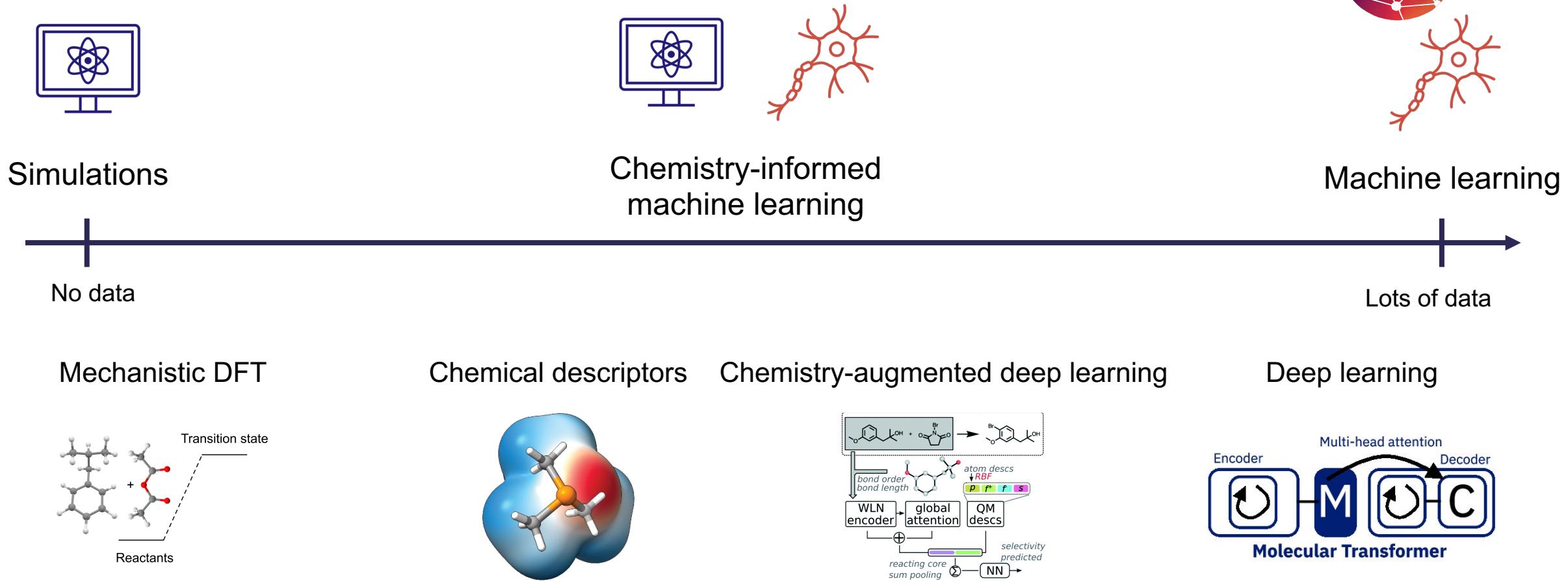
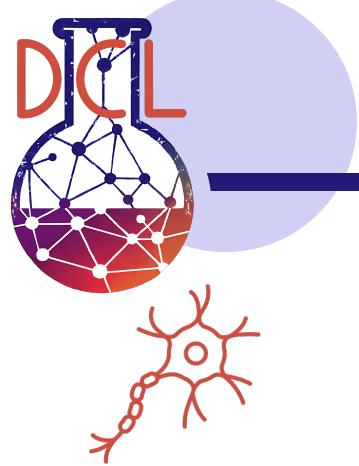
Organic electronic materials



Safe and sustainable by design



Our niche: Models in the low-data regime

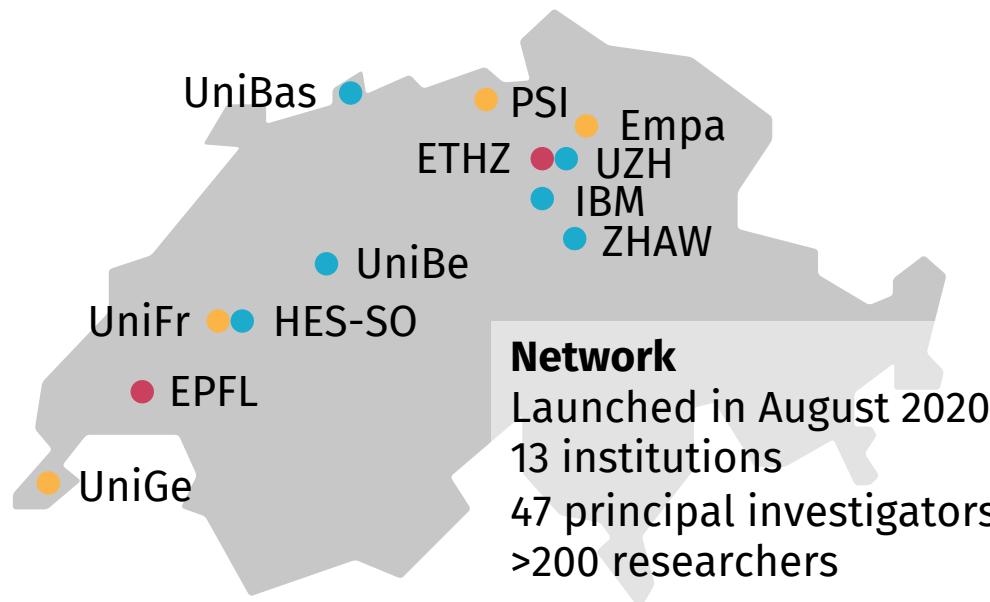




NCCR Catalysis

Research themes

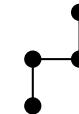
Small molecule activation
Fine chemicals
from renewable platforms
Advanced tools
Digitalization
Implementation



Network

Launched in August 2020
13 institutions
47 principal investigators
>200 researchers

Total investment
32 MCHF, 4 years



**Swiss National
Science Foundation**

The National Centres of Competence in Research (NCCRs) are a funding scheme of the Swiss National Science Foundation



**Prof. Javier
Pérez-Ramírez**
Director
ETH zürich



**Prof. Jérôme
Waser**
Co-Director
EPFL



nccr-catalysis.ch



NCCR Catalysis



@NCCR_Catalysis



@nccr_catalysis



ETH AI CENTER

114

Faculty members

41

Fellows

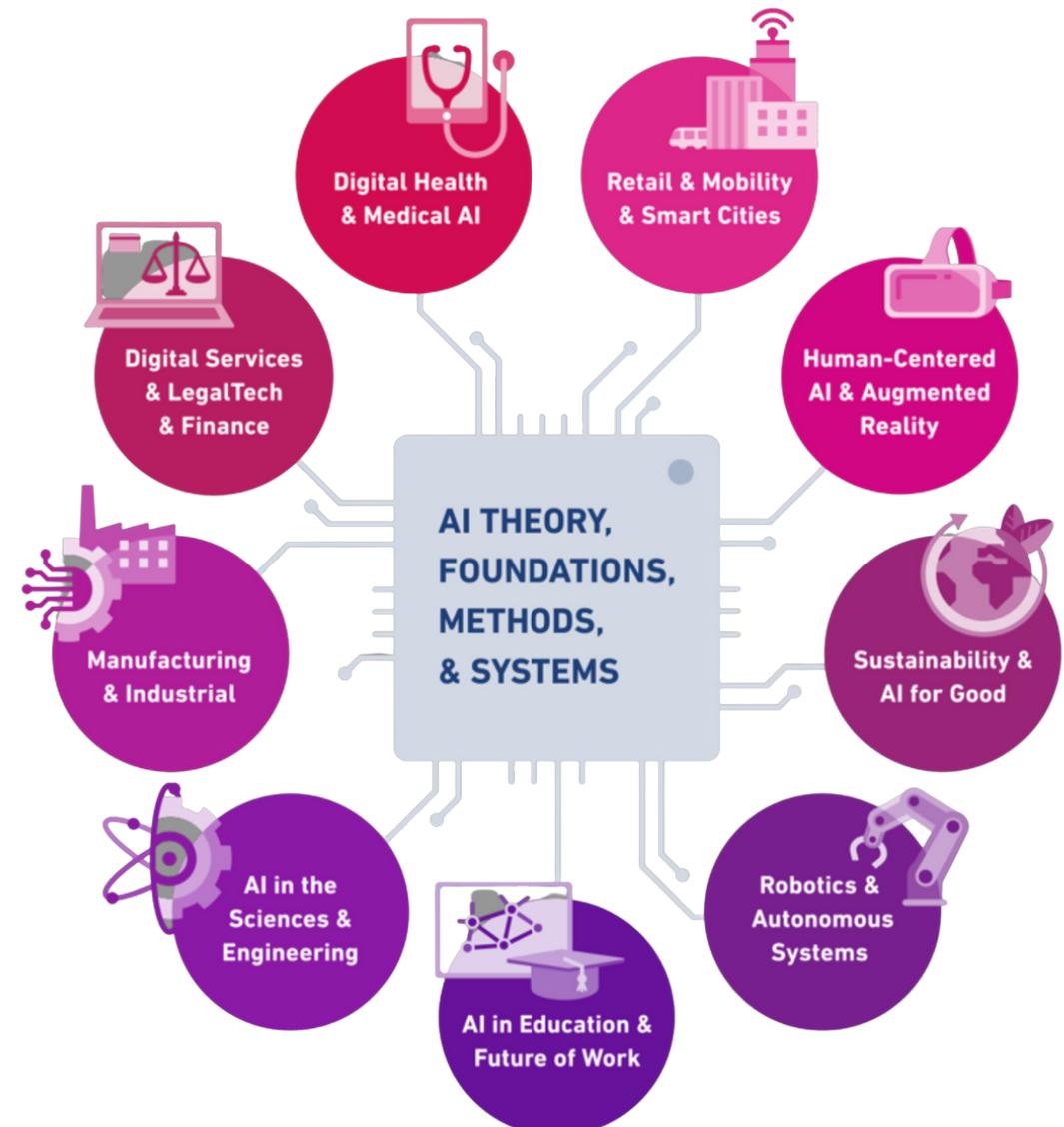


PD Dr. Alexander Illic
Executive Director



Prof. Andreas Krause
Chair of Steering Committee
and Core Faculty Member

As **ETH's central hub for artificial intelligence**, we bring together researchers of AI foundations, applications, and implications across all departments. We foster research excellence, industry innovation, and AI entrepreneurship to **promote trustworthy, accessible, and inclusive AI systems**.



ai.ethz.ch



eth_ai_center

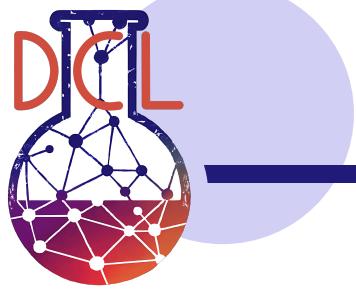


@ETH_AI_Center

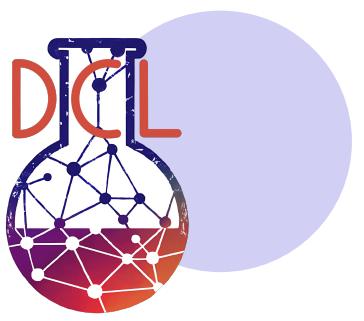
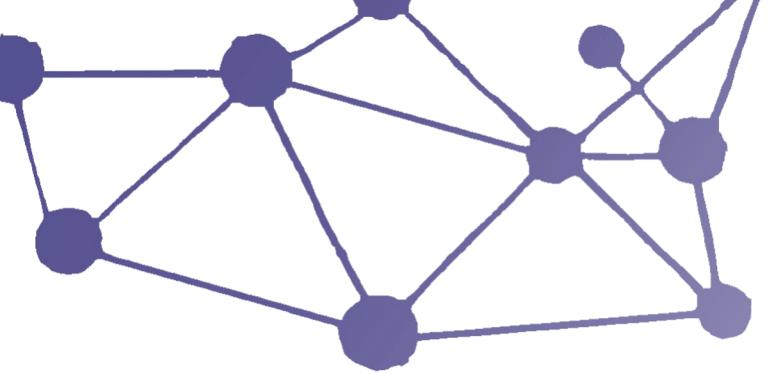


ETH AI Center

Outline

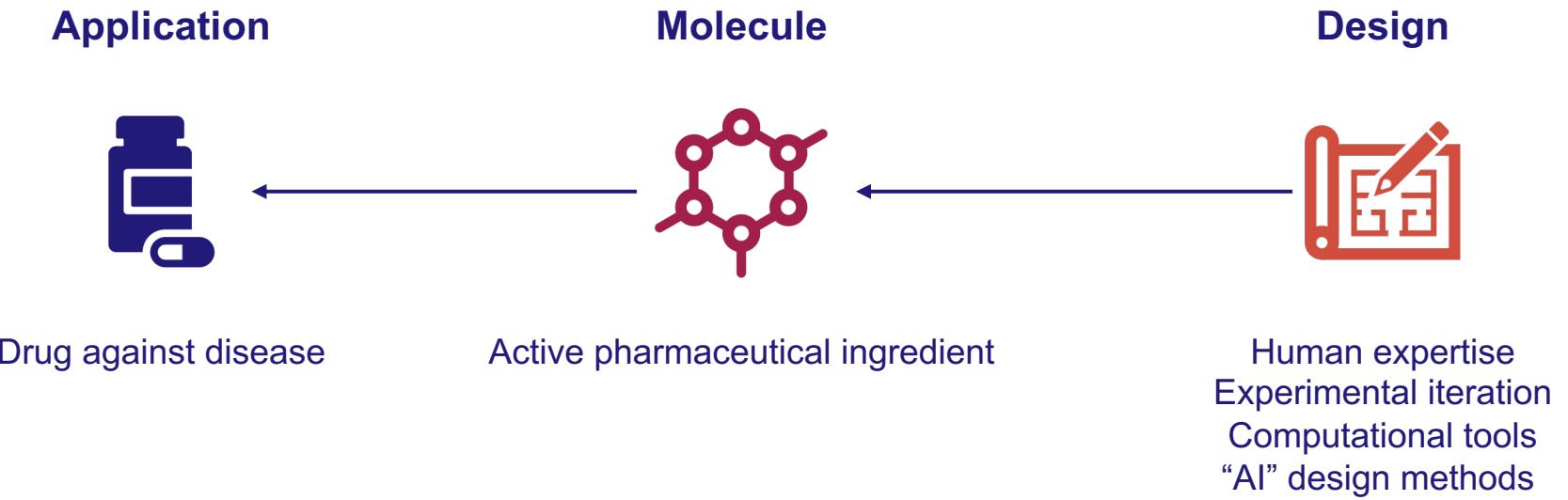
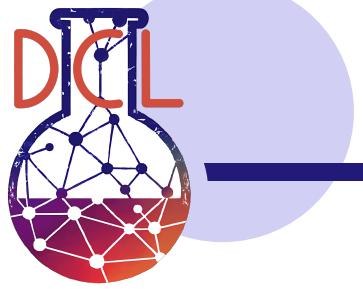


- Computer-aided molecular design
- Molecular design with Hückel theory
- Designing singlet fission materials with uncertainty-aware genetic algorithms
- Reaction design with reactive force field methods



Computer-aided molecular design

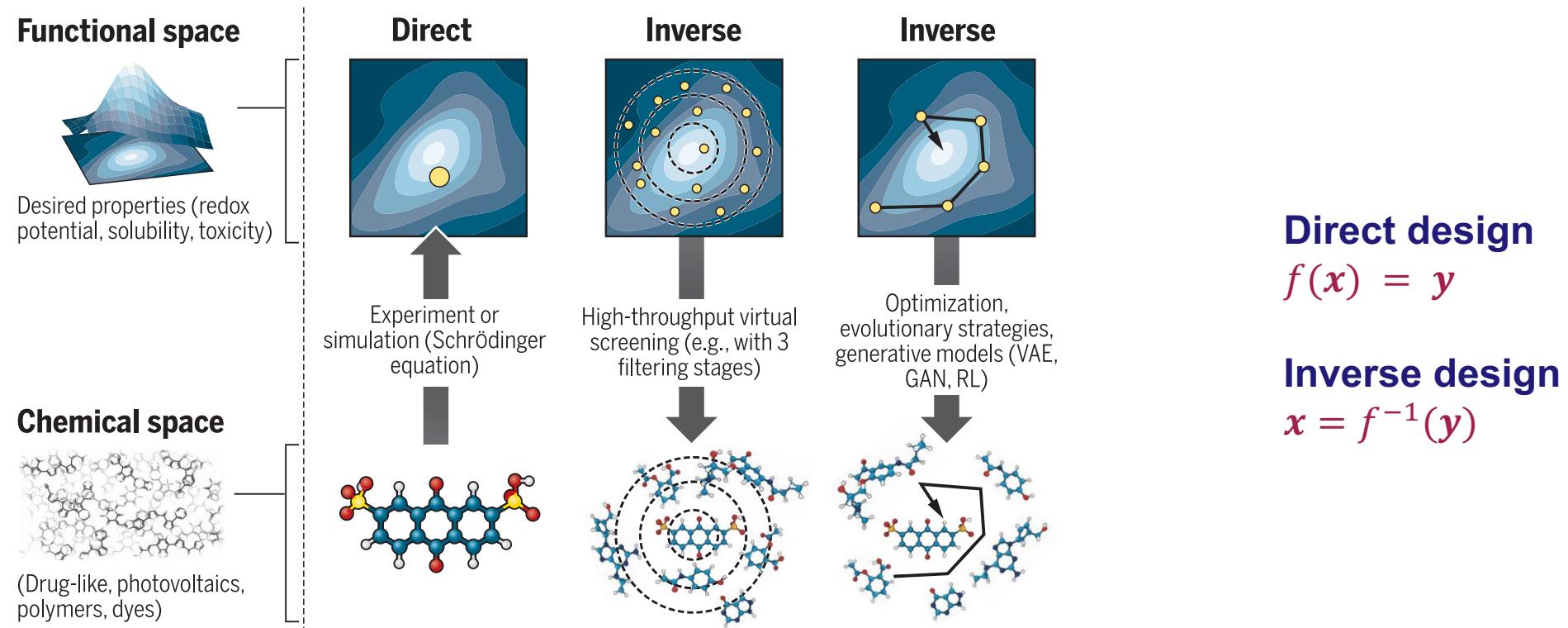
Principle of computer-aided molecular design



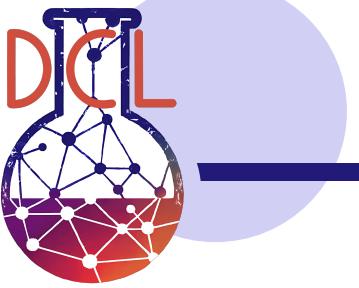
Inverse design



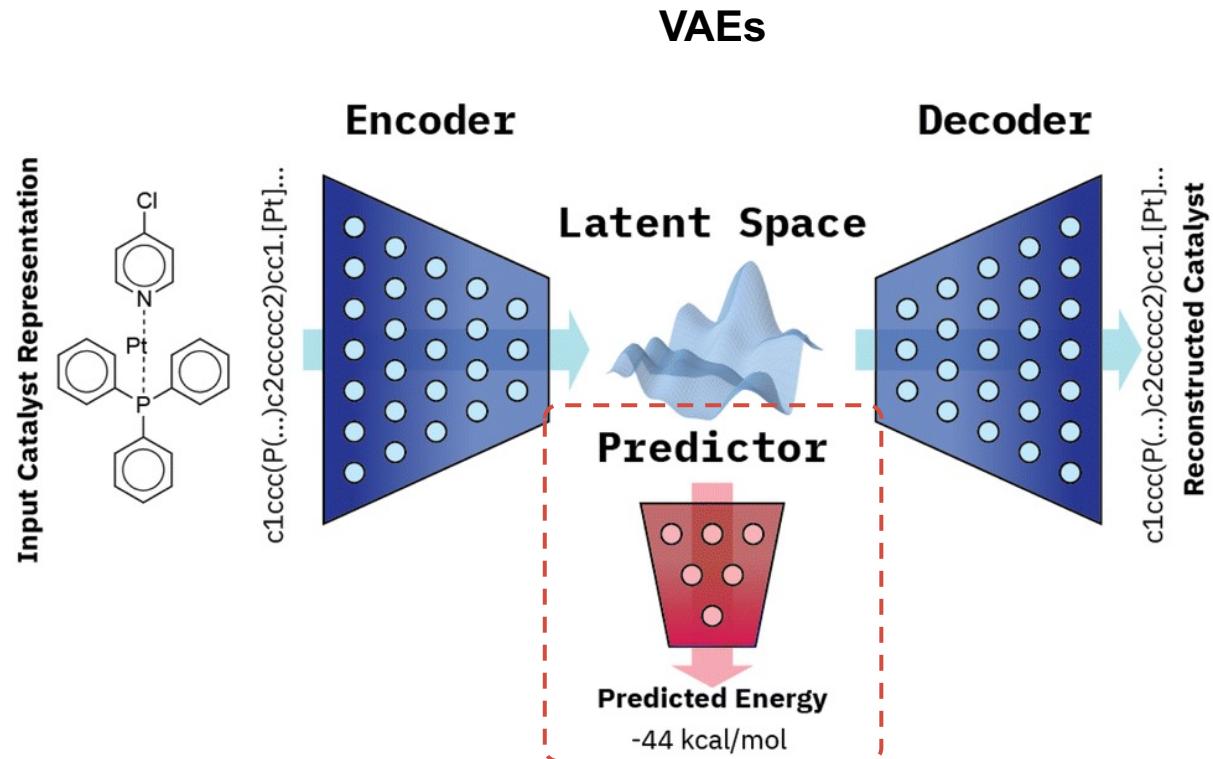
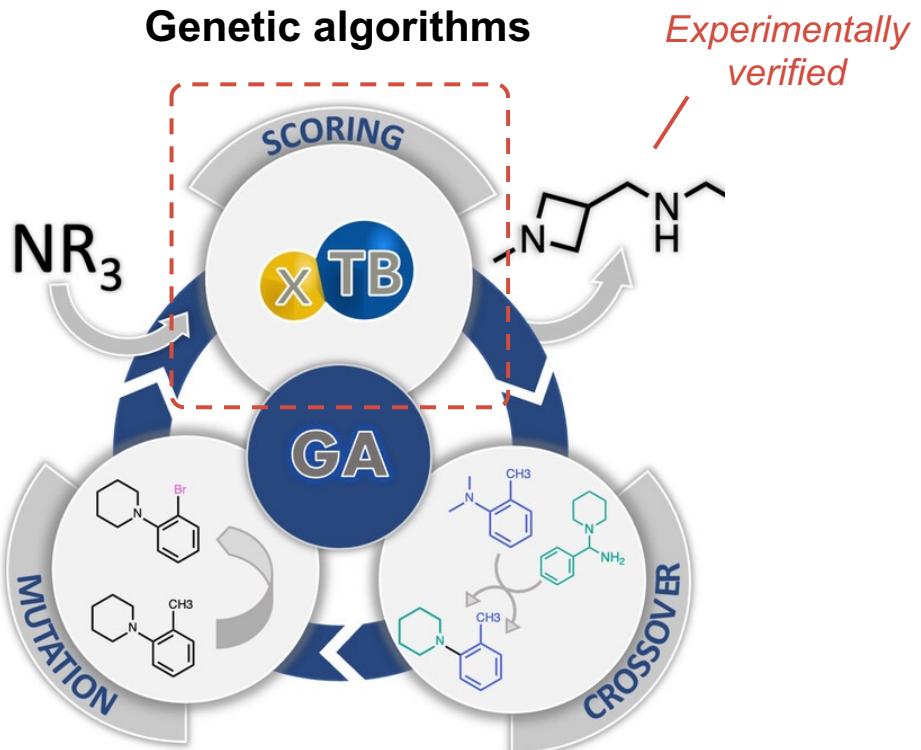
- Promises more effective searching of very large chemical design spaces (10^{60})
- Complementary to human design and older techniques such as virtual screening



