



BO for Guided Chemical Design

Austin Mroz

Imperial College London

a.mroz@imperial.ac.uk

tutorial notebooks



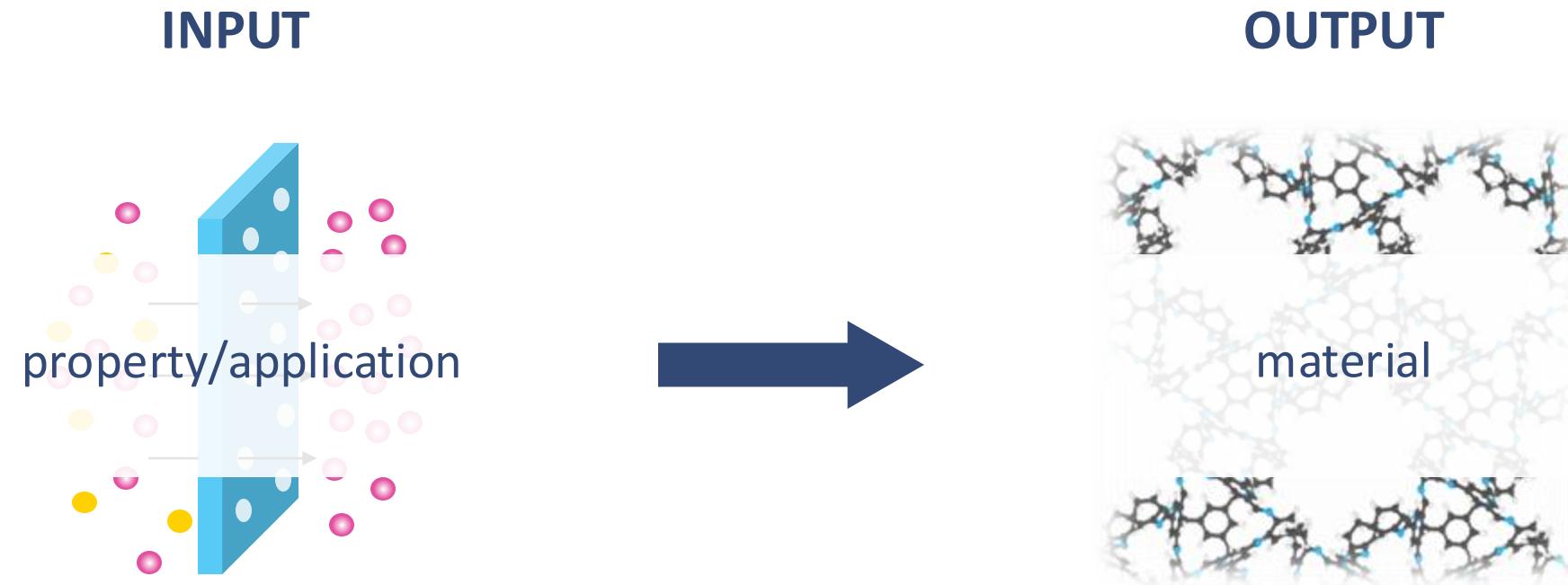
Agenda

time	subject	learning objectives
14:00-14:20	Why BO for chemistry	decision making in chemistry from OFAT & DOE to BO
14:20-14:30	Chemistry-specific considerations	surrogate model selection & chemical representations
14:30-14:40	Web-BO walkthrough	introduction to tool and notebooks we will be exploring
14:40-15:20	<i>interactive</i>	chemical representation comparison with SUMMIT csv
15:20-15:30	Web-BO: what's happening on the backend?	SOBO code walkthrough
15:30-16:00	<i>break</i>	

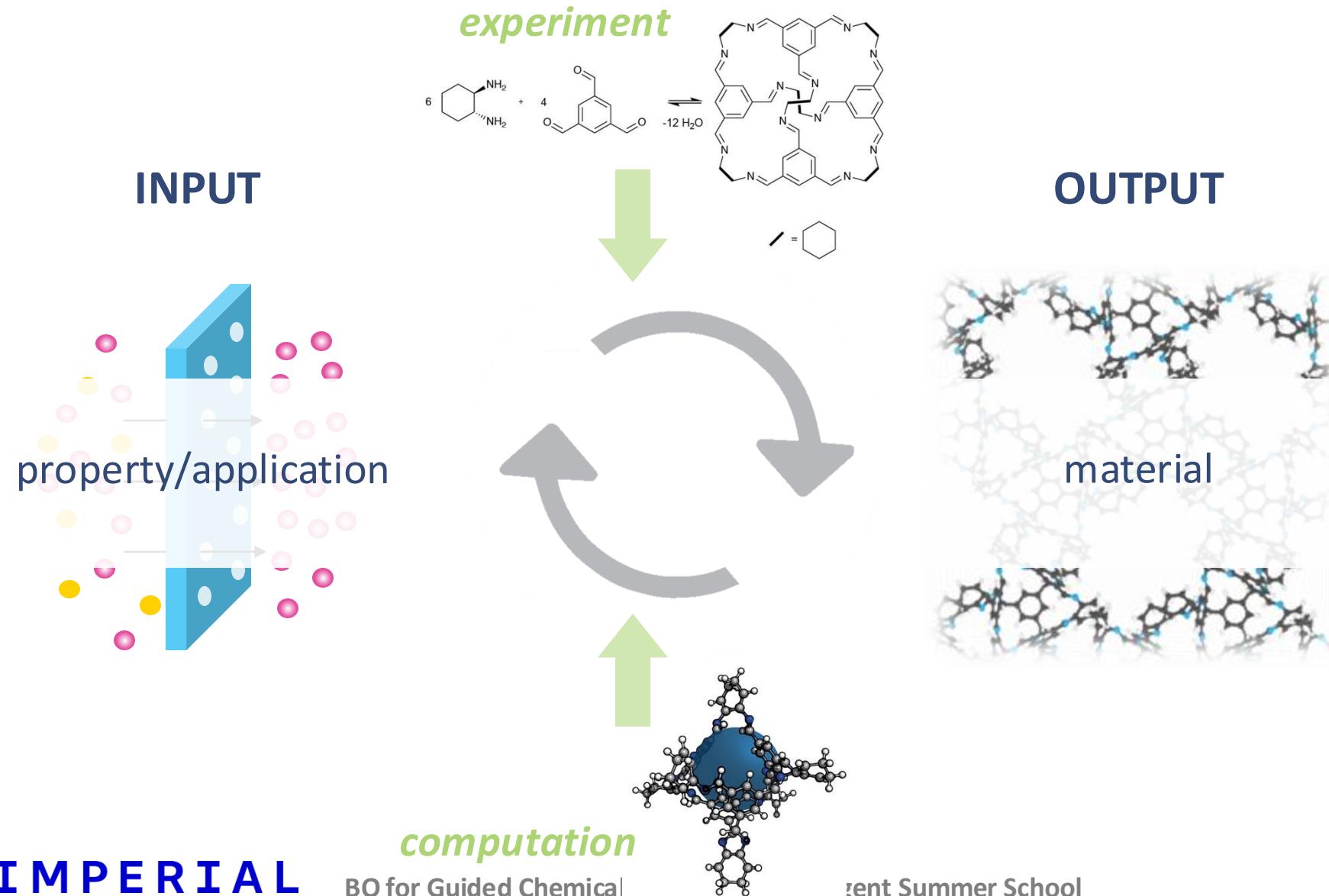
time	subject	learning objectives
16:00-16:15	Complex BO formulations in chemistry	MFBO, MOBO, MF-MO BO, etc.
16:15-17:00	<i>interactive</i>	code-your-own-adventure – <i>MOBO, MFBO, GPs for molecules</i>
17:00-17:20	Complex BO formulations in chemistry	Let's dig into the literature
17:20-17:30	wrap-up discussion	summary and additional resources