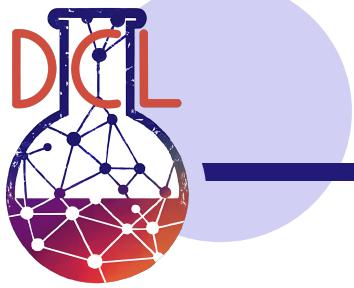
Inverse design methods



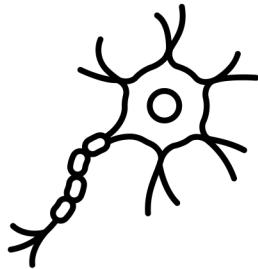
- Models that suggest new molecules with targeted properties
- Avoid exhaustive enumeration of chemical/materials space ($>10^{60}$)
- Rely on property predictors to assess molecular function



Property predictors

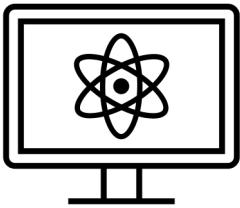


Machine learning



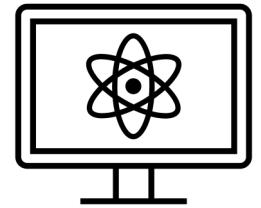
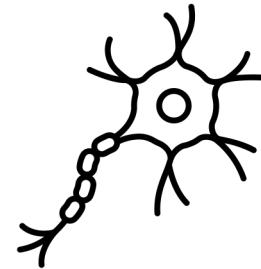
- ⊕ Fast
- ⊕ Accurate for “interpolation”
- ⊖ Limited applicability domain

Physics-based simulation

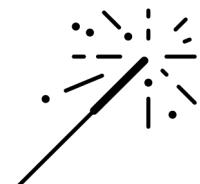


- ⊕ Wide applicability domain
- ⊖ Slow
- ⊖ Inaccurate for some problems

Hybrid models

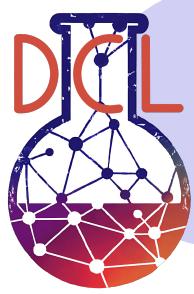


- ⊕ Fast
- ⊕ Wide applicability domain
- ⊖ Not generally available yet



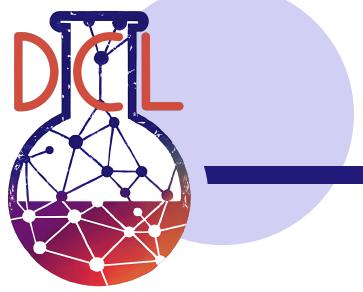
“Hybrid light”: Semi-empirical methods

Goodhart's law

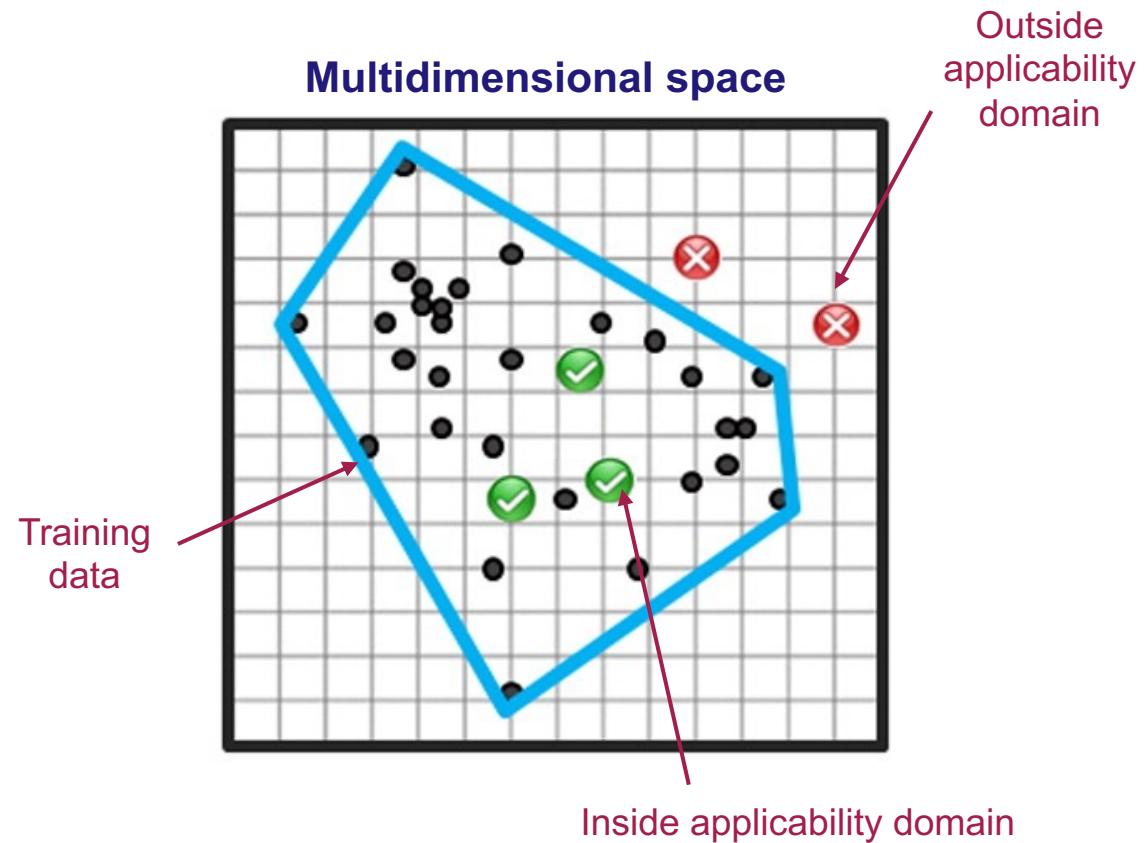
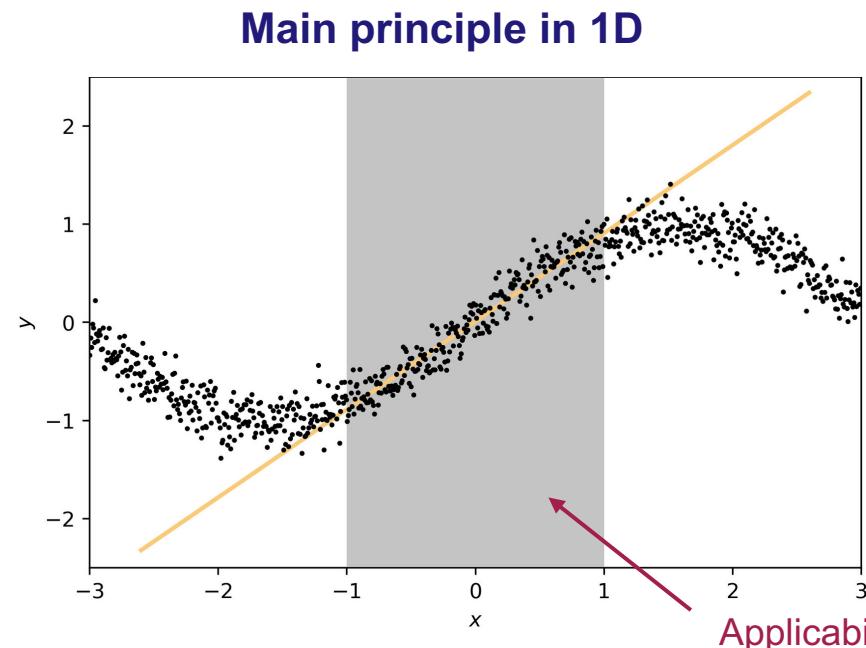


“When a measure becomes a target, it ceases to be a good measure.”

Applicability domain



- Abstract concept of when a (chemical) ML model can be used
- Grows with more data
- Physical models have wider domain
- Machine learning models have narrower domain

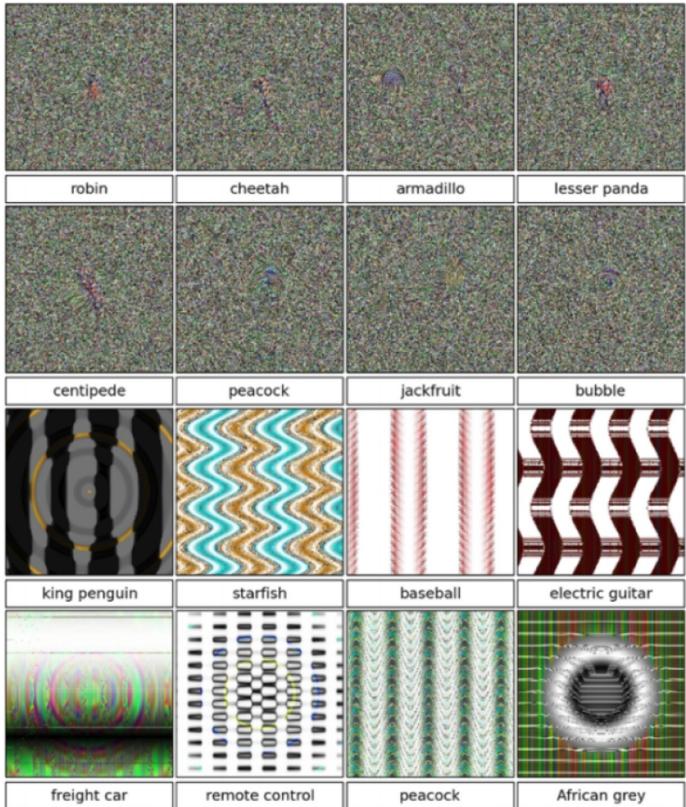


“Adversarial examples”

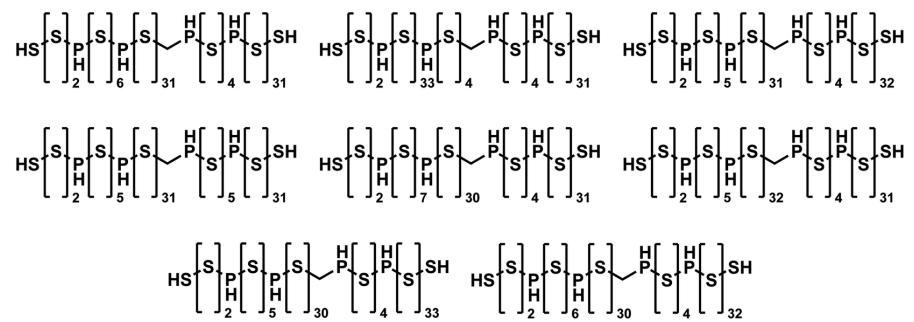


- Optimize the input to score well on classifier without connection to underlying task

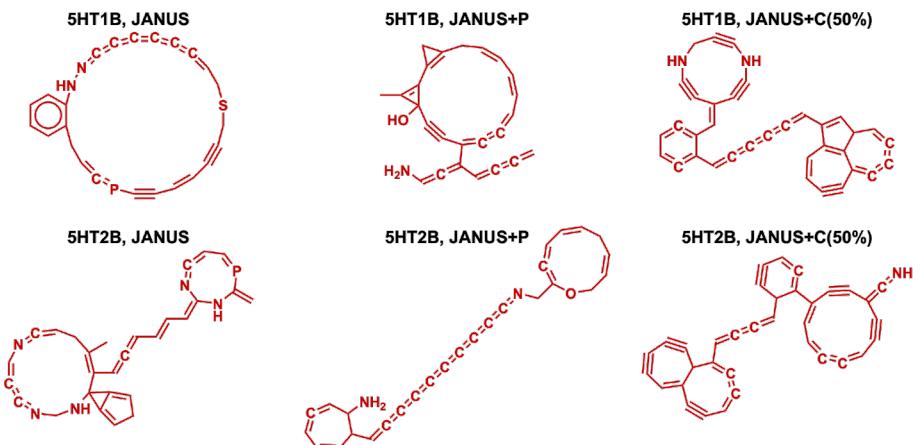
Image classification



Penalized logP



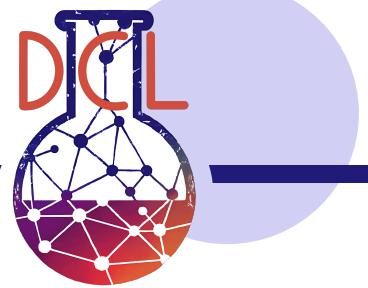
Docking



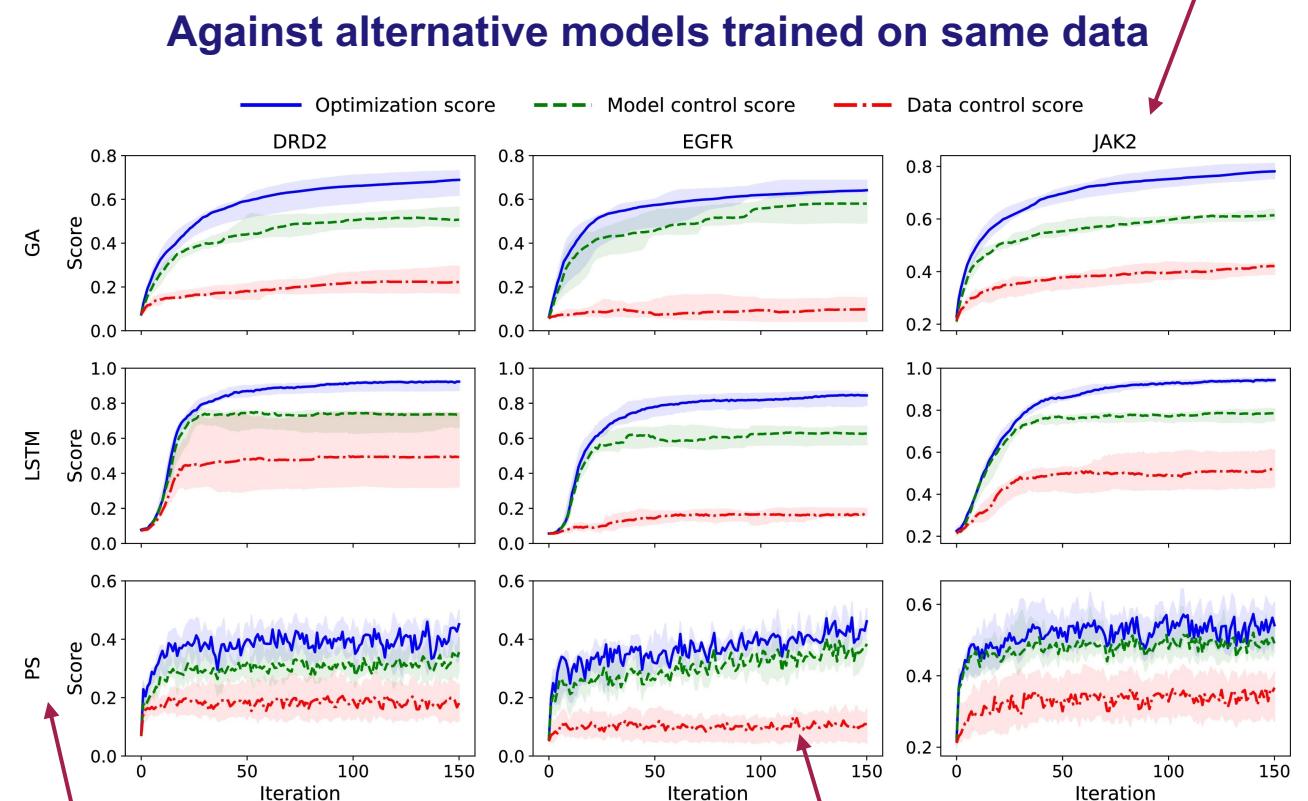
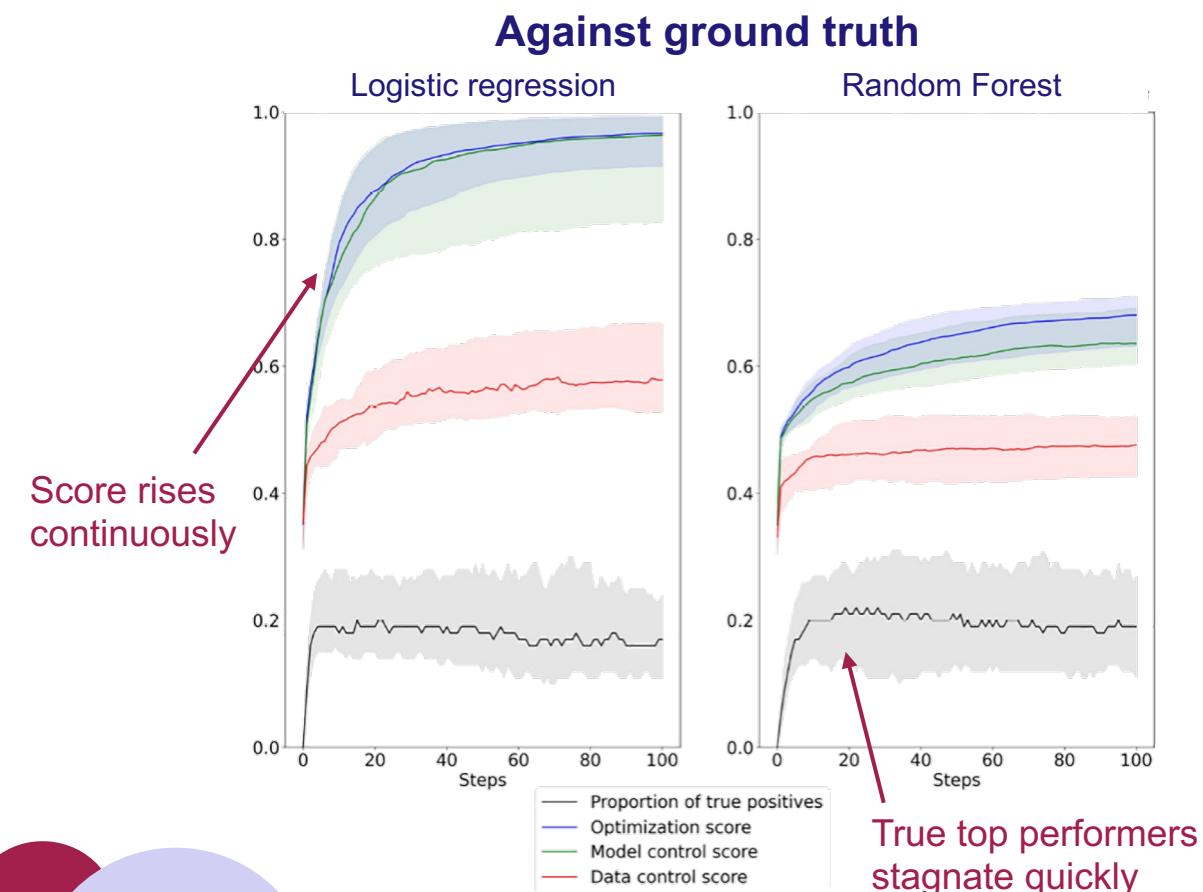
Molecular

(a) Gendreau, P.; Turk, J.-A.; Drizard, N.; Ribeiro Da Silva, V. B.; Descamps, C.; Gaston-Mathé, Y. *J. Chem. Inf. Model.* **2023**. 10.1021/acs.jcim.3c00195 (b) Nigam, A.; Pollice, R.; Aspuru-Guzik, A. *Digit. Discov.* **2022**. 10.1039/D2DD00003B