why BO for chemistry?

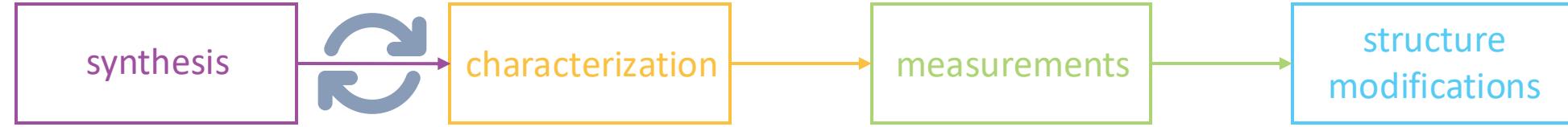
chemical design



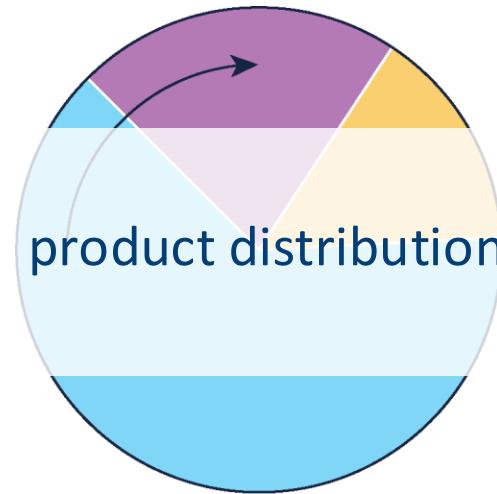
chemical design



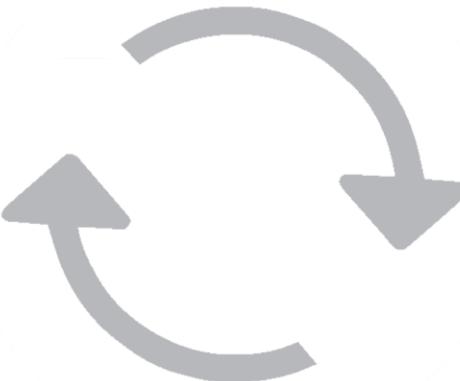
conventional *experimental* chemical discovery



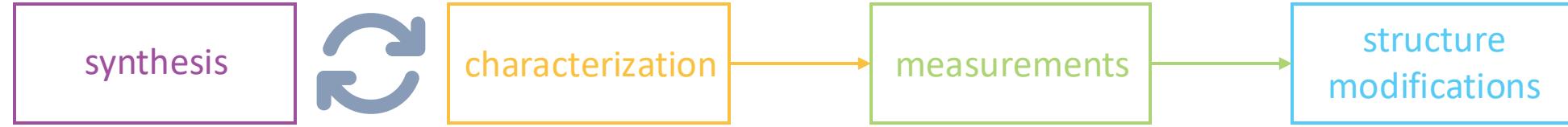
INPUT



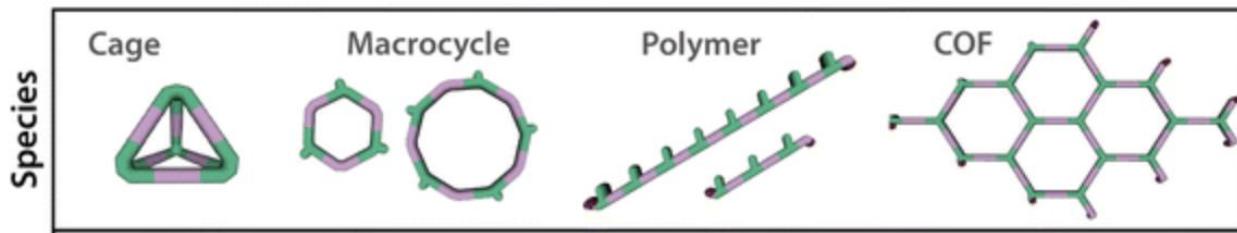
OUTPUT



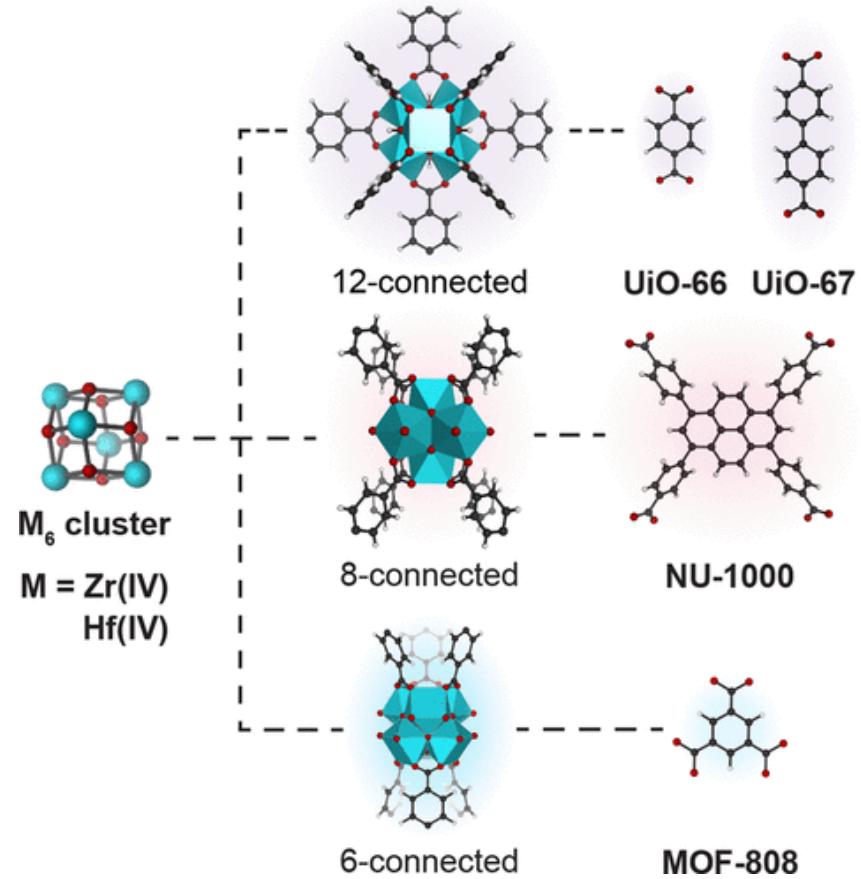
conventional *experimental* chemical discovery