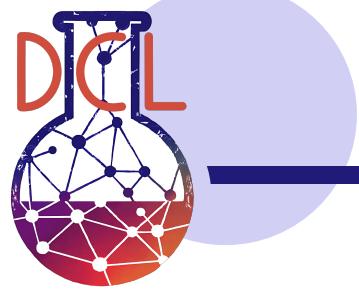
Hacking the property predictor



- Generative models quickly exploit deficiencies in property predictors

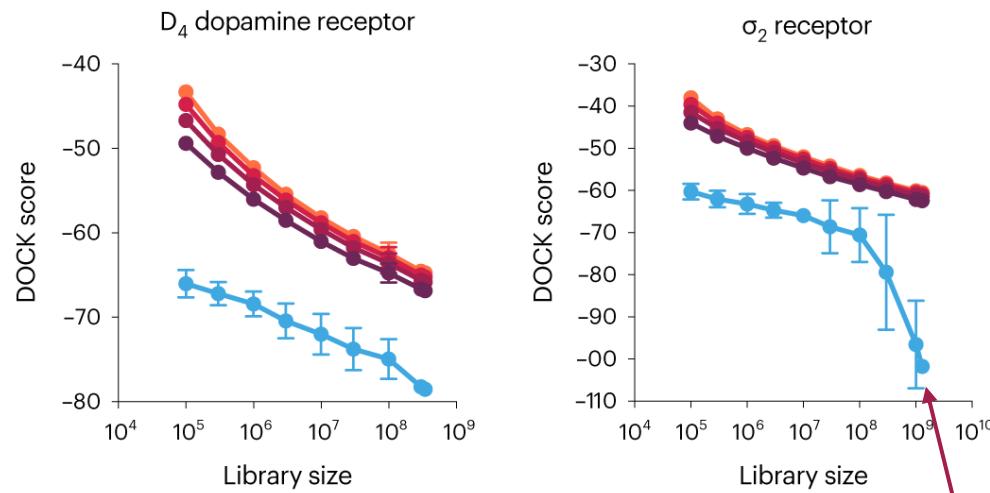


Problem worsens in larger chemical spaces

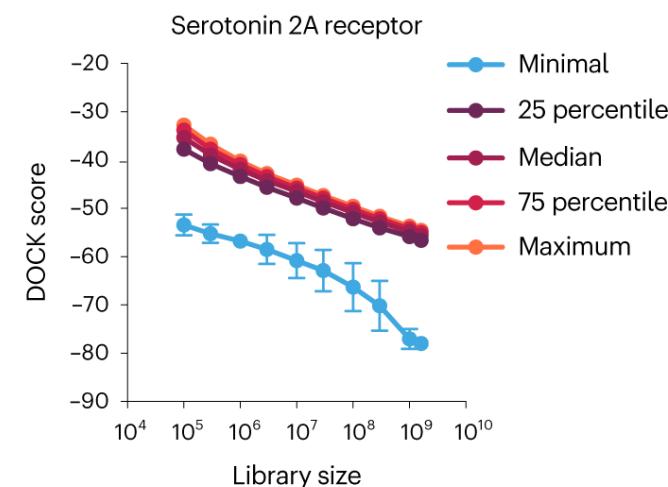


- “Artifacts” are particularly offensive molecules that exploit holes in simple predictors
- At larger chemical library sizes, they tend to dominate the top performers
- Goal-directed generative models easily latch onto these artifacts

Docking score as function of library size

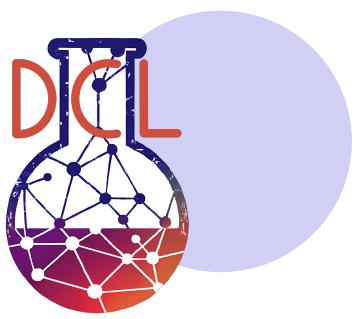
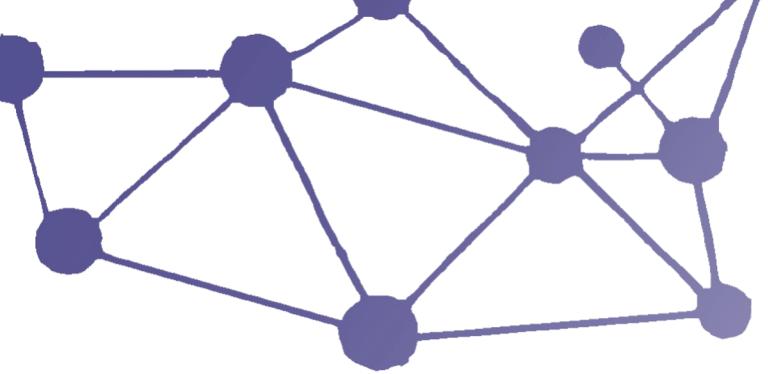


Artifacts dominate best-scoring molecules at larger sizes



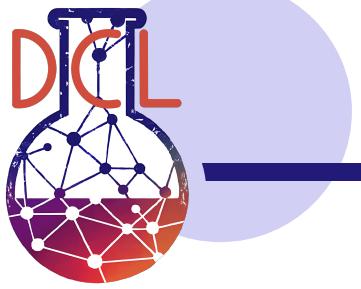
Percentage of artifacts

Library size	Top 100		Top 1K	
	10^{-6}	10^{-5}	10^{-6}	10^{-5}
1B	100	100	42.4	99.7
	42	100	5.9	56.2
	9	57	1.0	7.8
	0	9	0.1	0.9
100M				
10M				
1M				
Top 10K				
Top 100K				

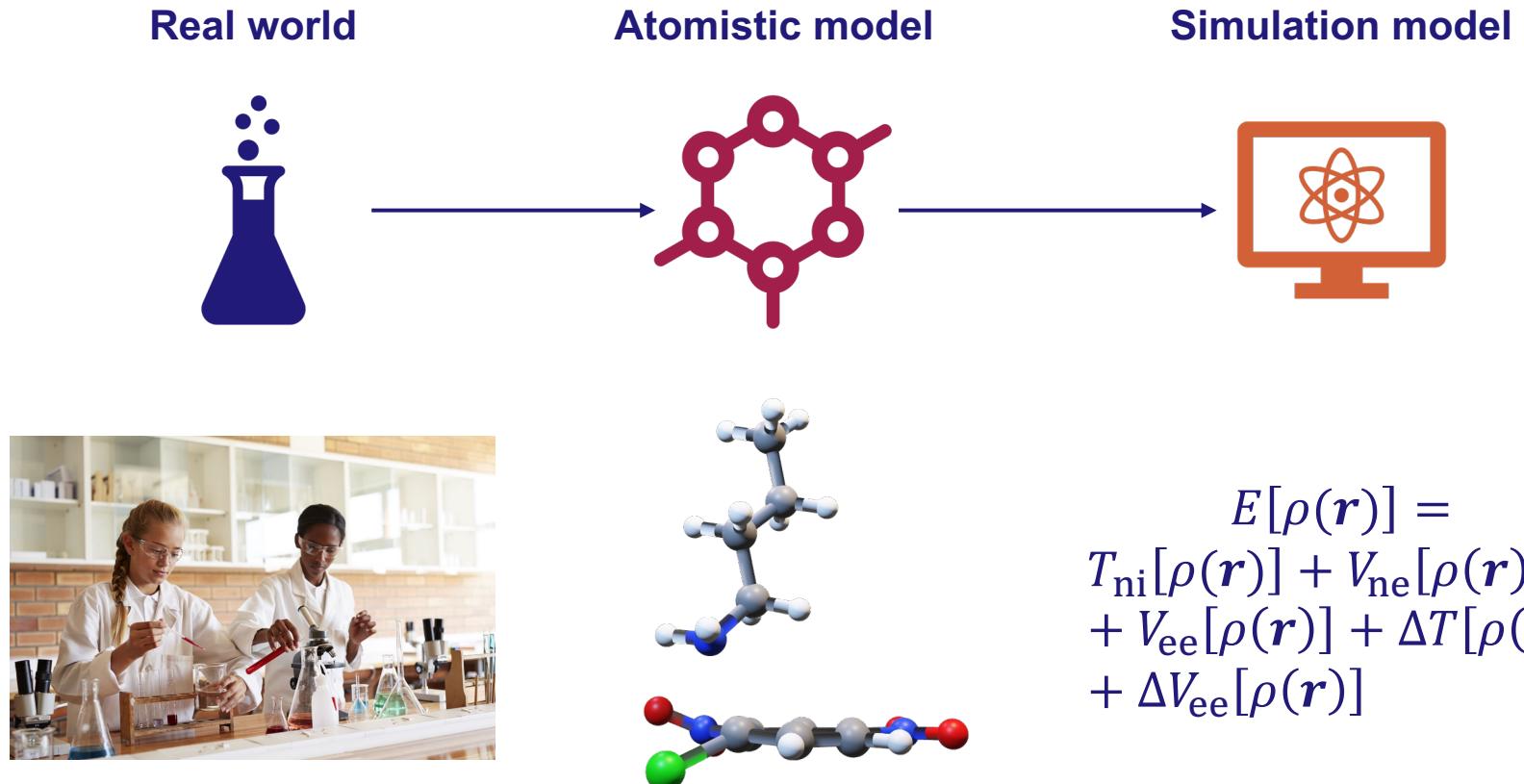


Molecular design with Hückel theory

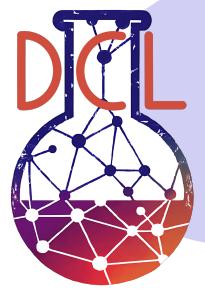
Computational modeling of chemistry



- Translate experiment into mathematical model of atomistic model

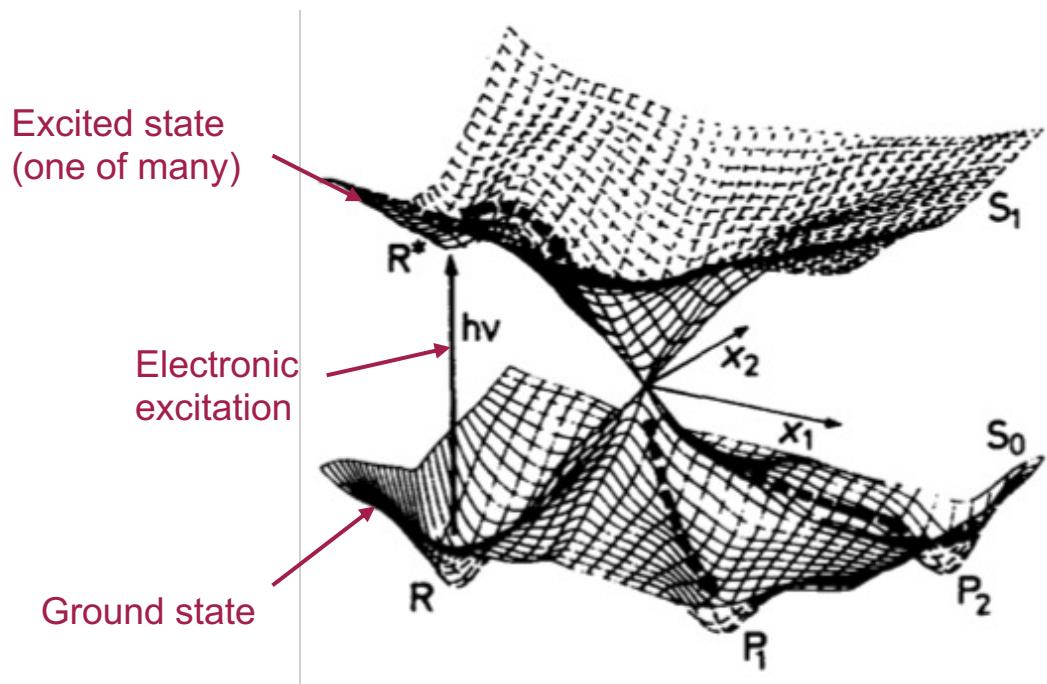


The potential energy surface (PES)

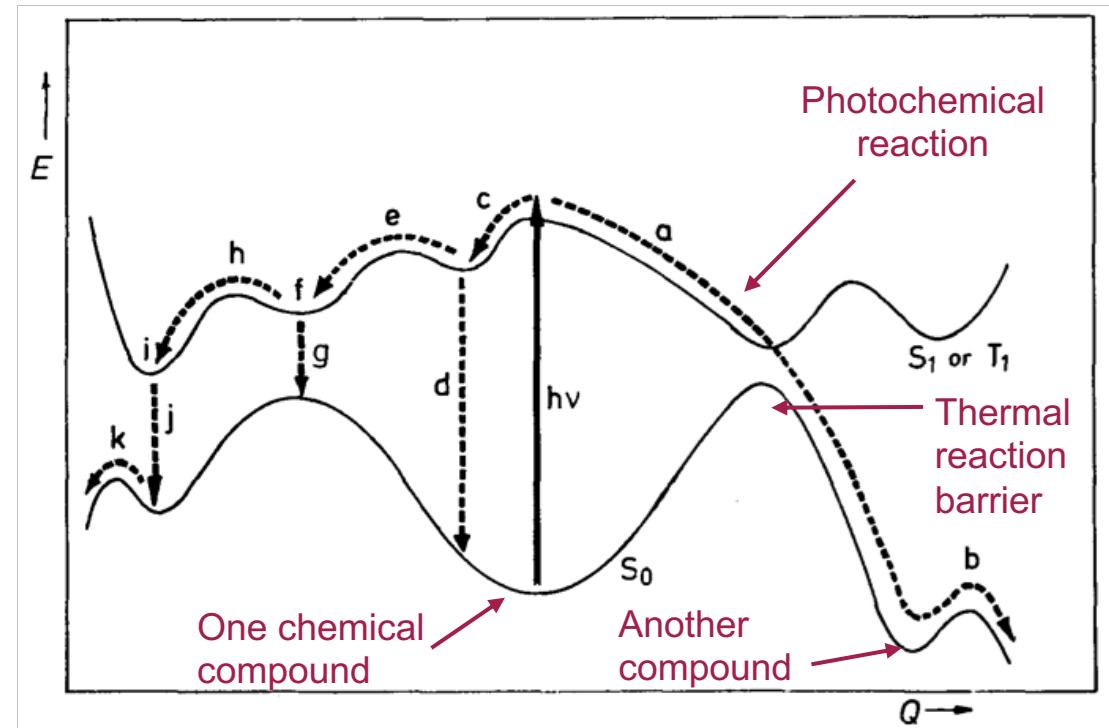


- Molecules live in $3N$ -dimensional space (N number of atoms)
- The energy as a function of the coordinates is called the PES

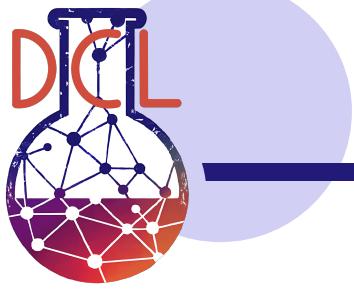
Multi-dimensional space



Reduced-coordinate “slice”

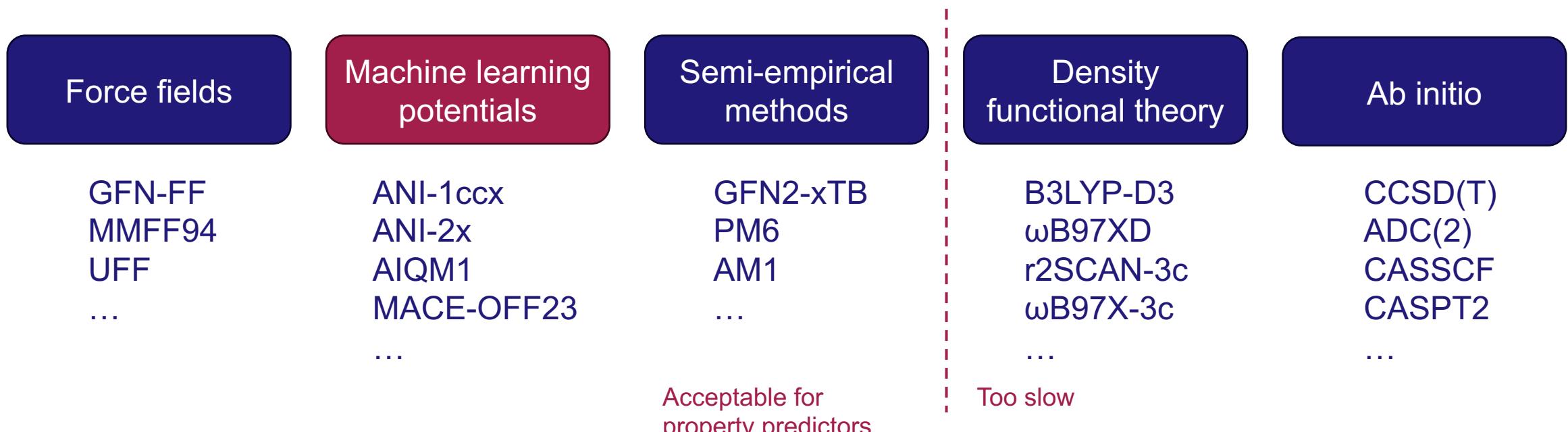


Computational methods

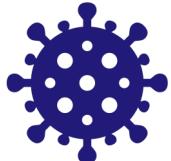
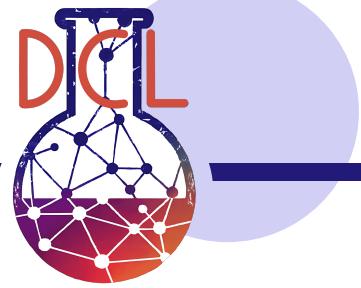


- Calculate the energy as function of coordinates
- A sequence of increasingly time-consuming physical methods

Increasing computational cost



What can we simulate?



Interactions with surroundings

- Drug binding energies
- Solubility
- Redox potentials
- Thermodynamic properties of solutions
- ...



Structure & stability

- Crystal structure polymorphism
- Protein structure
- Molecular conformers
- Tautomer/protonation state
- ...



Chemical reactivity

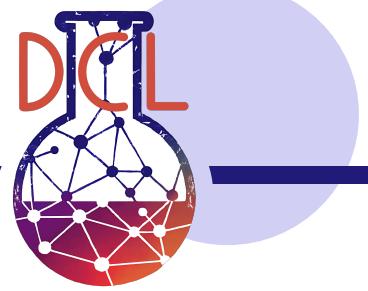
- Rate constants
- Selectivity
- Bond dissociation energies
- Kinetic isotope effects
- Photochemistry
- ...



Spectroscopic properties

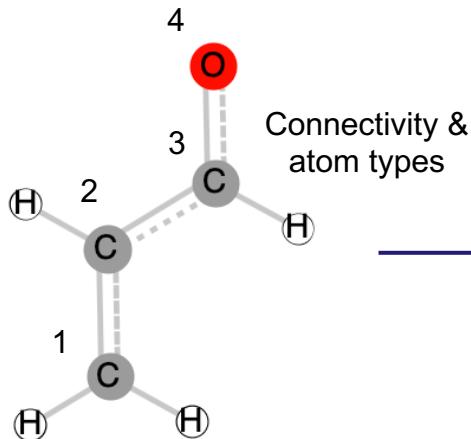
- UV/vis
- IR
- Rotational
- X-ray
- Ionization potentials and electron affinities
- NMR & EPR
- Optical rotation & circular dichroism
- Mass spectra
- ...

Hückel theory – extremely semi-empirical



- One of the first molecular orbital theories
- Schrödinger equation reduces to a very simple form

Molecule



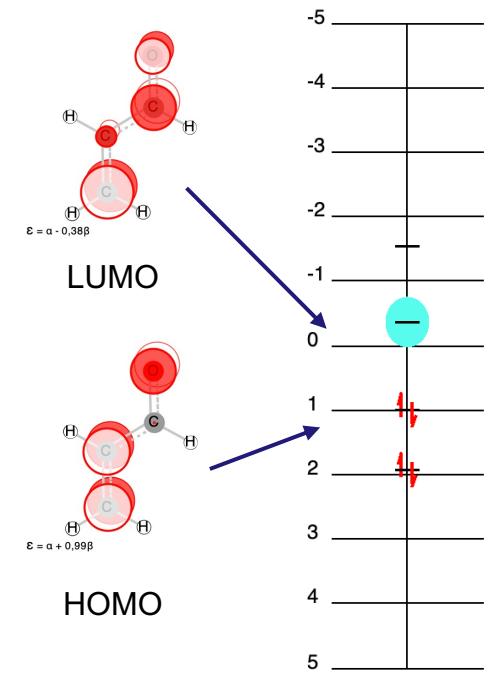
Hückel matrix

$$\begin{bmatrix} h_C & k_{CC} & 0 & 0 \\ k_{CC} & h_C & k_{CC} & 0 \\ 0 & k_{CC} & h_C & k_{CO} \\ 0 & 0 & k_{CO} & h_O \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1.06 \\ 0 & 0 & 1.06 & 0.97 \end{bmatrix}$$

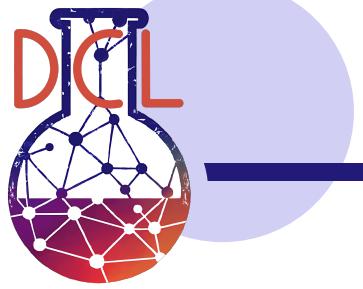
One-electron integrals

Diagonalize

Orbitals & energies

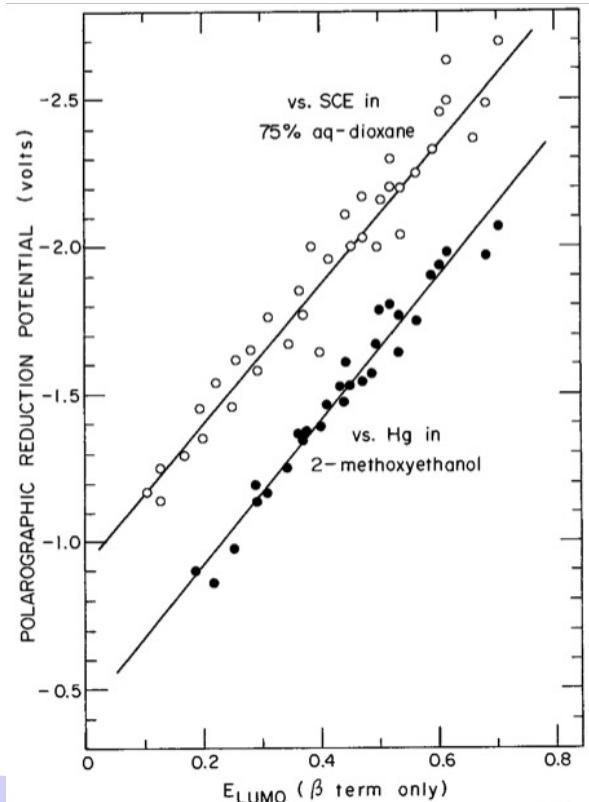


Indication of performance of Hückel theory

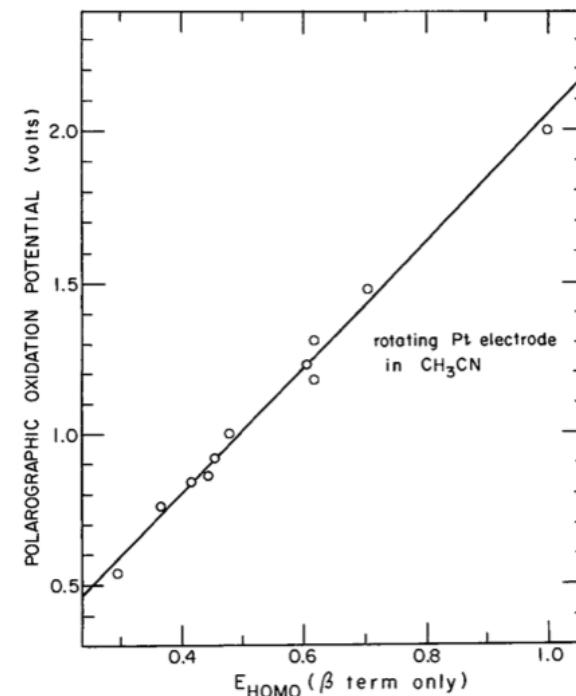


- Gives favorable correlations with experiment within certain compound classes

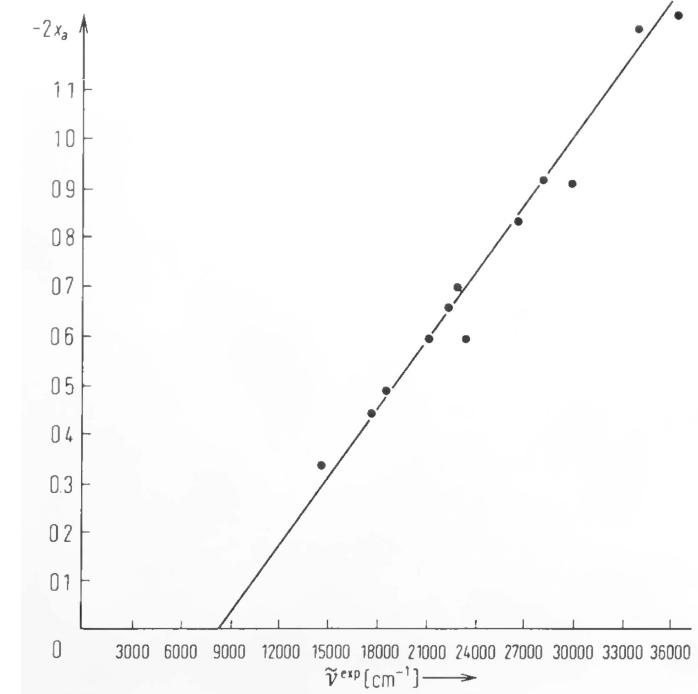
Reduction potential vs. LUMO



Oxidation potential vs. HOMO



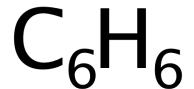
Photoexcitation vs. HOMO-LUMO gap



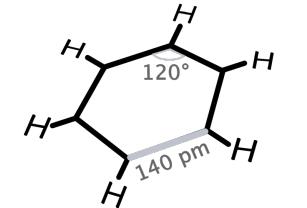
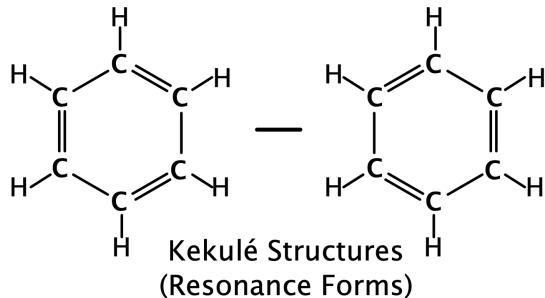
Aromaticity



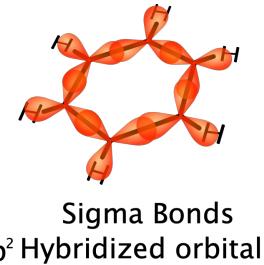
- Extra stabilization due to cyclic delocalization of π -electrons
- Applications to organic electronic materials, reactivity etc.



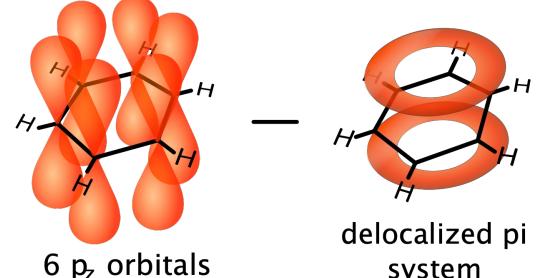
Benzene
Molecular formula



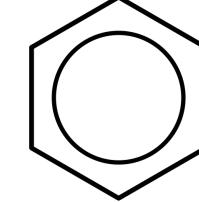
Planar Hexagon
Bond Length 140 pm



Sigma Bonds
 sp^2 Hybridized orbitals

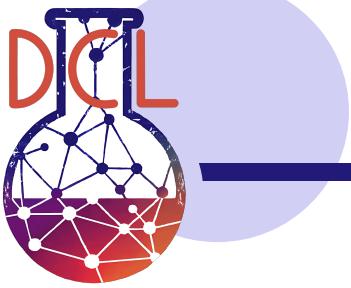


delocalized pi
system

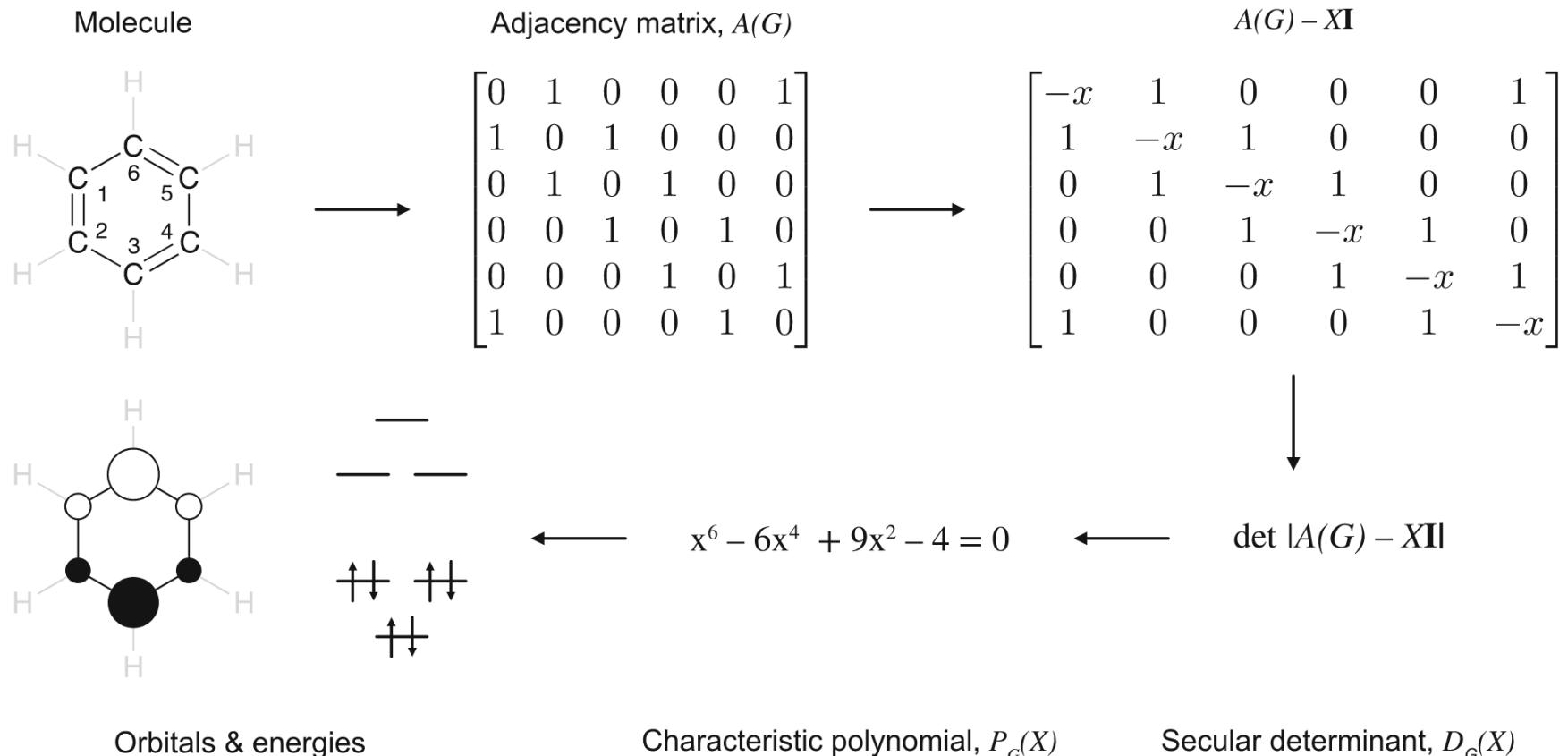


Benzene ring
Simplified depiction

The characteristic polynomial



- Solving Hückel theory in the case of only carbon atoms



Graph theory of aromaticity



- Invented independently by Aihara in Japan and Trinajstić & Gutman in Croatia
- Builds on Coulson's graph treatment of Hückel molecular orbital theory
- Calculates aromatic properties as difference between
 - E_π : Energy of actual π system
 - E_π^{ref} : Energy of “olefinic” reference system (lacking cyclic conjugation)



Jun-ichi Aihara



Nenad Trinajstić



Ivan Gutman

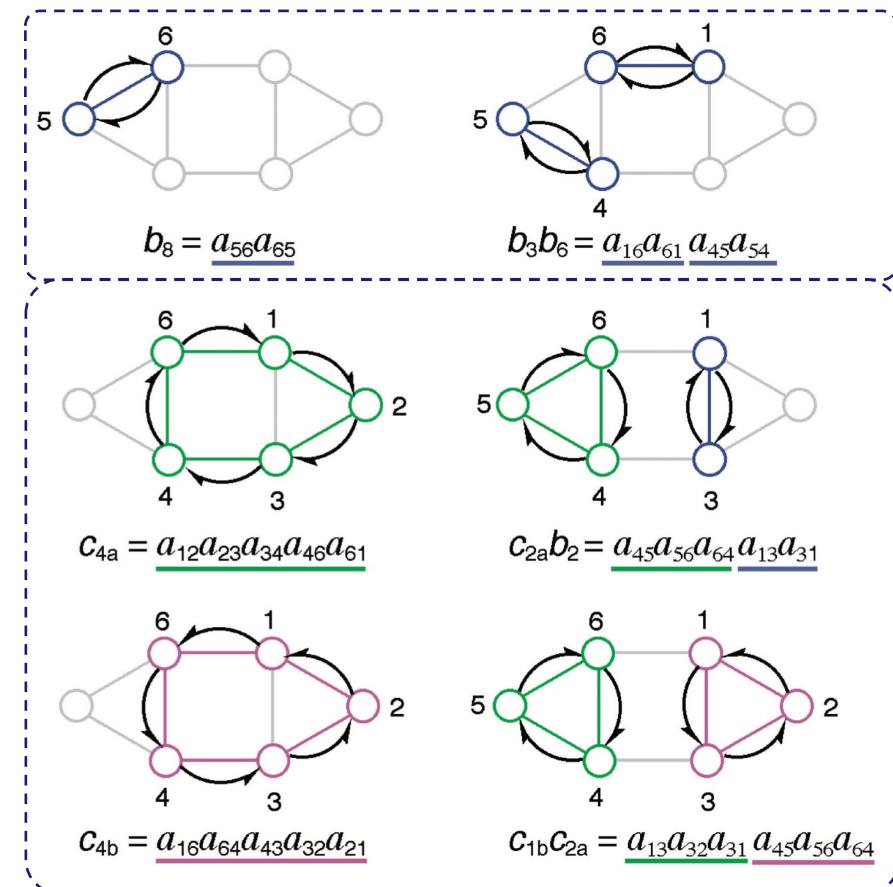
Constructing the reference polynomial



Expand the characteristic polynomial of the Hückel matrix

$$\begin{aligned} P_G(X) = & X^6 + (-\mathbf{b}_1 - \mathbf{b}_2 - \mathbf{b}_3 - \mathbf{b}_4 - \mathbf{b}_5 - \mathbf{b}_6 - \mathbf{b}_7 - \mathbf{b}_8)X^4 \\ & + (-\mathbf{c}_{1a} - \mathbf{c}_{1b} - \mathbf{c}_{2a} - \mathbf{c}_{2b})X^3 + (\mathbf{b}_1\mathbf{b}_5 + \mathbf{b}_1\mathbf{b}_6 \\ & + \mathbf{b}_1\mathbf{b}_7 + \mathbf{b}_1\mathbf{b}_8 + \mathbf{b}_2\mathbf{b}_6 + \mathbf{b}_2\mathbf{b}_7 + \mathbf{b}_2\mathbf{b}_8 + \mathbf{b}_3\mathbf{b}_4 \\ & + \mathbf{b}_3\mathbf{b}_5 + \mathbf{b}_3\mathbf{b}_6 + \mathbf{b}_4\mathbf{b}_6 + \mathbf{b}_4\mathbf{b}_7 + \mathbf{b}_4\mathbf{b}_8 + \mathbf{b}_5\mathbf{b}_8 \\ & - \mathbf{c}_{3a} - \mathbf{c}_{3b})X^2 + (\mathbf{c}_{1a}\mathbf{b}_6 + \mathbf{c}_{1a}\mathbf{b}_7 + \mathbf{c}_{1a}\mathbf{b}_8 + \mathbf{c}_{1b}\mathbf{b}_6 \\ & + \mathbf{c}_{1b}\mathbf{b}_7 + \mathbf{c}_{1b}\mathbf{b}_8 + \mathbf{c}_{2a}\mathbf{b}_1 + \mathbf{c}_{2a}\mathbf{b}_4 + \mathbf{c}_{2a}\mathbf{b}_2 + \mathbf{c}_{2b}\mathbf{b}_1 \\ & + \mathbf{c}_{2b}\mathbf{b}_4 + \mathbf{c}_{2b}\mathbf{b}_2 - \mathbf{c}_{4a} - \mathbf{c}_{4b} - \mathbf{c}_{5a} - \mathbf{c}_{5b})X \\ & + (-\mathbf{b}_1\mathbf{b}_5\mathbf{b}_8 - \mathbf{b}_4\mathbf{b}_6\mathbf{b}_3 + \mathbf{c}_{1a}\mathbf{c}_{2a} + \mathbf{c}_{1a}\mathbf{c}_{2b} + \mathbf{c}_{1b}\mathbf{c}_{2a} \\ & + \mathbf{c}_{1b}\mathbf{c}_{2b} - \mathbf{c}_{6a} - \mathbf{c}_{6b}) \end{aligned}$$

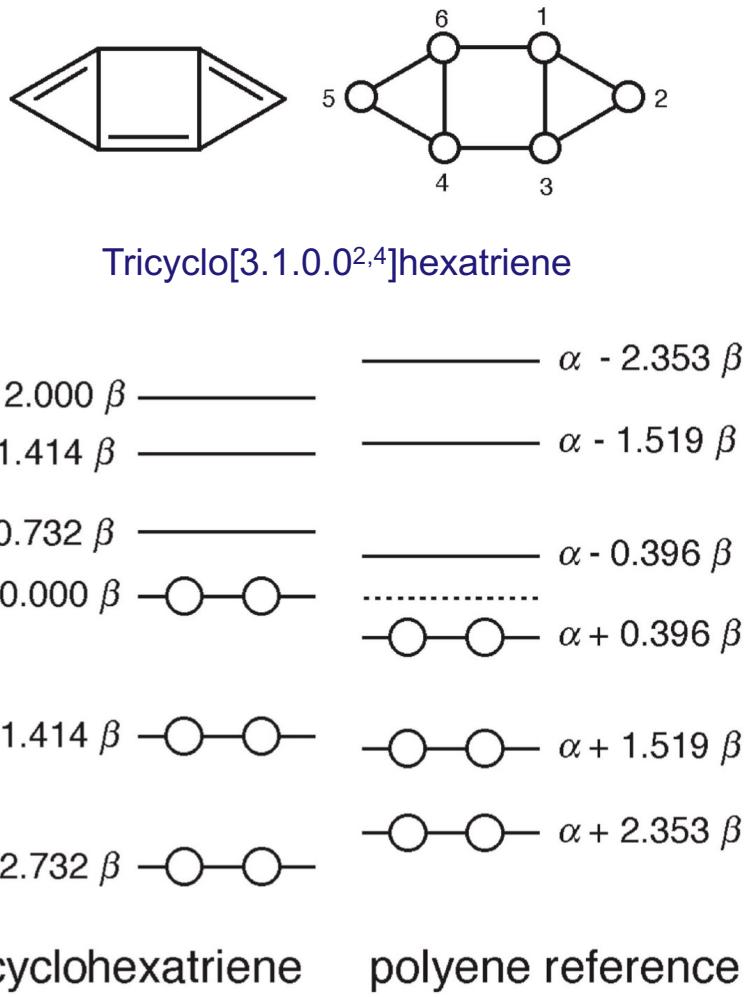
Sachs subgraphs



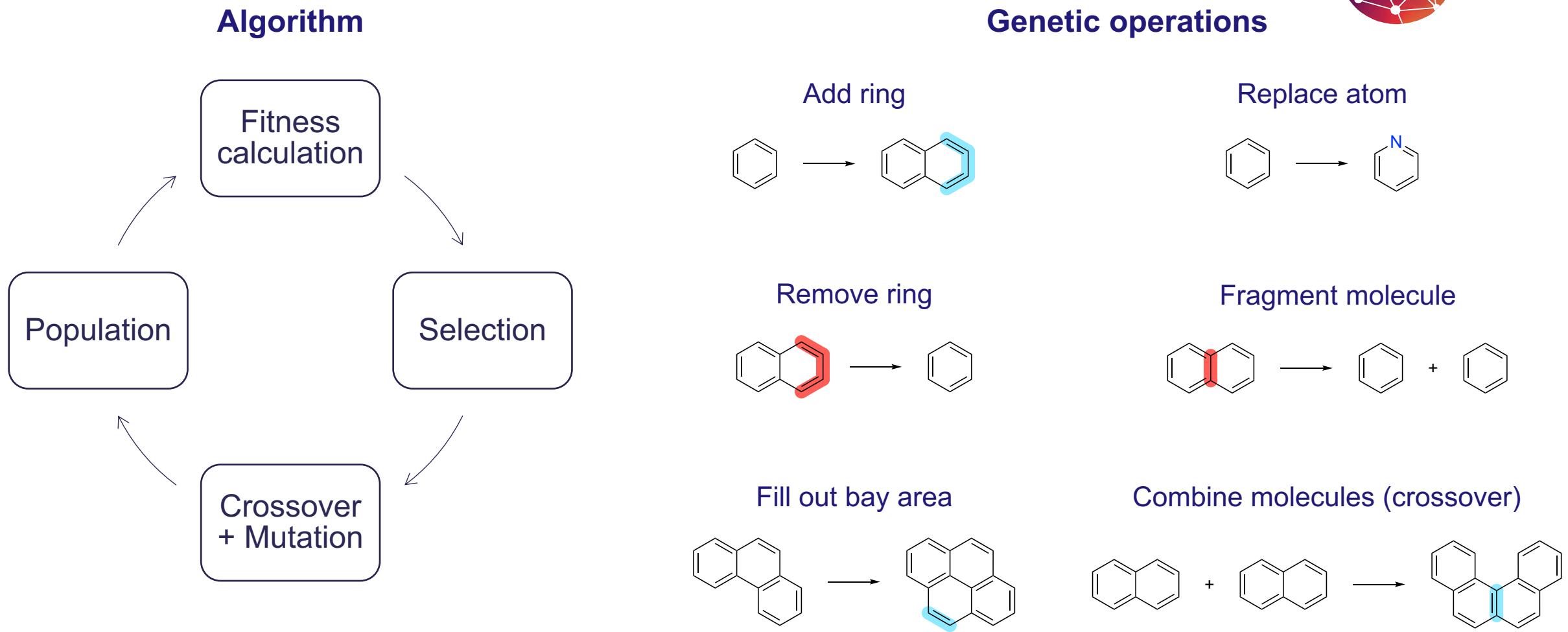
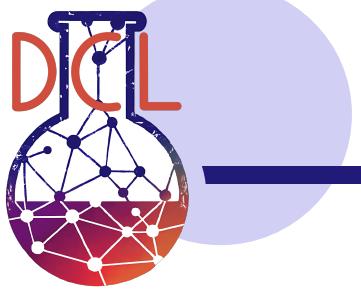
Topological resonance energy