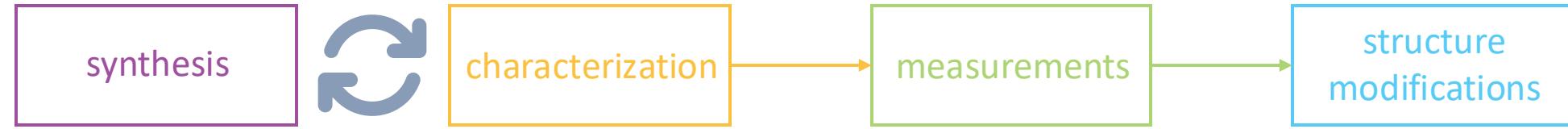
10^{60}
Scale



Chem. Sci., 2024, 15, 6331

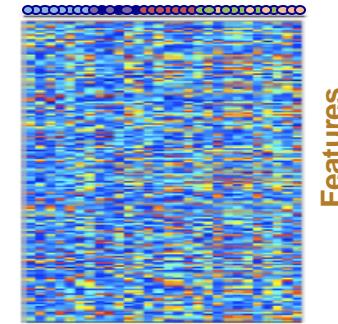


conventional *experimental* chemical discovery



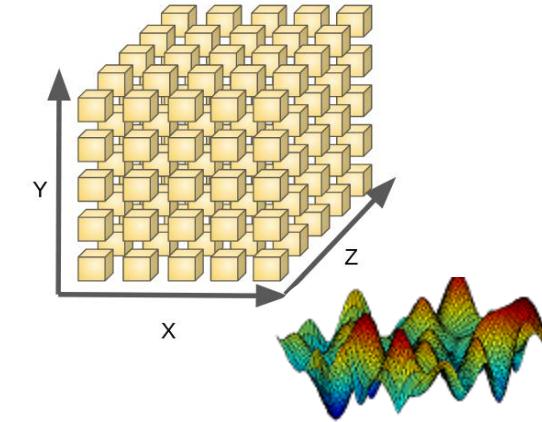
10^{60}

Scale



10^6

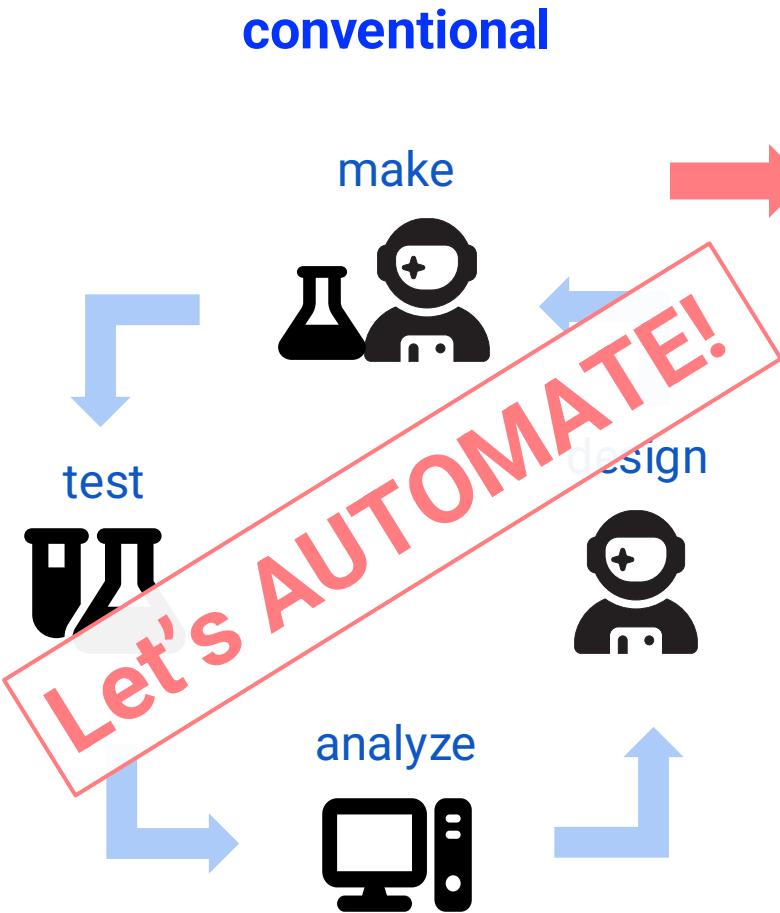
High-throughput
technologies



“Curse of dimensionality”

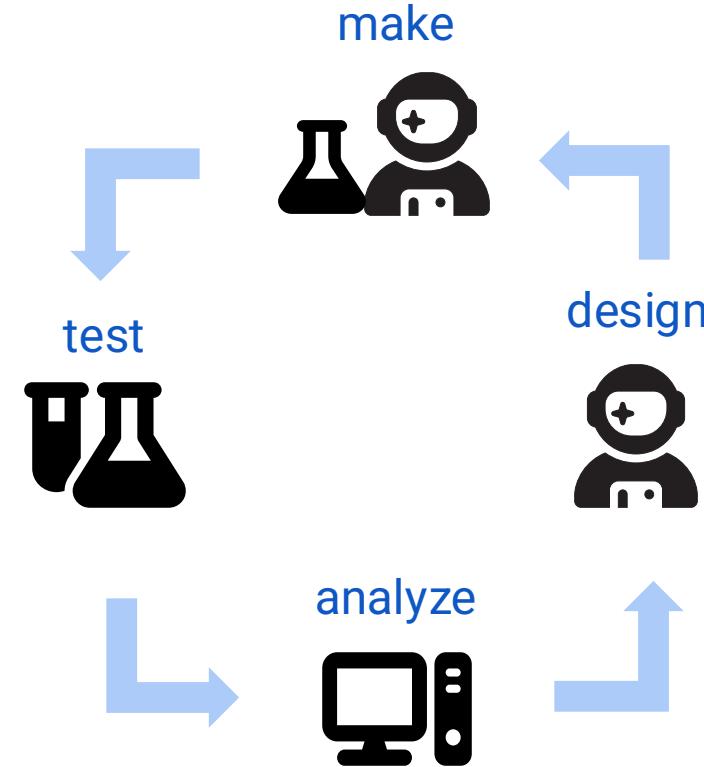
but given these near-infinite landscapes and limited resources,
how do we efficiently move towards desired maxima?