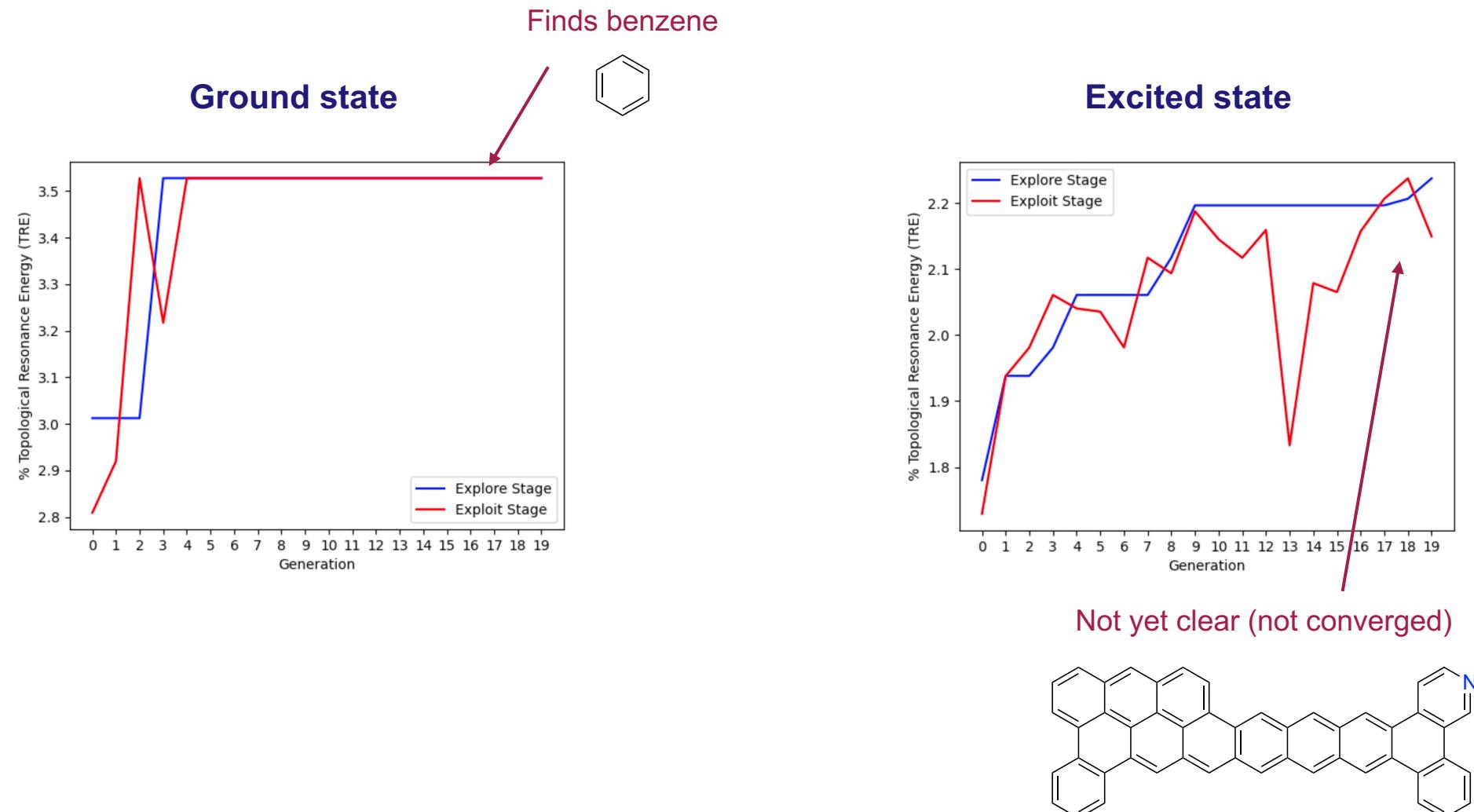
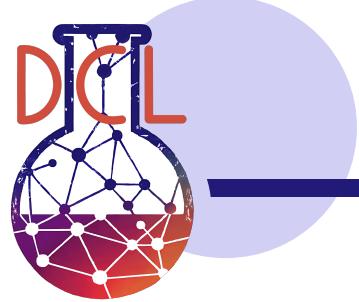
- $P_G(X) \Rightarrow$ roots $X \Rightarrow$ orbital energies ϵ
- $R_G(X) \Rightarrow$ roots $X \Rightarrow$ orbital energies ϵ
- $\text{TRE} = -(E_\pi - E_\pi^{\text{ref}}) = -0.2425|\beta|$
 - Positive TRE \Rightarrow aromatic
 - Negative TRE \Rightarrow anti-aromatic
- Normalized metric
 - $\% \text{TRE} = \frac{\text{TRE}}{E_\pi^{\text{ref}}}$



Genetic algorithm



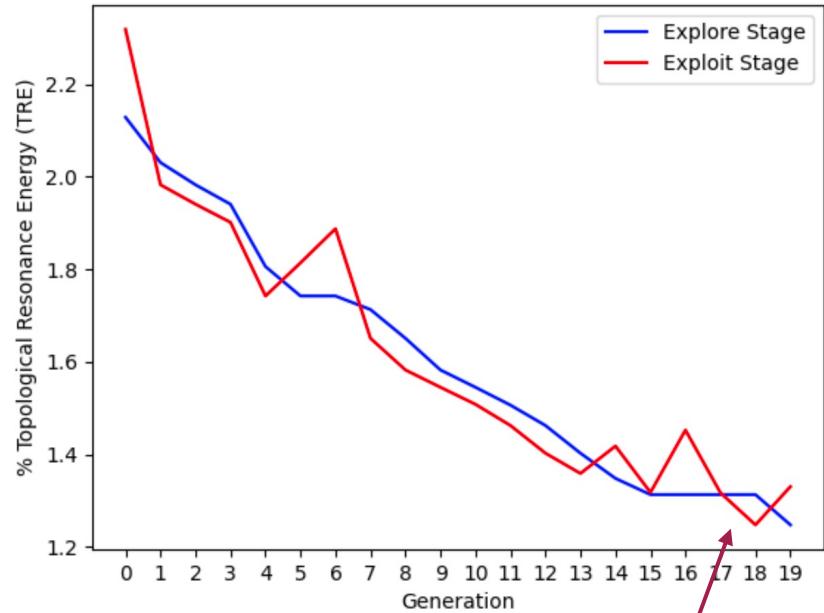
Maximizing %TRE – more aromatic



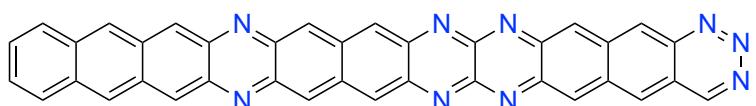
Minimizing %TRE – less aromatic



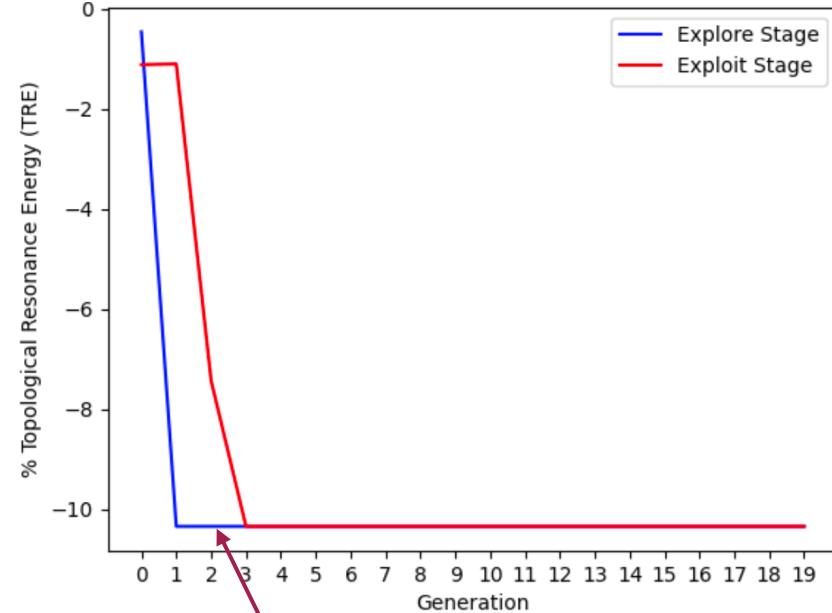
Ground state



Finds longer acenes (not converged)



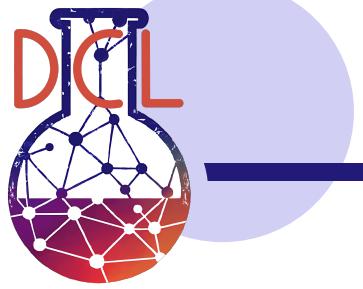
Excited state



Finds benzene

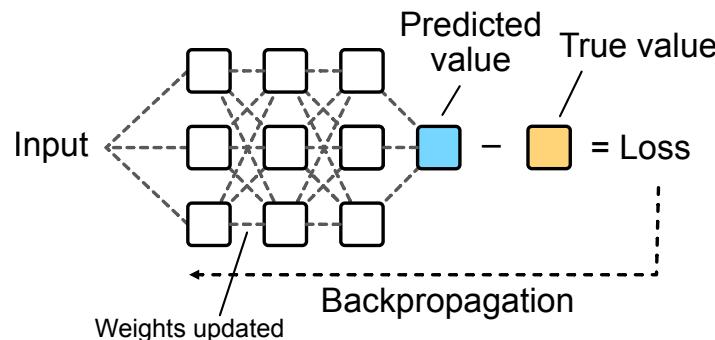


Hybrid ML & physical models

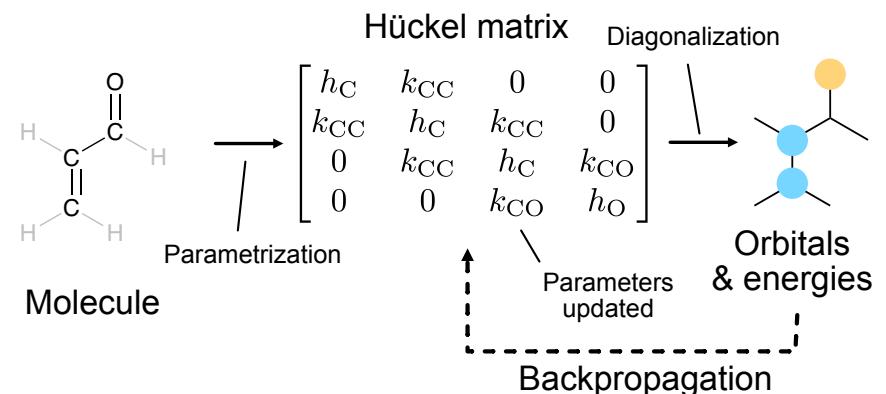


- Physical models can be coded within auto-differentiable frameworks
- Hybrids are regularized by the physical model
- Can in principle be legacy code via Enzyme LLVM compiler

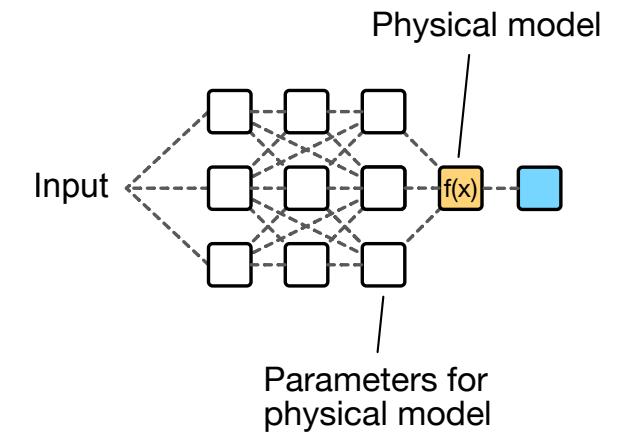
Neural network



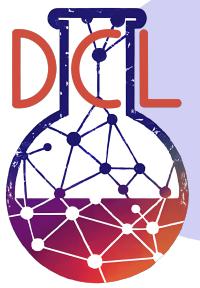
Differentiable model



Hybrid

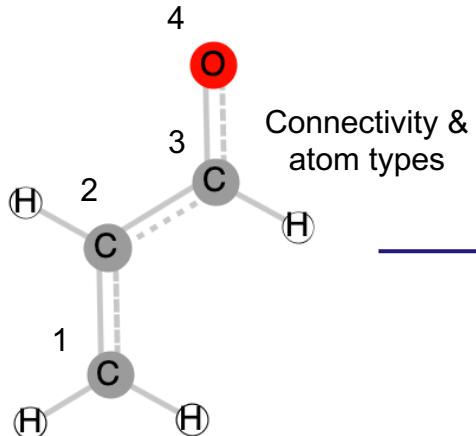


Auto-differentiable Hückel theory



- One of the first molecular orbital theories
- Schrödinger equation reduces to a very simple form

Molecule

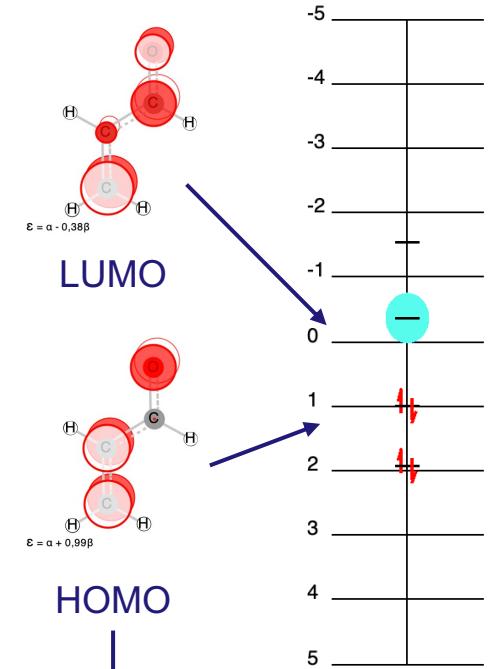


Hückel matrix

$$\begin{bmatrix} h_C & k_{CC} & 0 & 0 \\ k_{CC} & h_C & k_{CC} & 0 \\ 0 & k_{CC} & h_C & k_{CO} \\ 0 & 0 & k_{CO} & h_O \end{bmatrix} = \begin{bmatrix} 0 & 1 & 0 & 0 \\ 1 & 0 & 1 & 0 \\ 0 & 1 & 0 & 1.06 \\ 0 & 0 & 1.06 & 0.97 \end{bmatrix}$$

Diagonalize

Orbitals & energies



Backpropagate error to parameters

Automatic differentiation

Translating NumPy code to JAX



- jax.numpy & jax.scipy

COULSON

```
import numpy as np
eig_vals, eig_vects = np.linalg.eigh(huckel_matrix)
```



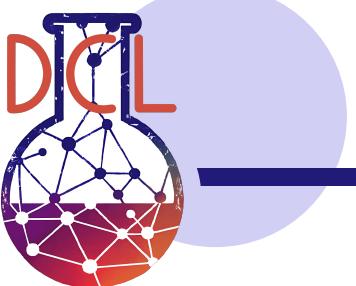
HÜXEL

```
import jax.numpy as jnp
eig_vals, eig_vects = jnp.linalg.eigh(huckel_matrix)
```

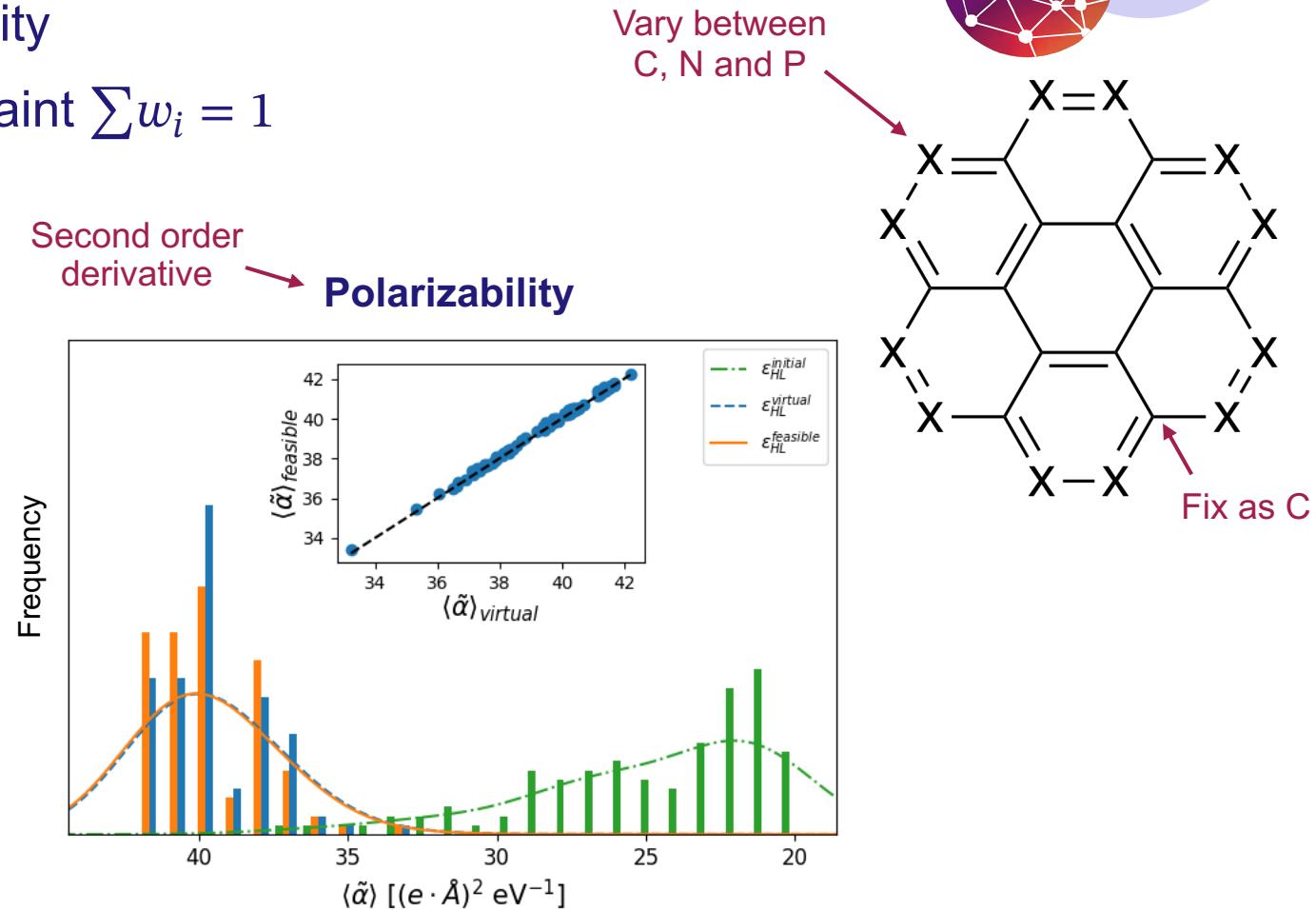
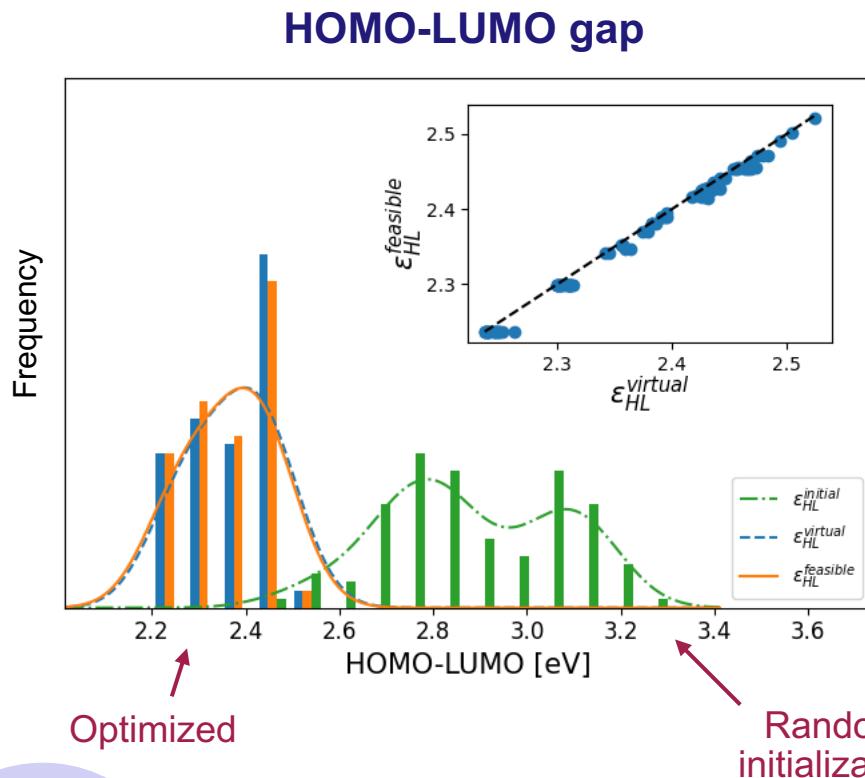
+