Conventional design in chemistry

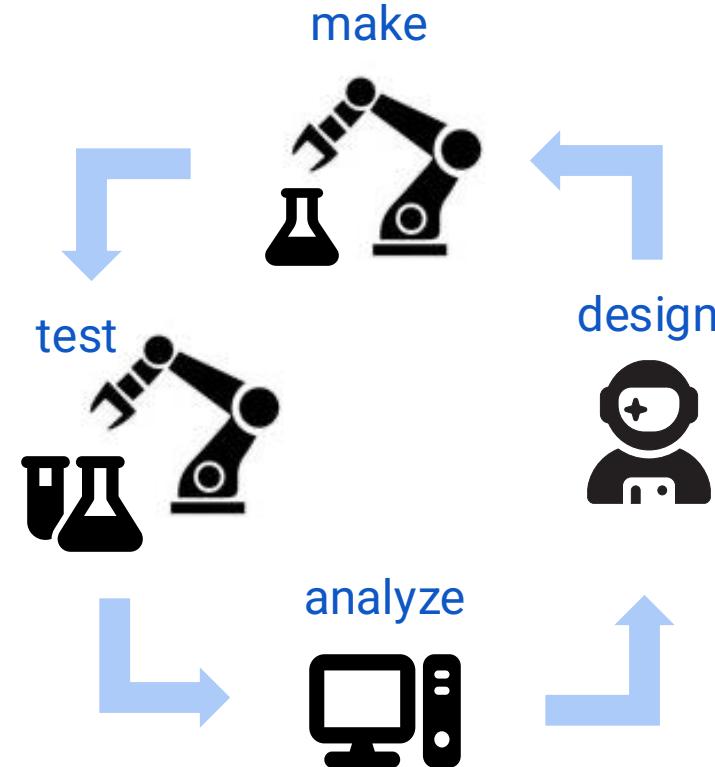


Design-make-test-analyze

conventional



automated



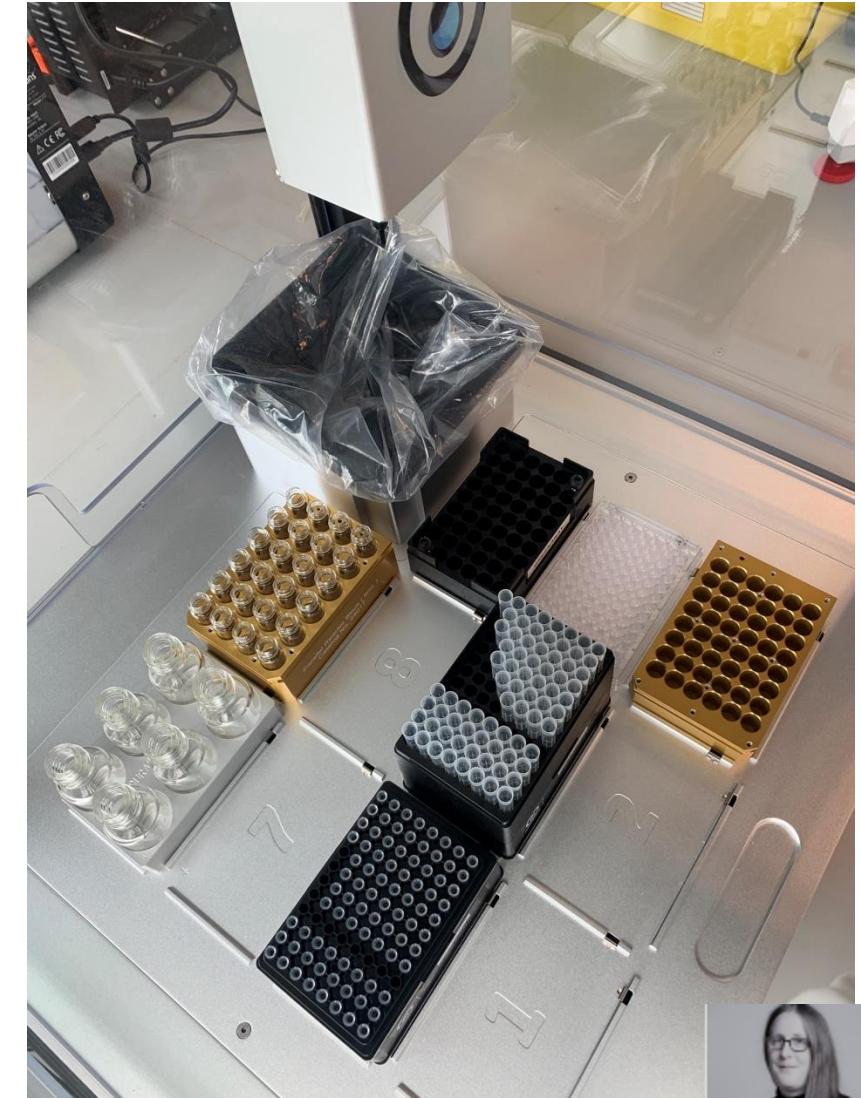
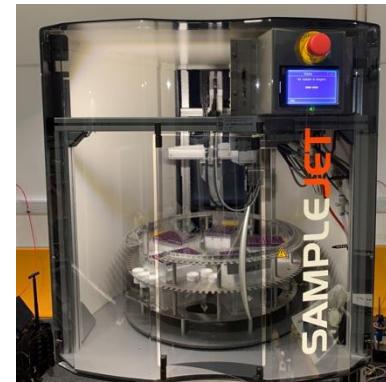
where are the bottlenecks?

1. manual synthesis, characterization and target property measurements are time and resource intensive
2. chemical intuition is largely limited to *small* changes to either experimental conditions or structure

how does this help us?

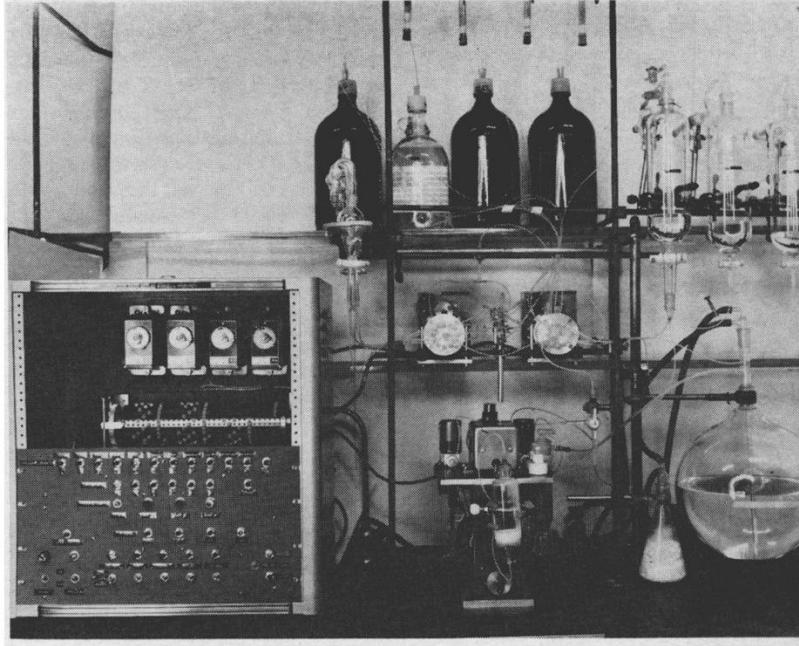
1. chemistry research 24/7!
2. improved reproducibility!
3. large batches!

Automation in chemical sciences



Dr. Becky
Greenaway

Automation in chemistry is not new

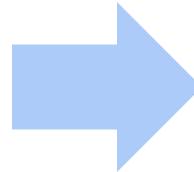


Science, 1965, 150, 178

Automated Synthesis of Peptides

Solid-phase peptide synthesis, a simple and rapid synthetic method, has now been automated.