Code for parameter optimization

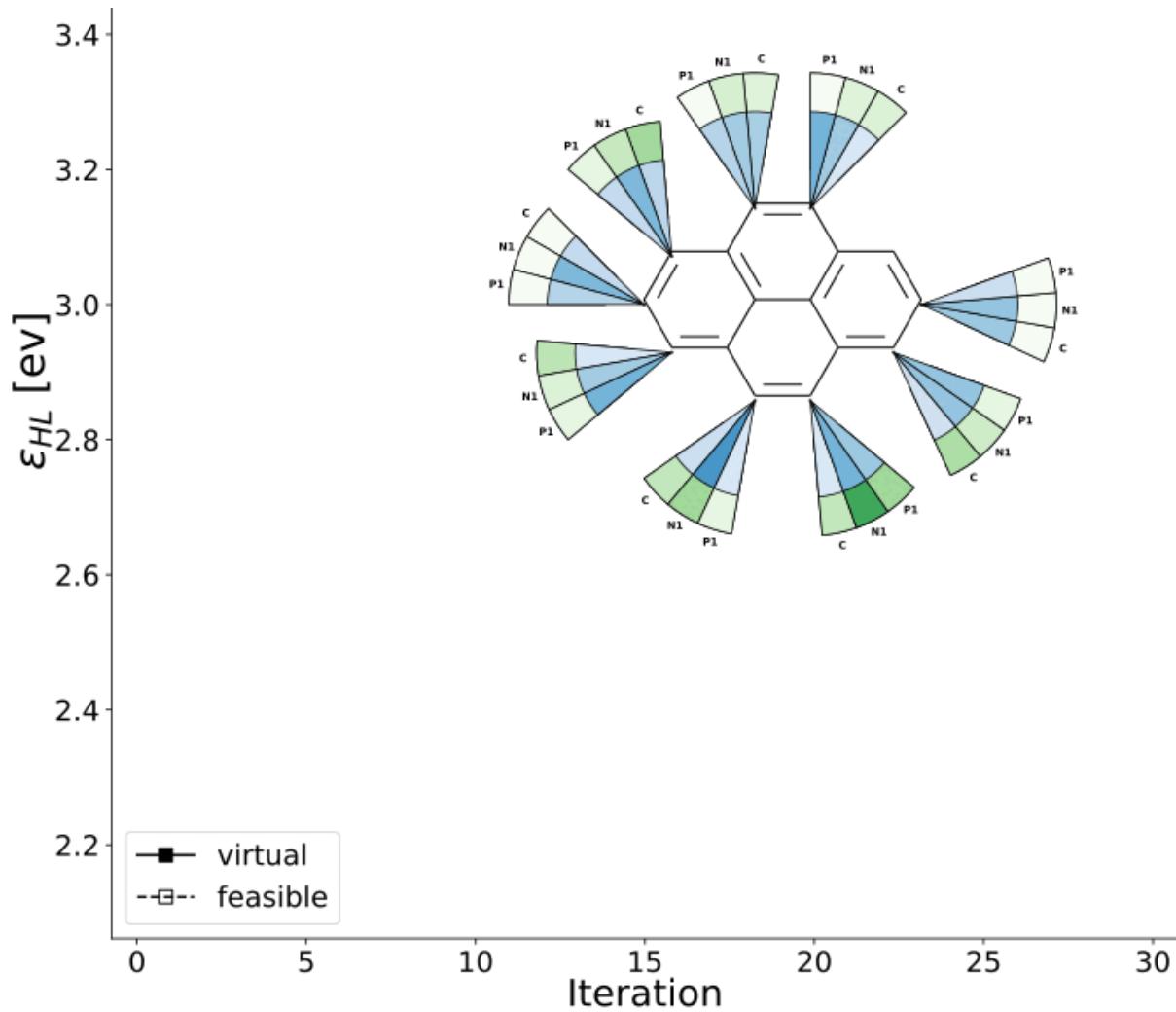
Inverse design

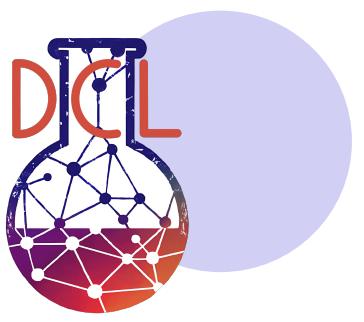
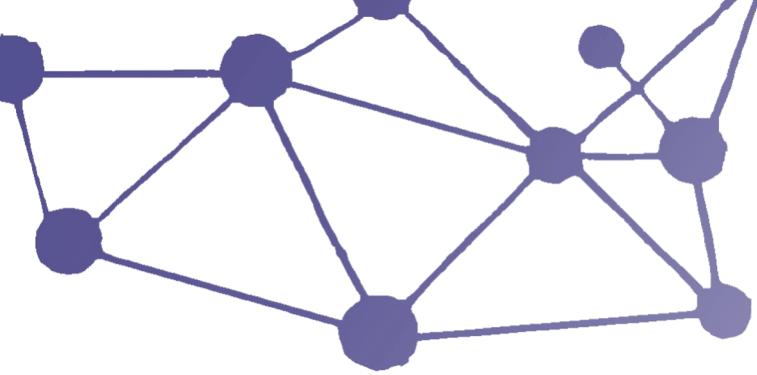


- “Alchemical” optimization of atom identity
- Optimize atom “weights” w_i with constraint $\sum w_i = 1$



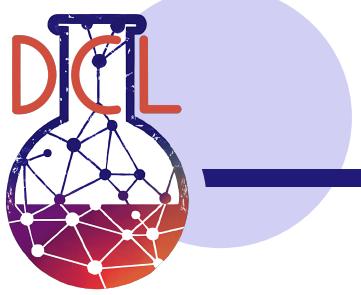
Inverse design in action



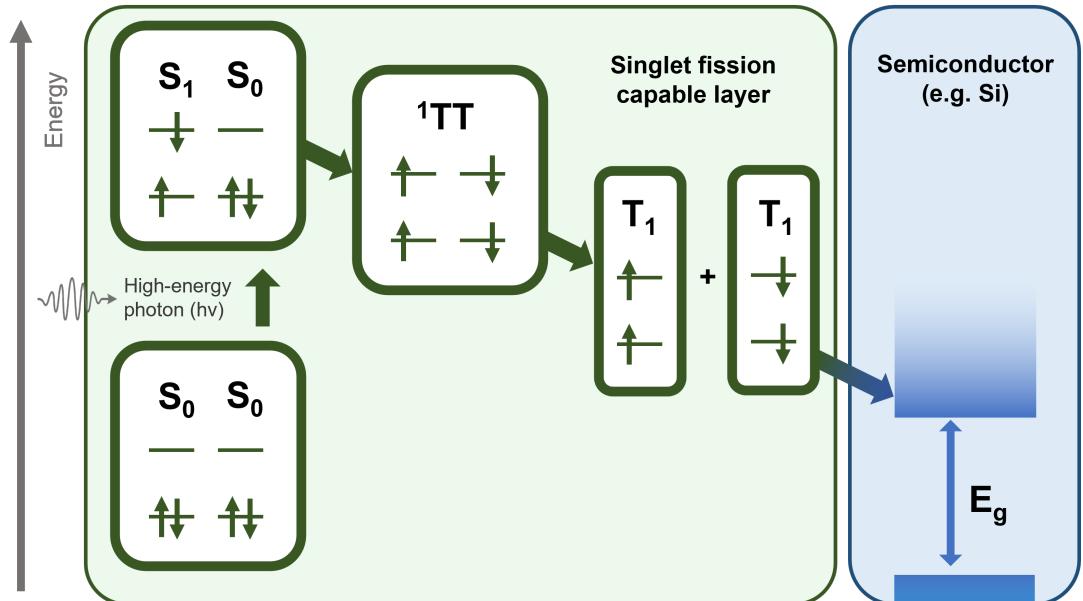


Designing singlet fission materials with uncertainty-aware genetic algorithms

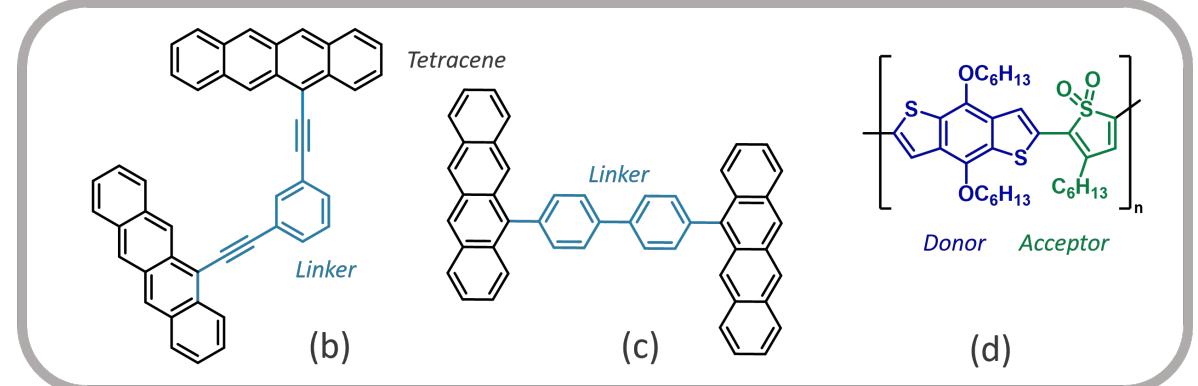
Singlet fission



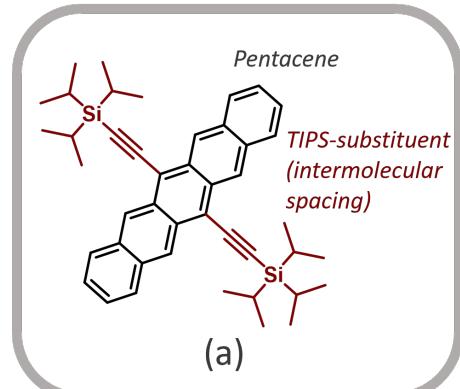
- Allow harvesting of high-energy photons to match band-gap of solar cell
- Relatively few molecules known



Intramolecular SF



Intermolecular SF

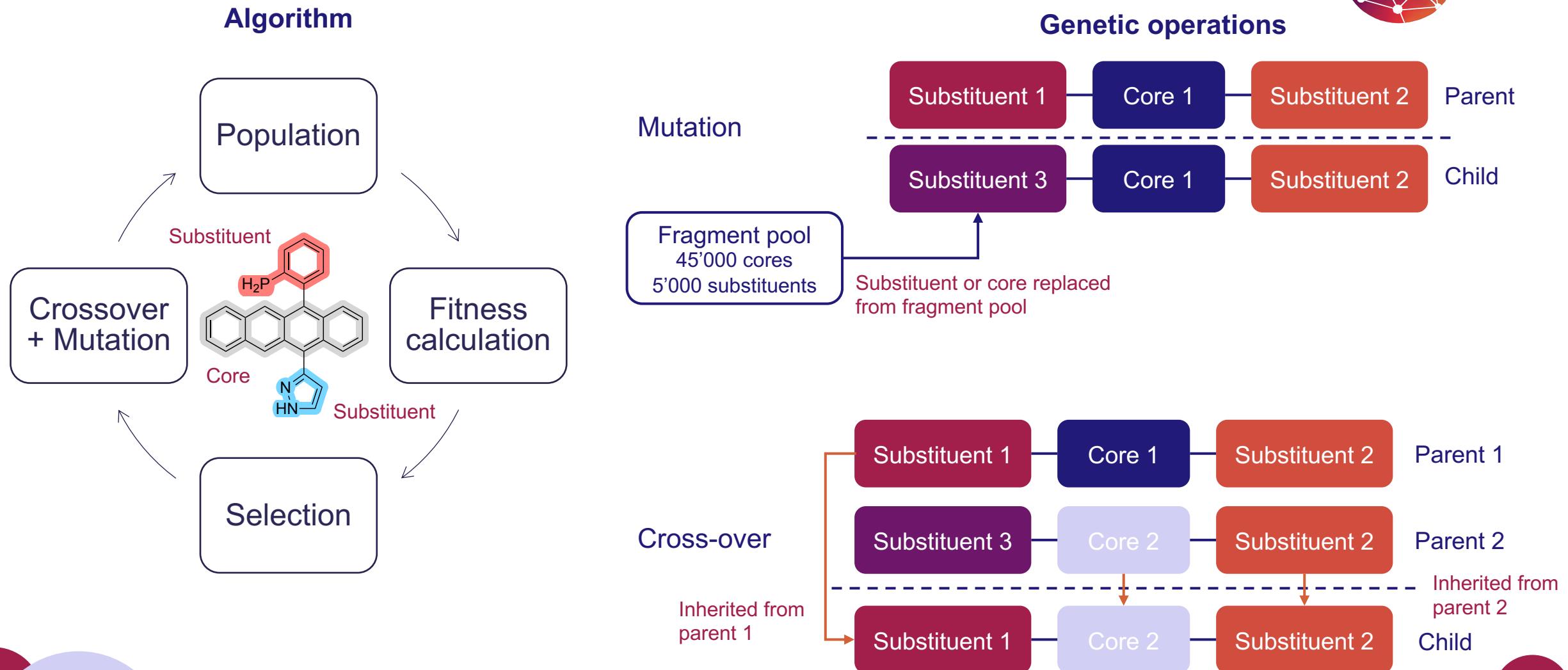
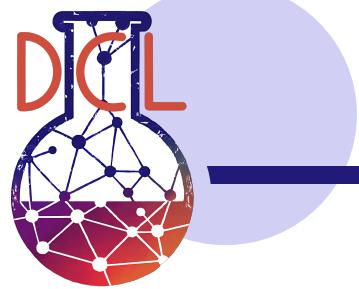


Luca Schaufelberger

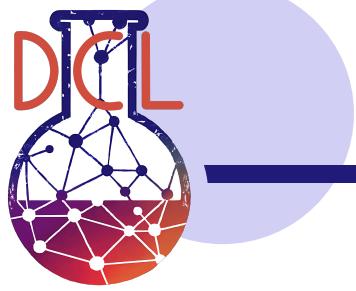


Clémence Corminboeuf

Genetic algorithm



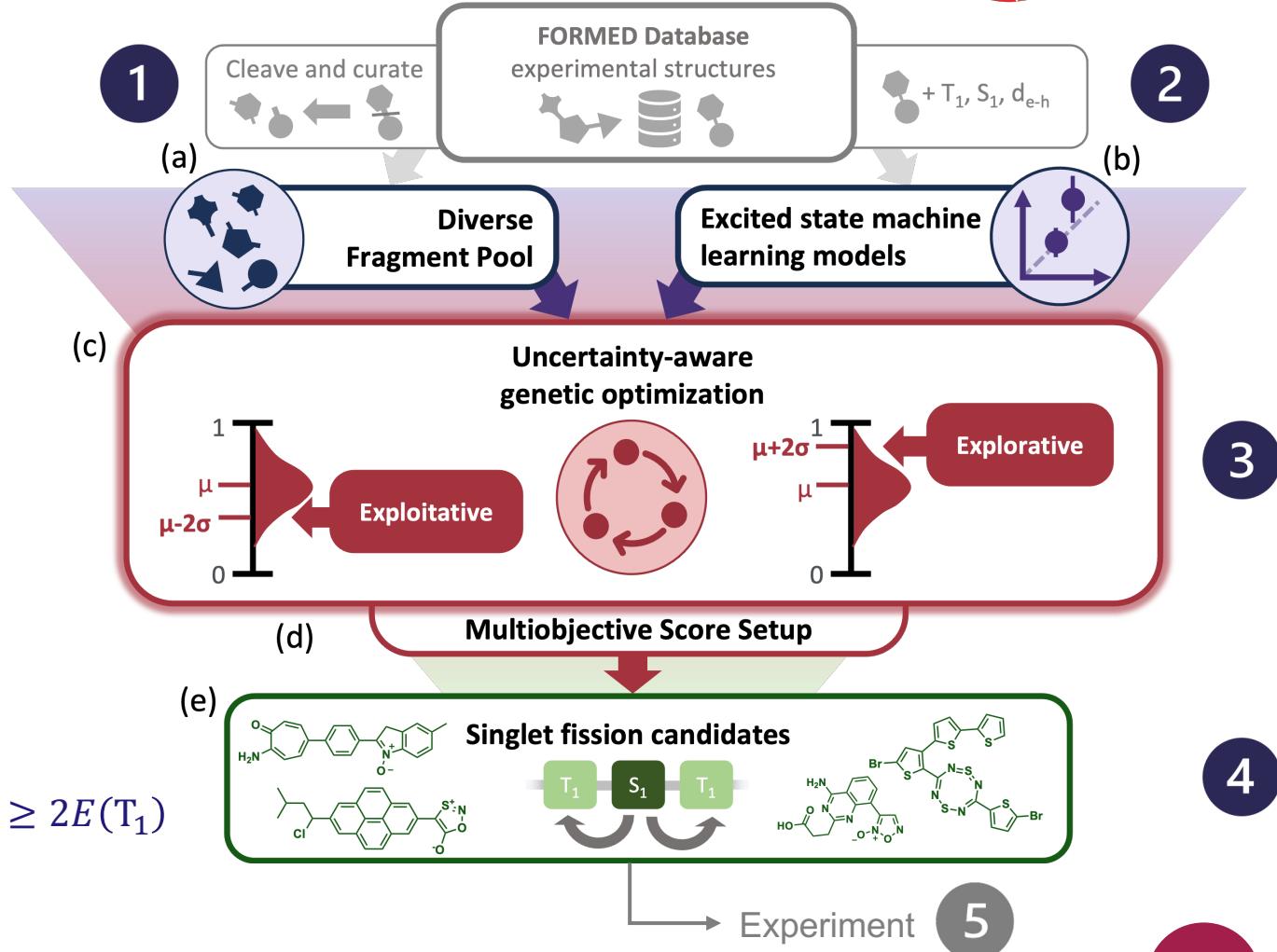
Optimization with genetic algorithm



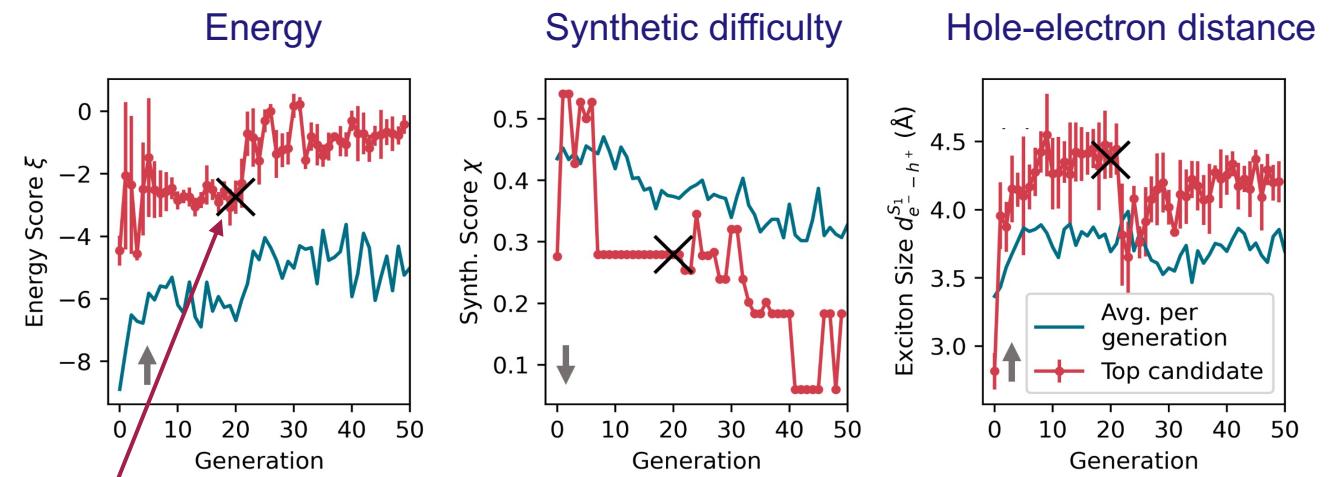
1. Create a pool of fragments from FORMED database
2. Build machine learning models for molecular properties
3. Run genetic algorithm in **explorative** or **exploitative** mode
4. Collect candidates and further study them with DFT simulations
5. (Experimental testing)

Optimization criteria

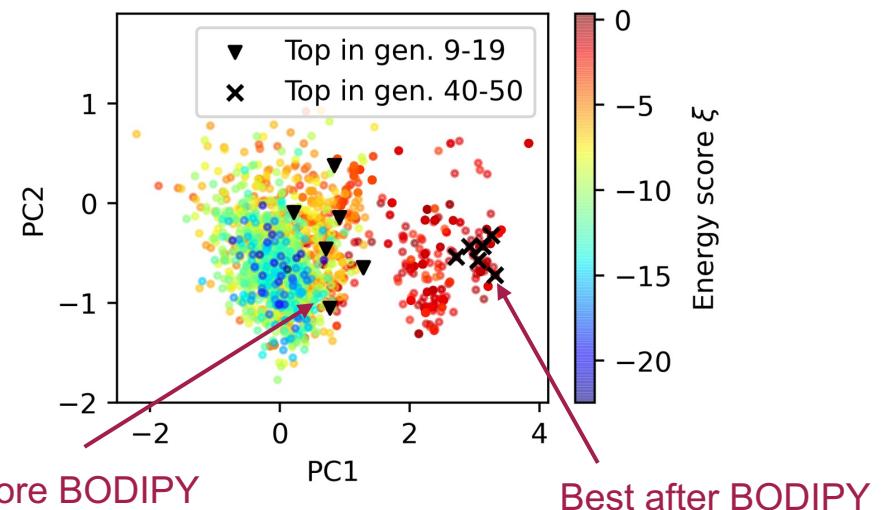
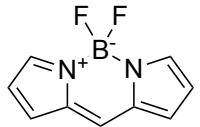
1. Singlet fission energy score $E(S_1) \geq 2E(T_1)$
2. Synthetic difficulty
3. Electron-hole separation



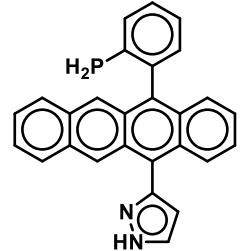
Results



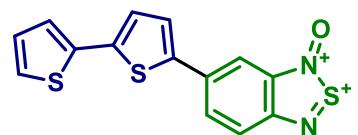
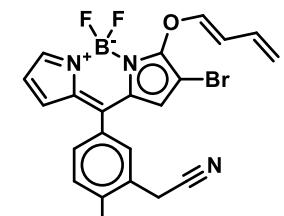
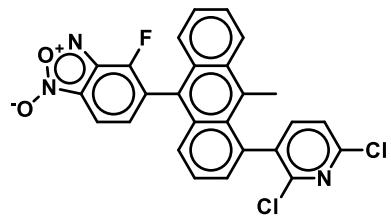
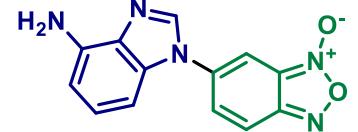
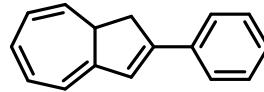
“Discovery” of BODIPY core



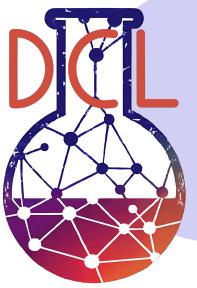
Rediscovered
Exploitation



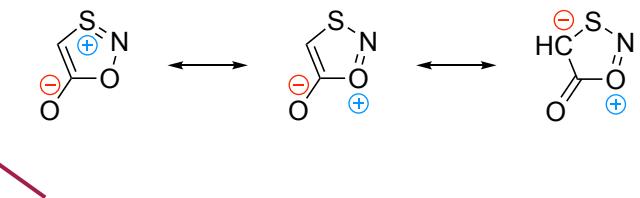
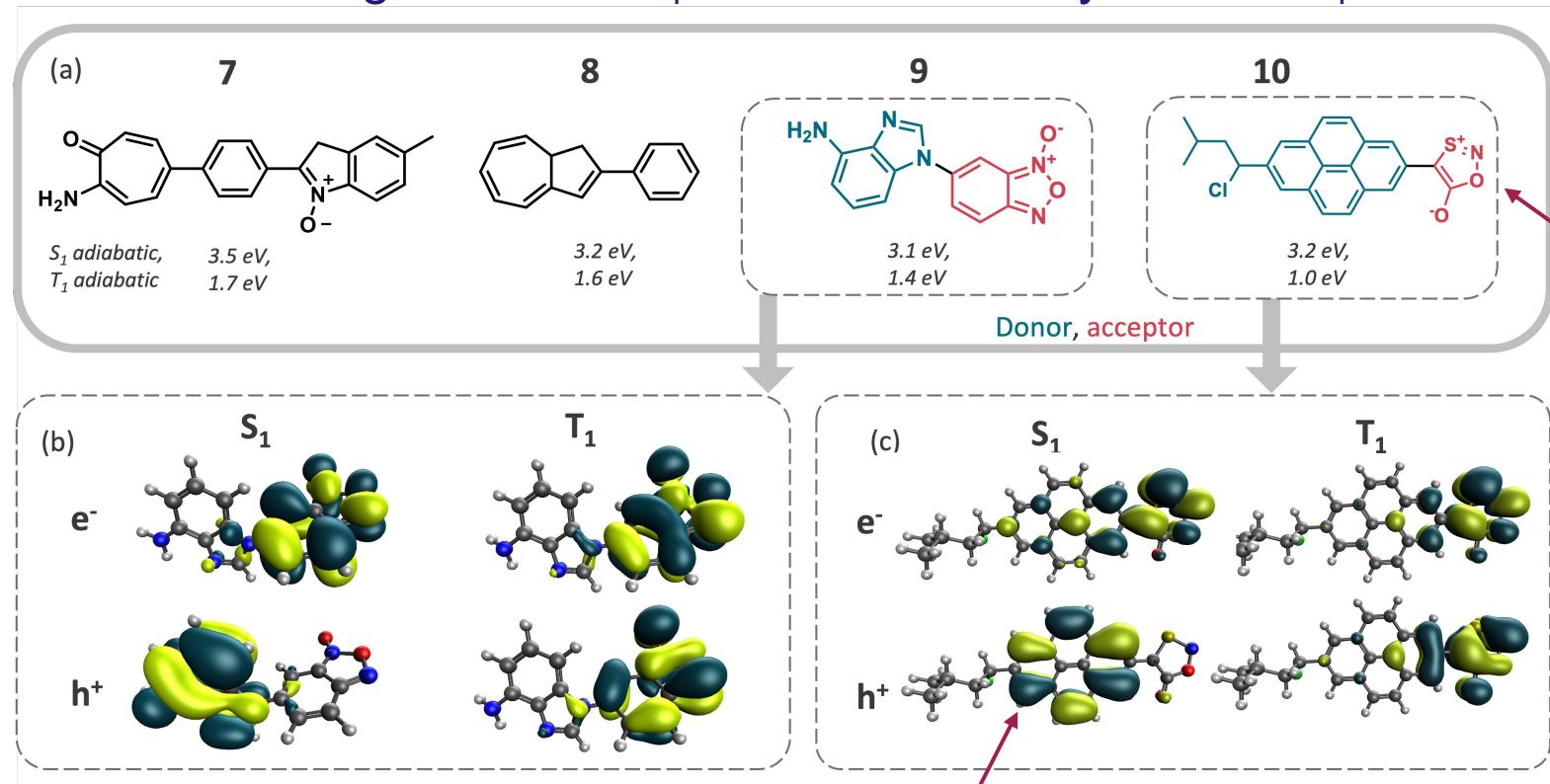
Previously unknown (?)
Exploration



Mesoionic acceptors



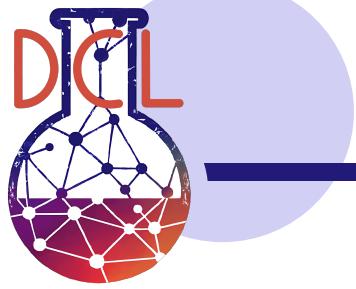
- Exotic electronic structure where no charge-neutral resonance forms can be drawn
- Feature charge-transfer S_1 states and locally excited T_1 states



"Mesoionic compounds are one in which a heterocyclic structure is dipolar and where both the negative and the positive charges are delocalized"

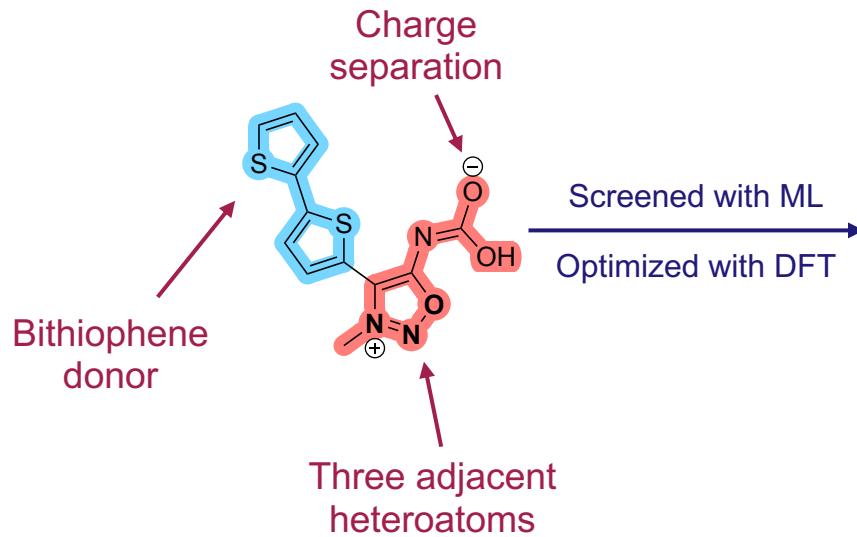
Electron and hole densities in T_1 are localized on same fragment (local excitation)

Further exploration from identified trends

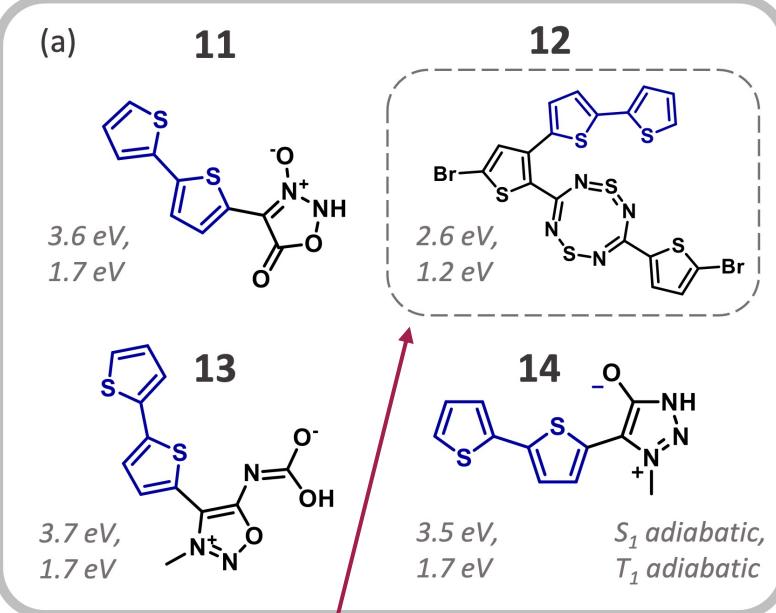


- Exploiting all mesoionic-like fragments in the original FORMED database
- Combine them with bithiophene donors and screening with machine learning and DFT

Combination of fragments

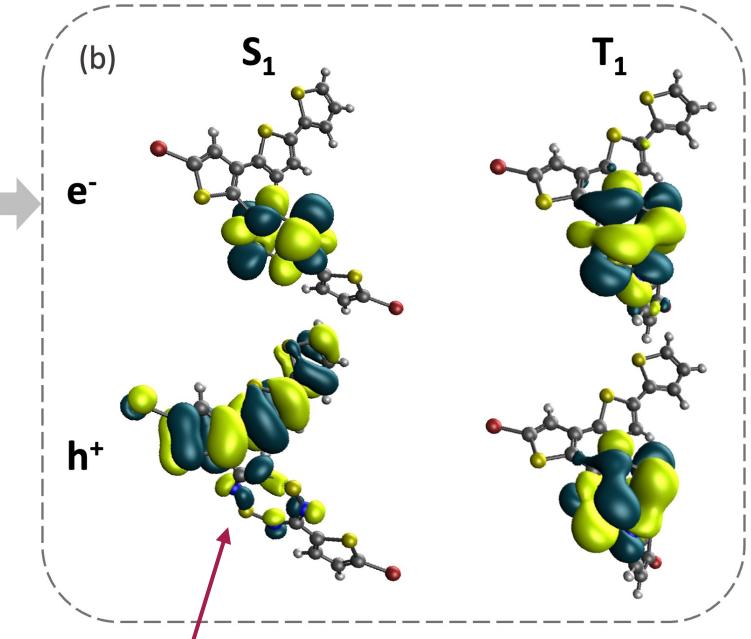


Most promising candidates

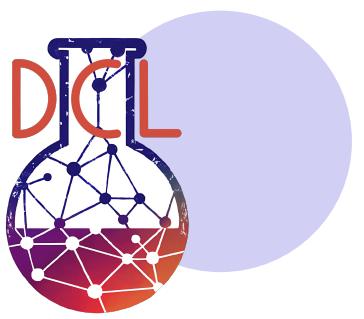
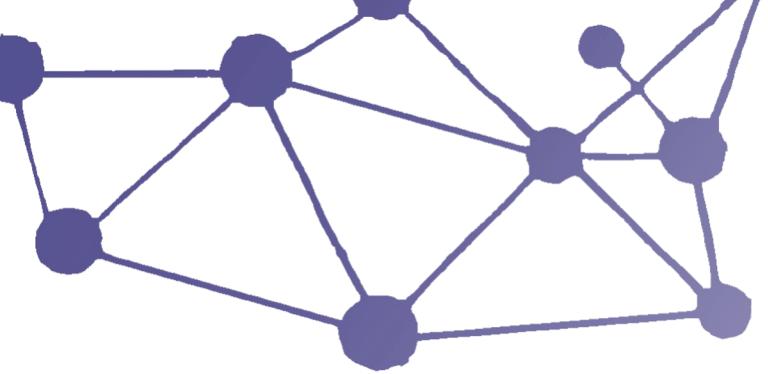


Collaboration for experimental testing of candidate

Excited state characterization

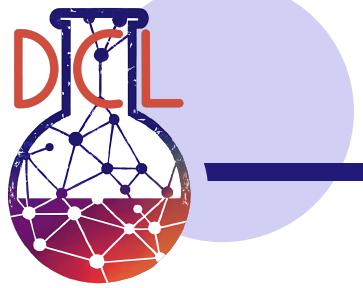


Also displays charge-transfer S₁ and locally excited T₁



Reaction design with reactive force field methods

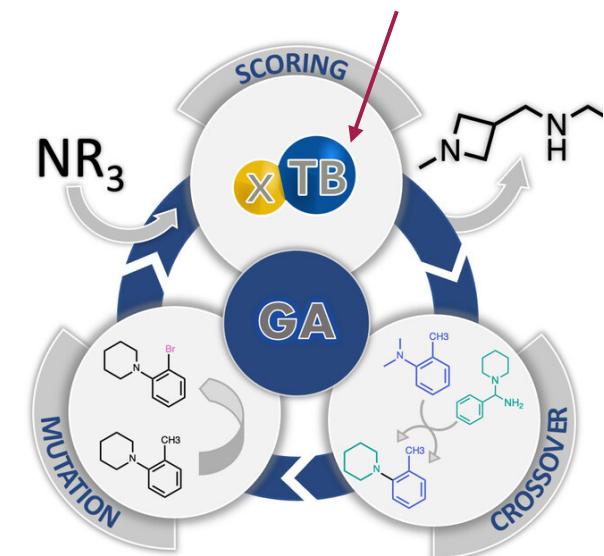
Simulating reactive chemistry



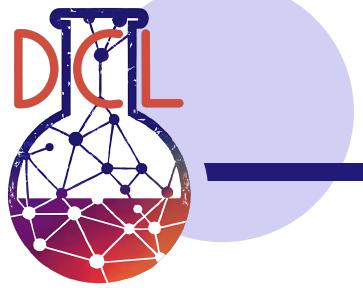
- Density functional theory is the go-to-method for simulations and dataset generation
- Computational cost is prohibitively high to use as property predictor for generative models
- Current work use reactive force fields or semi-empirical methods

Typical reaction simulations	This work
High computational cost	Fast force fields
High failure rate (20–70%)	Robustness of >99.9%
High variance of outputs	Exhaustive sampling
Doubtful correlation to experiment	Initial validation to experiment

Semi-empirical methods can be used in favorable cases

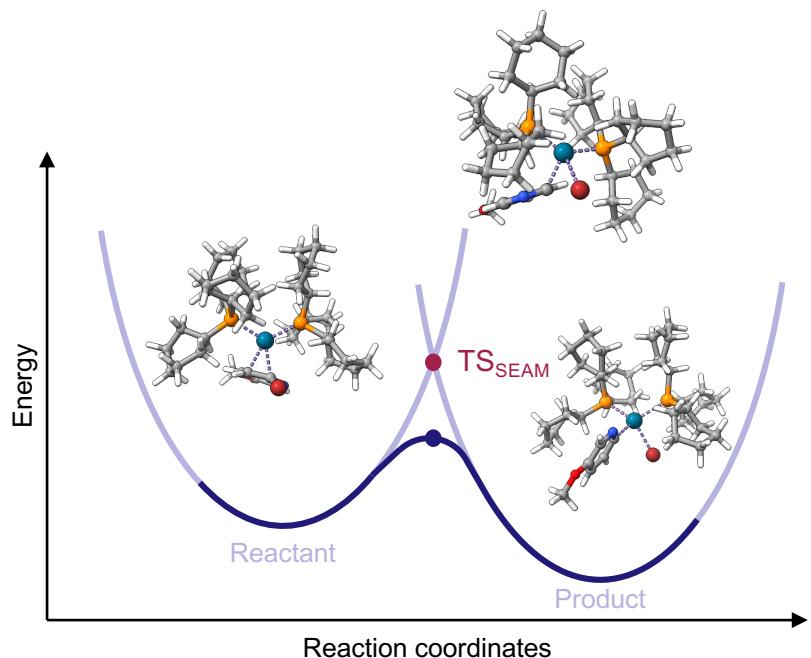


Reactive force fields



- Simulate chemical reactivity approximately with the cost of force fields
- The SEAM method avoids reaction-specific parametrization

SEAM method



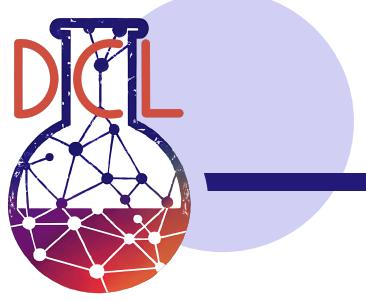
Simplified code with our POLANYI package

```
from polanyi.workflow import opt_ts_python
from polanyi.io import read_xyz

elements_r, coordinates_r = read_xyz("reactant.xyz")
elements_p, coordinates_p = read_xyz("product.xyz")

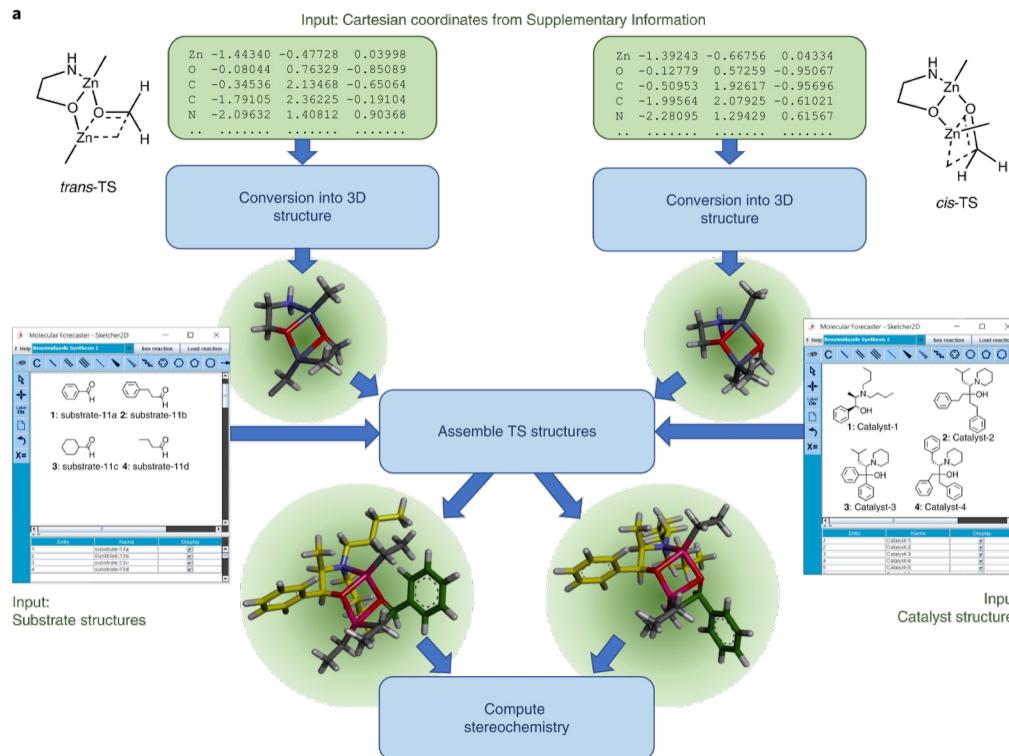
results = opt_ts_python(
    elements_r,
    [coordinates_r, coordinates_p]
)
coordinates_ts = results.coordinates_opt
```

Similar methods

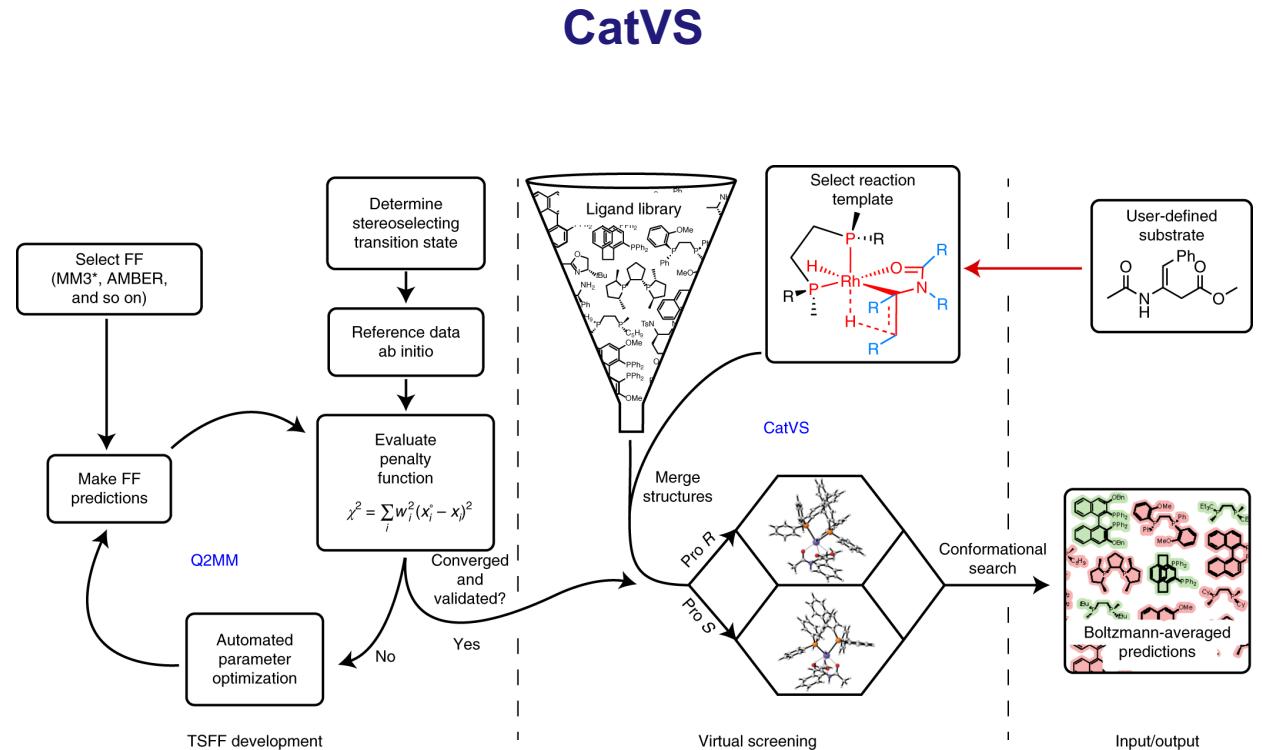


- Asymmetric catalyst evaluation (ACE) from Moitessier
- Quantum-Guided Molecular mechanics (Q2MM) from Norrby, Wiest