R. B. Merrifield



Science, 2020, 368, 980

PEPTIDE SYNTHESIS

Synthesis of proteins by automated flow chemistry

N. Hartrampf^{1*}, A. Saebi^{1†}, M. Poskus^{1†}, Z. P. Gates^{1‡}, A. J. Callahan¹, A. E. Cowfer¹, S. Hanna¹, S. Antilla^{1,2}, C. K. Schissel¹, A. J. Quartararo¹, X. Ye¹, A. J. Mijalis^{1§}, M. D. Simon^{1¶}, A. Loas¹, S. Liu^{1#}, C. Jessen³, T. E. Nielsen³, B. L. Pentelute^{1**}

Levels of automation in chemistry

manual

single **step** at a time
carried out by hand

standalone

single **process** at a time
automated on one machine

integrated

multiple **processes**
automated workflow
no human intervention

autonomous

fully integrated workflow
closed-loop
algorithmic feedback

Levels of automation in chemistry

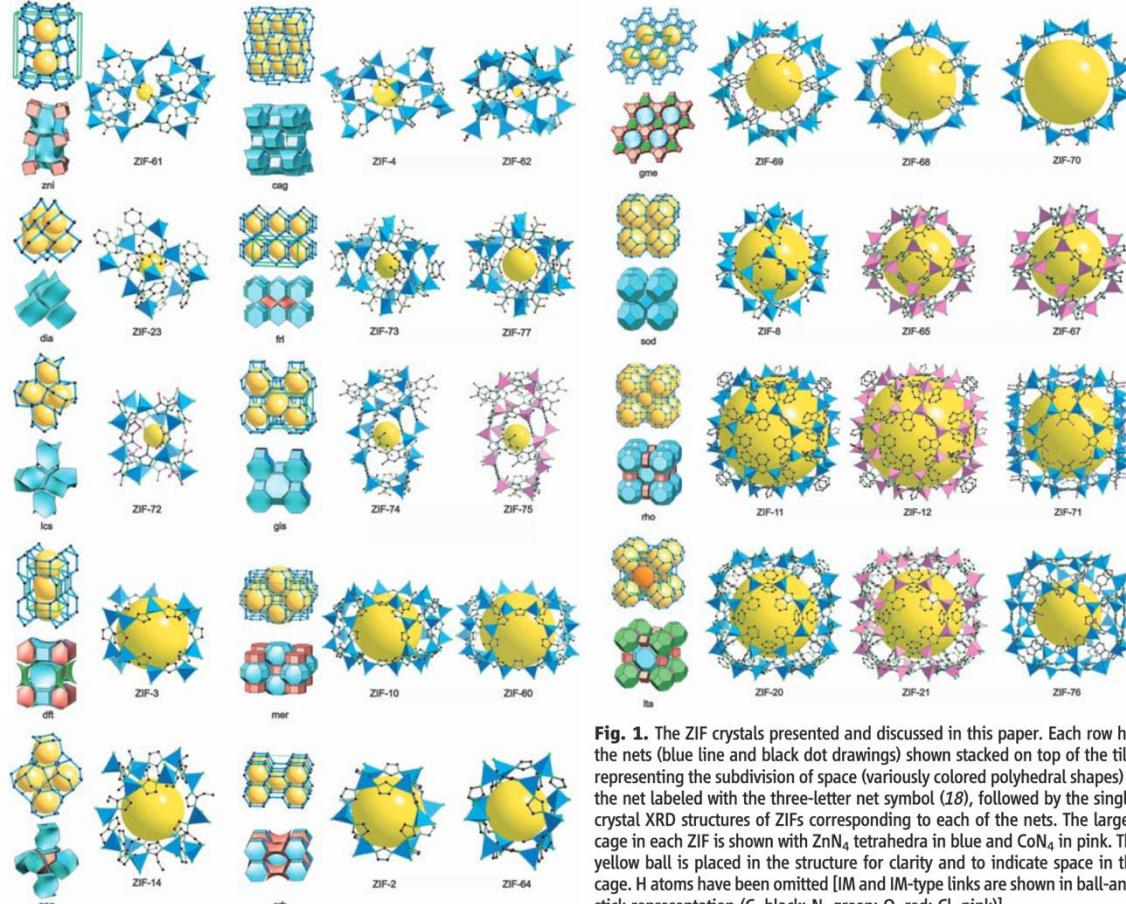
TIMESCA
LE

manual

standalone

integrated

autonomous



Zeolitic imidazole framework (ZIF) discovery

aim – high throughput synthesis of ZIFs

how – 96-well glass plate to screen 9600 microreactions

so? – 10 new ZIFs, featuring several high performing candidates



Science, 2008, 319, 939

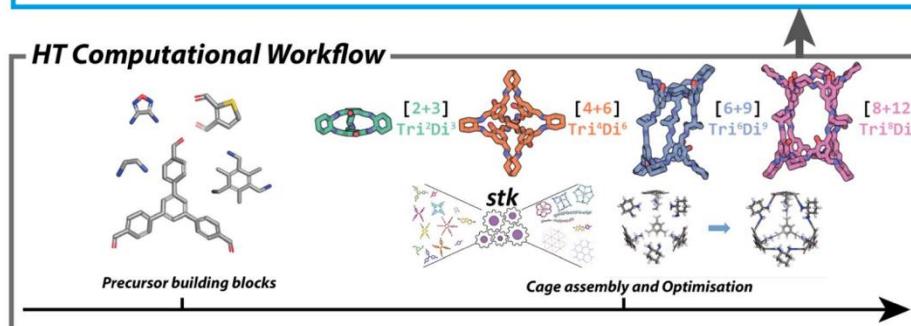
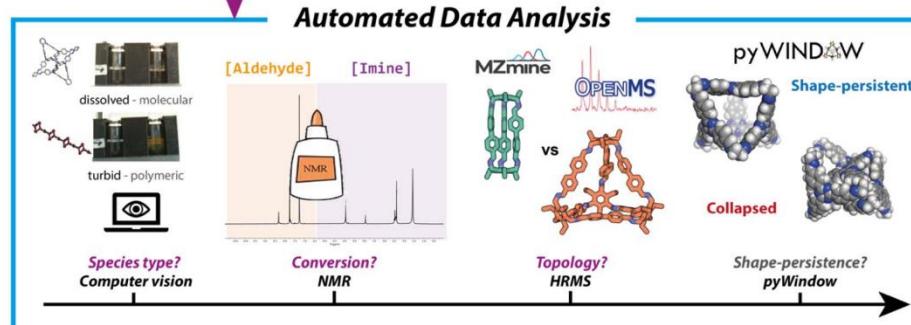
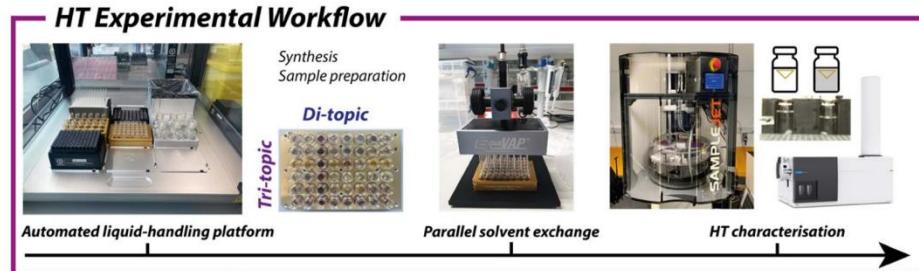
Levels of automation in chemistry