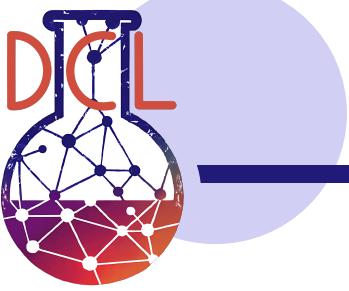
Virtual Chemist



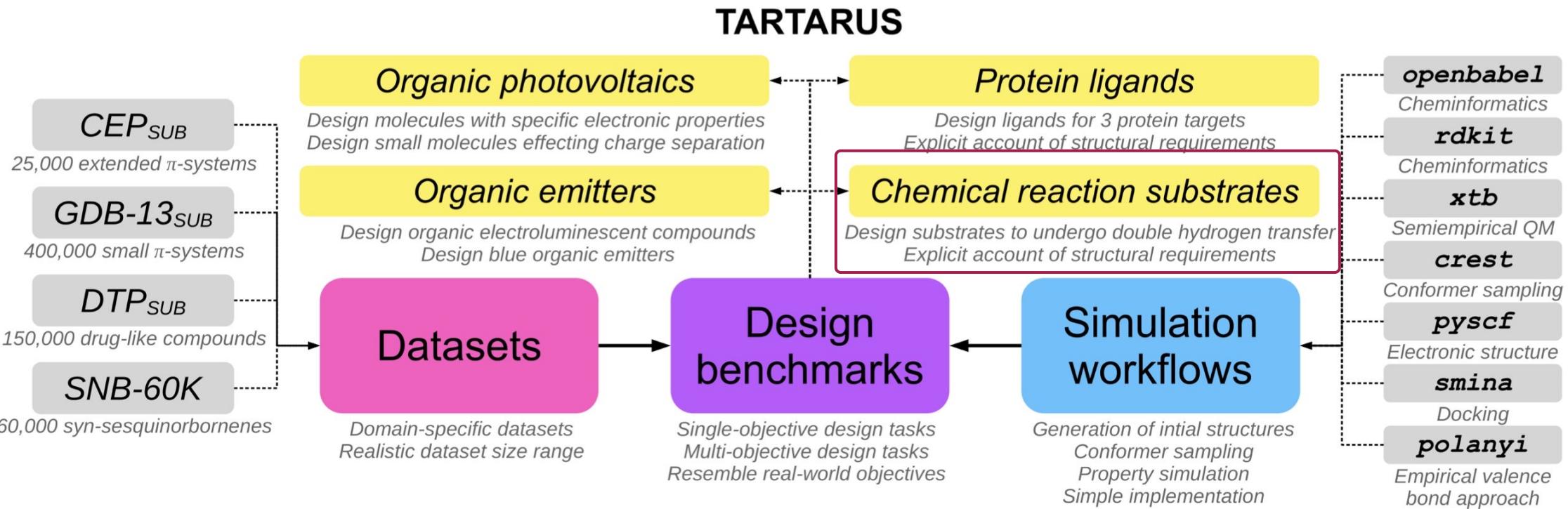
CatVS



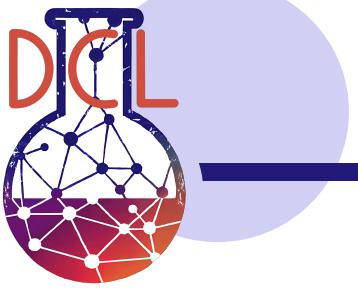
“Realistic” benchmark set for generative models



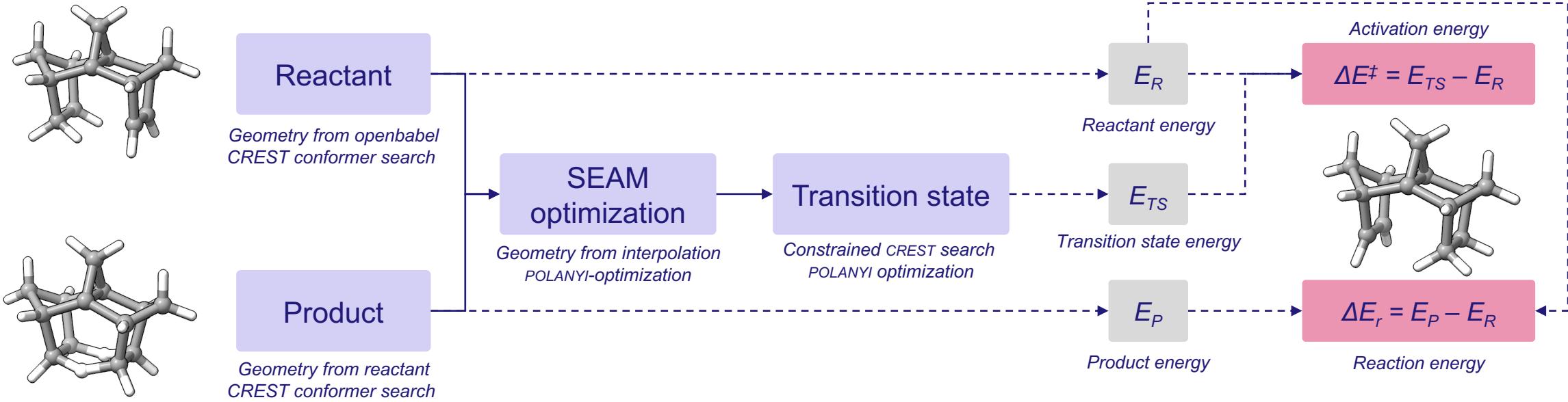
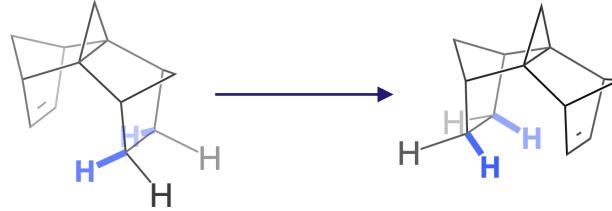
- Moderately complex tasks related to important applications
- Benchmarks based on fast simulations rather than surrogate ML models



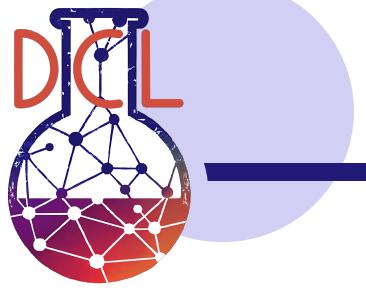
Reaction benchmark for generative models



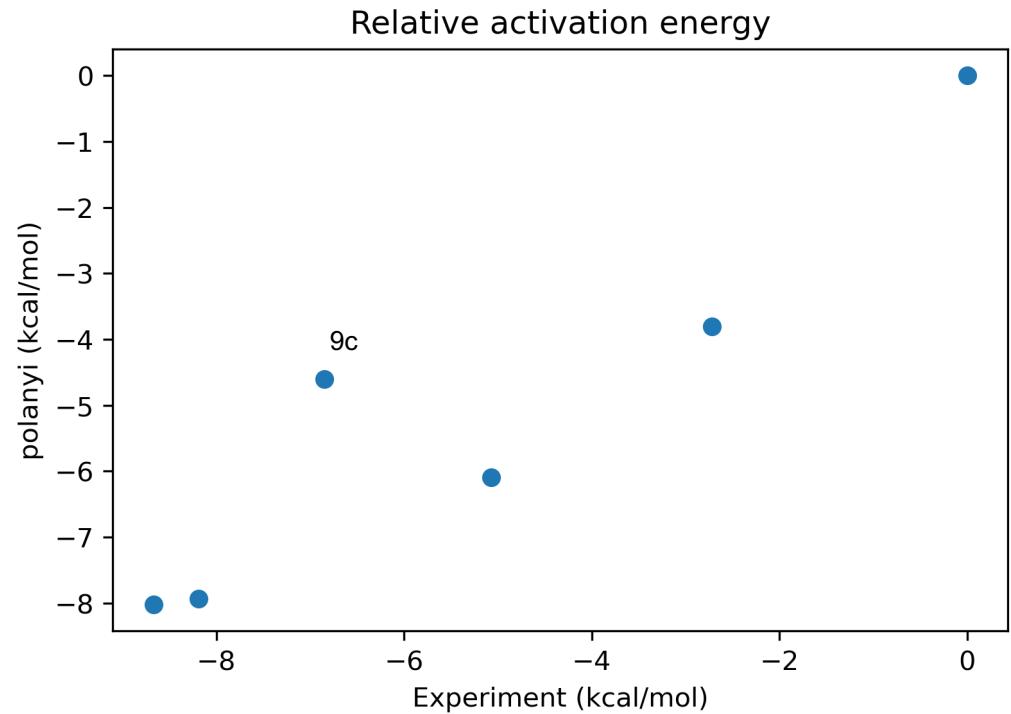
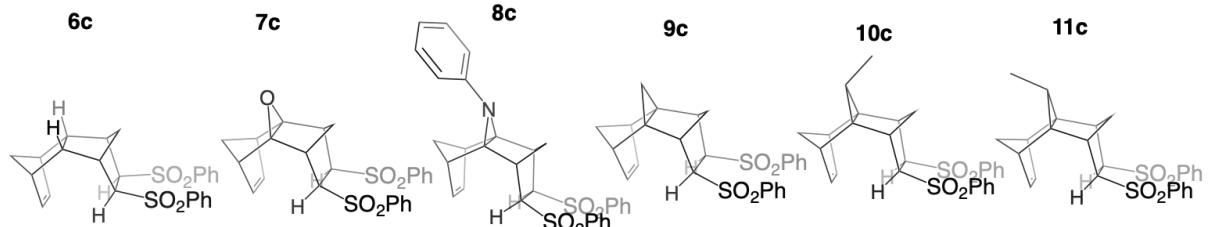
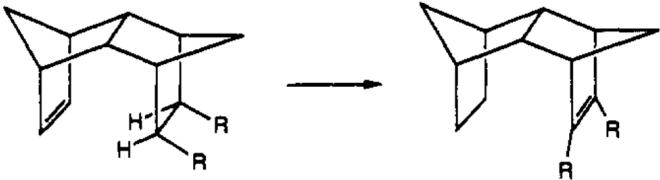
- Double hydrogen atom transfer of *syn*-sesquinorbornenes
- Benchmarked against experiment with SEAM
- Very robust optimization



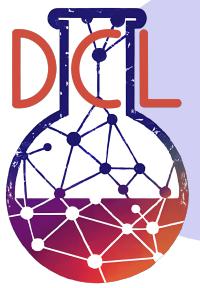
Reaction benchmark for generative models



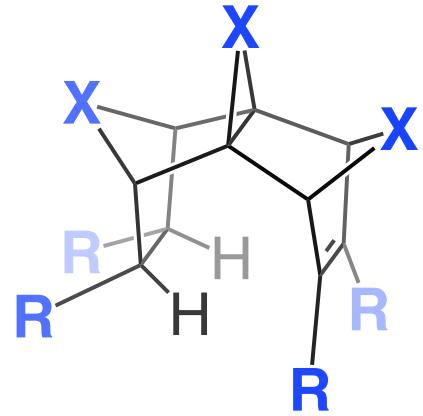
- Decent correlation with (very limited) experimental data
- No solvent model (exp: bromobenzene)



Dataset generation



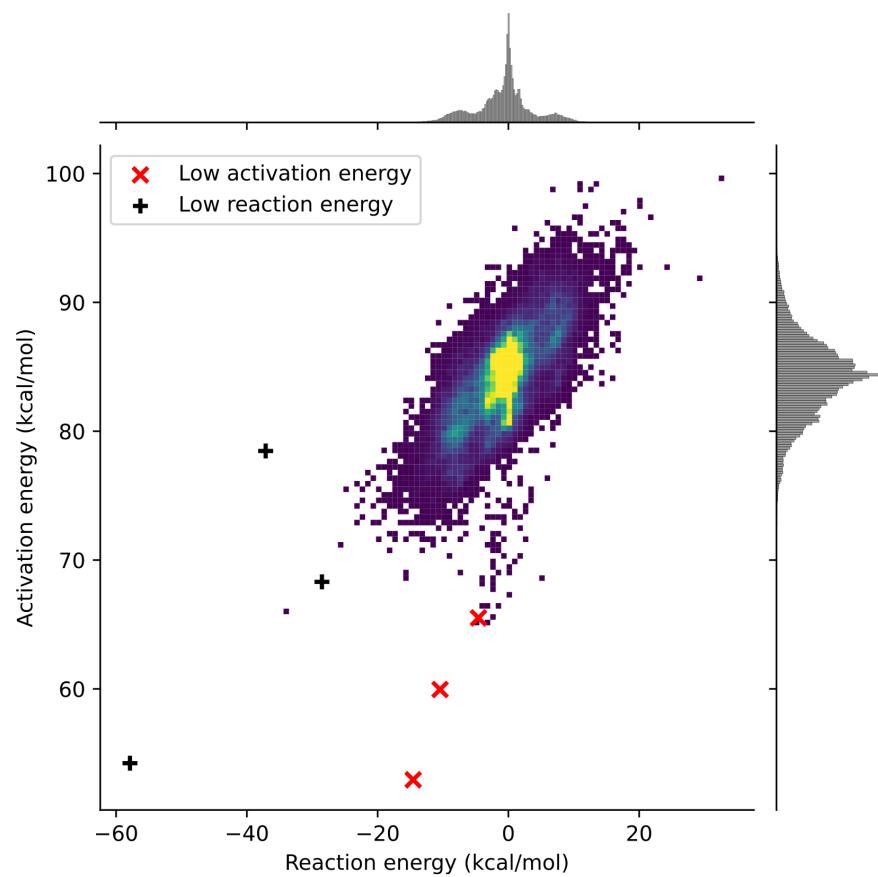
- Start from unsubstituted compound
- Perform STONED SELFIES mutations
- Filter for original fragment
- Obtain 60'850 examples
- Filter out "unwanted" fragments



```
[C-], [S-], [O-], [N-], '[*+], [*-] [PH], [pH], [N&X5], *= [S,s;!R], [S&X3], [S&X4],  
[S&X5], [S&X6], [P,p], [B,b,N,n,O,o,S,s]~[F,Cl,Br,I], *====, *#*,  
[O,o,S,s]~[O,o,S,s], [N,n,O,o,S,s]~[N,n,O,o,S,s]~[N,n,O,o,S,s],  
[N,n,O,o,S,s]~[N,n,O,o,S,s]~[C,c]=, :[O,o,S,s,N,n;!R], *=N-[*;!R], *~[N,n,O,o,S,s]-  
[N,n,O,o,S,s;!R]
```

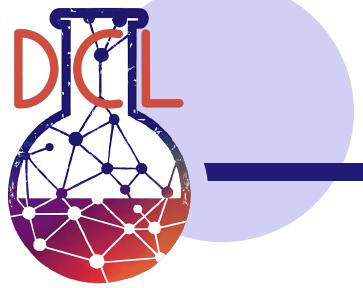
Dataset & generative modelling

- Reference dataset of 60,850 reactions
- 99.94% success rate
- ~ 12 min/molecule/CPU
- Limited by conformational sampling
- Interfaced with generative models

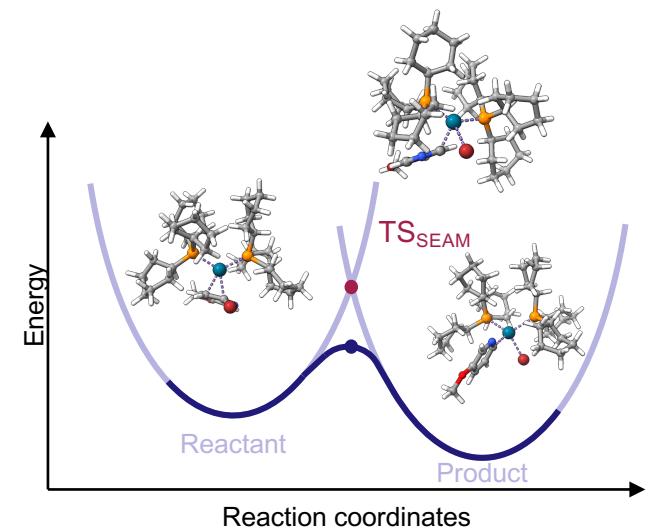
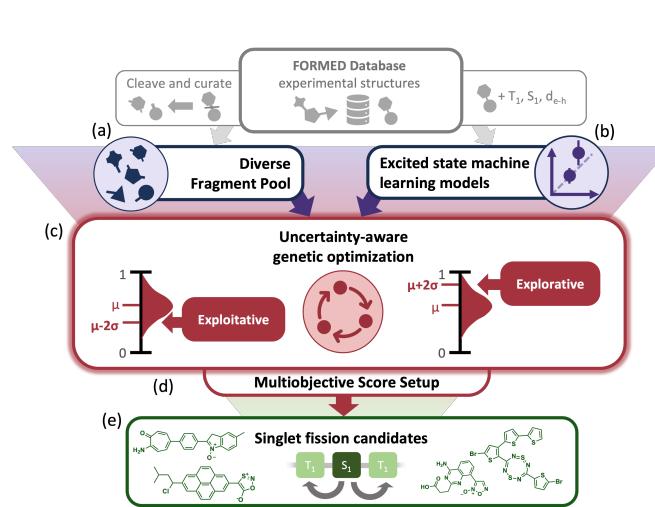
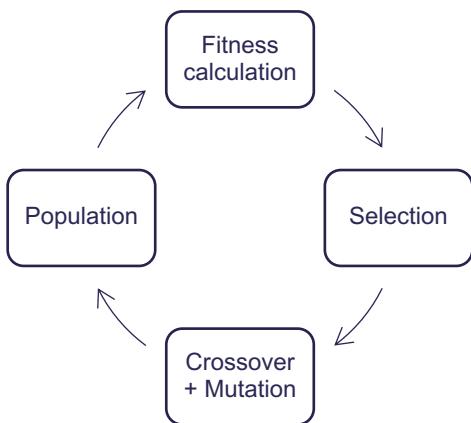
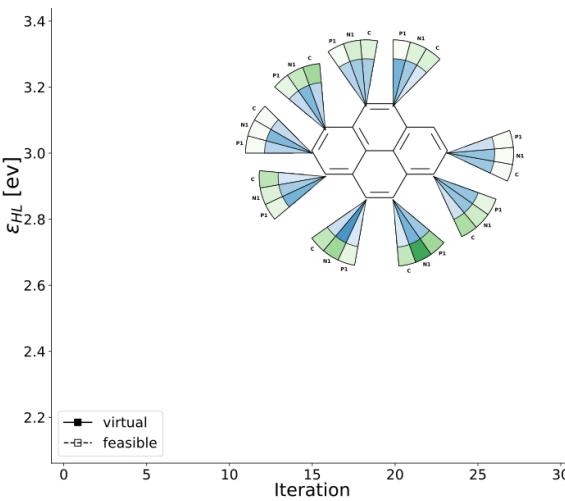


	ΔE^\ddagger	ΔE_r	$\Delta E^\ddagger + \Delta E_r$	$-\Delta E^\ddagger + \Delta E_r$
PARENT SUBSTRATE	85.16	0.00	85.16	-85.16
DATASET	64.94	-34.39	56.48	-95.25
SMILES-VAE	76.81 ± 0.25	-10.96 ± 0.71	71.01 ± 0.62	-90.94 ± 1.04
SELFIES-VAE	72.45 ± 3.79	-10.45 ± 3.83	72.05 ± 0.00	-87.82 ± 2.13
SMILES-LSTM-HC	59.64 ± 4.10	-31.03 ± 16.15	71.81 ± 1.56	-91.58 ± 2.14
SELFIES-LSTM-HC	63.17 ± 4.34	-21.02 ± 4.95	68.06 ± 5.74	-96.59 ± 4.59
REINVENT	68.38 ± 2.00	-24.35 ± 6.46	55.25 ± 5.88	-94.52 ± 1.20
GB-GA	56.04 ± 3.07	-41.39 ± 5.76	45.20 ± 6.78	-100.07 ± 1.35
JANUS	47.56 ± 2.19	-45.37 ± 7.90	39.22 ± 3.99	-97.14 ± 1.13

Conclusions and outlook



- Computer-aided molecular design promises to accelerate chemical discovery
- Reliable property predictors that generalize out-of-domain are essential
- Cheap physics-based predictors with limited number of parameters work well
- Alternative design algorithms work in-domain by, e.g., fragment approaches



Acknowledgements



University of Toronto



Alán Aspuru-Guzik



Gary Tom



Rajvi Rana



Cyrille Lavigne

Stanford University



Akshat Kumar Nigam



Anshul Kundaje



Swedish
Research
Council



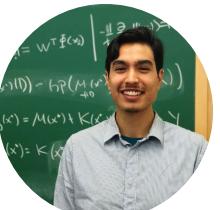
Robert Pollice



Luca Thiede



John Willes



Rodrigo Vargas

EPFL



Clémence
Corminboeuf



Rubén Laplaza



Terence Blaskovits

ETH



Luca Schaufelberger