manual

standalone

integrated

autonomous



porous organic cage (POC) discovery

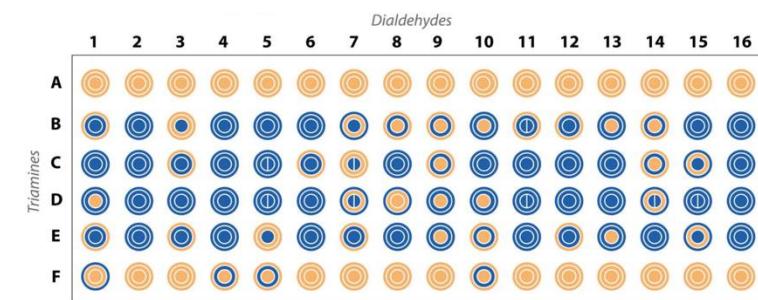
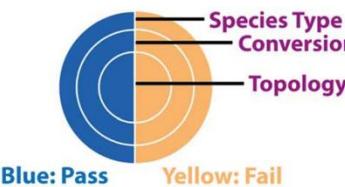
aim – high throughput synthesis of POCs

how – combined experimental and computational workflows

- automated data analysis
- computer vision for on-the-fly reaction feedback

so? – 225 POCs identified

- 350-fold decrease in time required for data analysis



Chem. Sci., 2024, 15, 6331

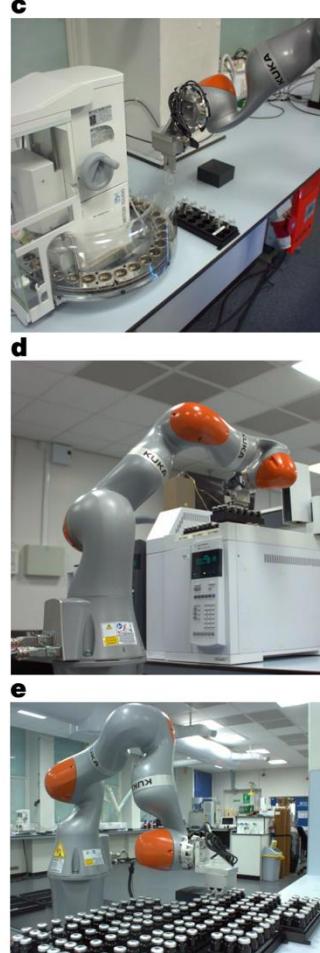
Levels of automation in chemistry

manual

standalone

integrated

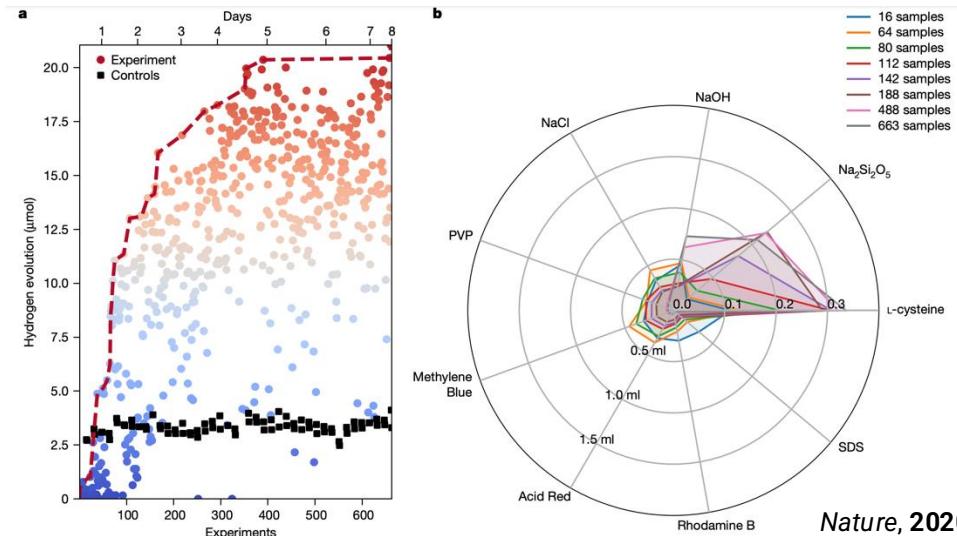
autonomous



mobile robotic chemist

aim – automating the researcher instead of the instruments

how – 8 days
688 experiments
10-variable parameter space



Nature, 2020, 583, 237

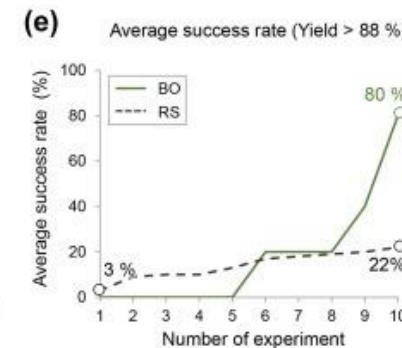
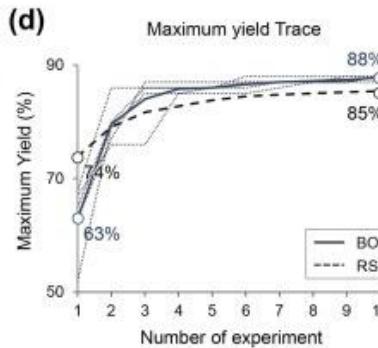
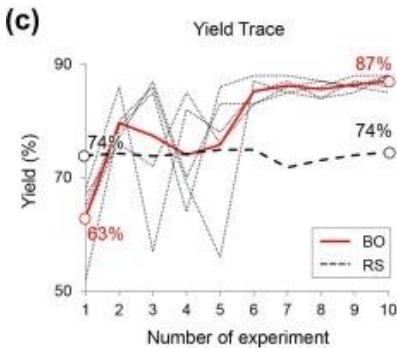
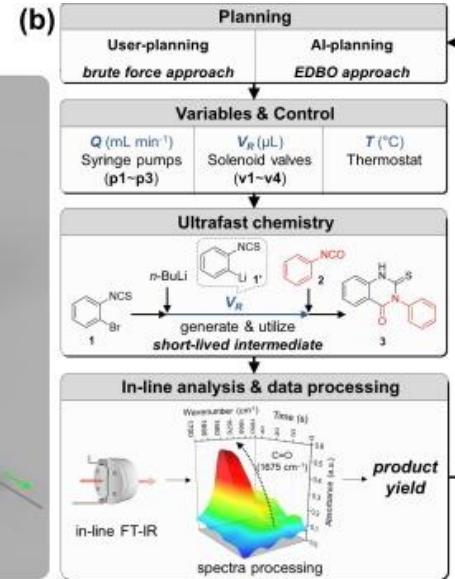
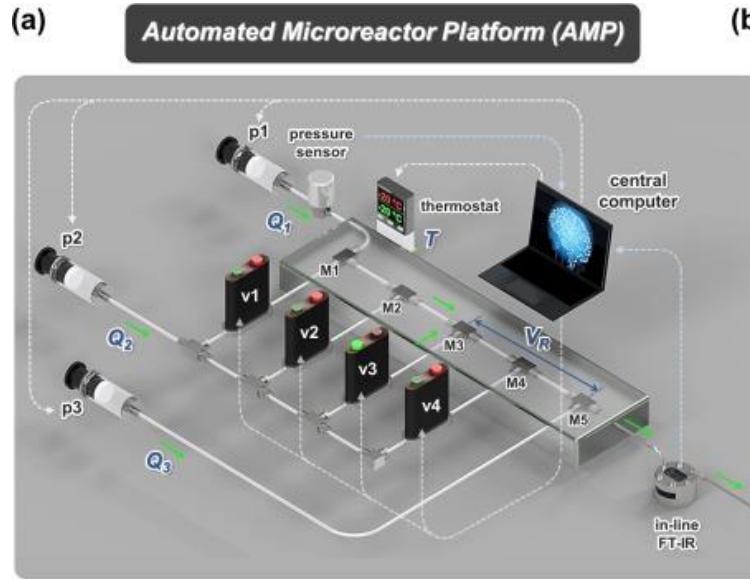
Levels of automation in chemistry

manual

standalone

integrated

autonomous



FLOW CHEMISTRY

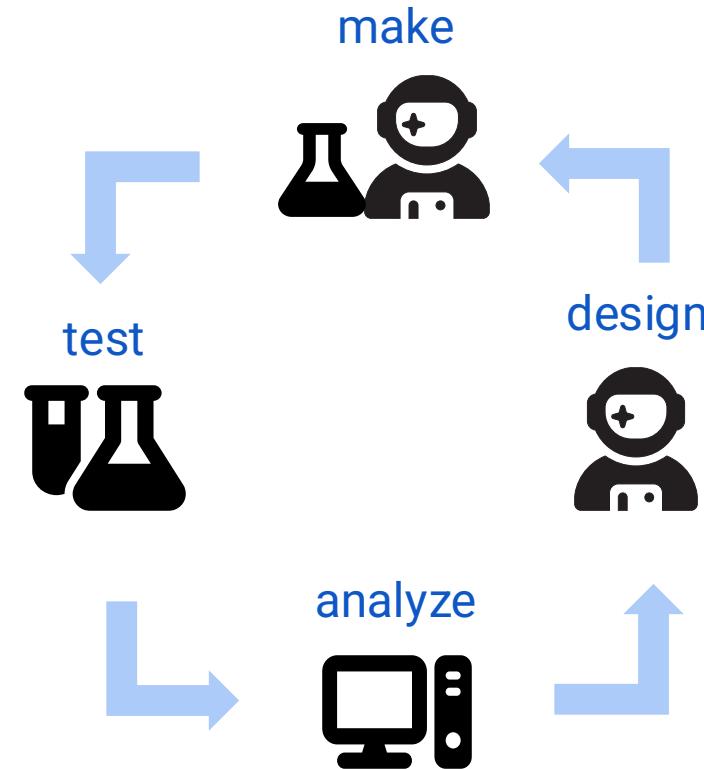
aim – ultrafast synthesis of biologically active compounds

how – exploits user-planned and AI-planned modes to accelerate synthesis

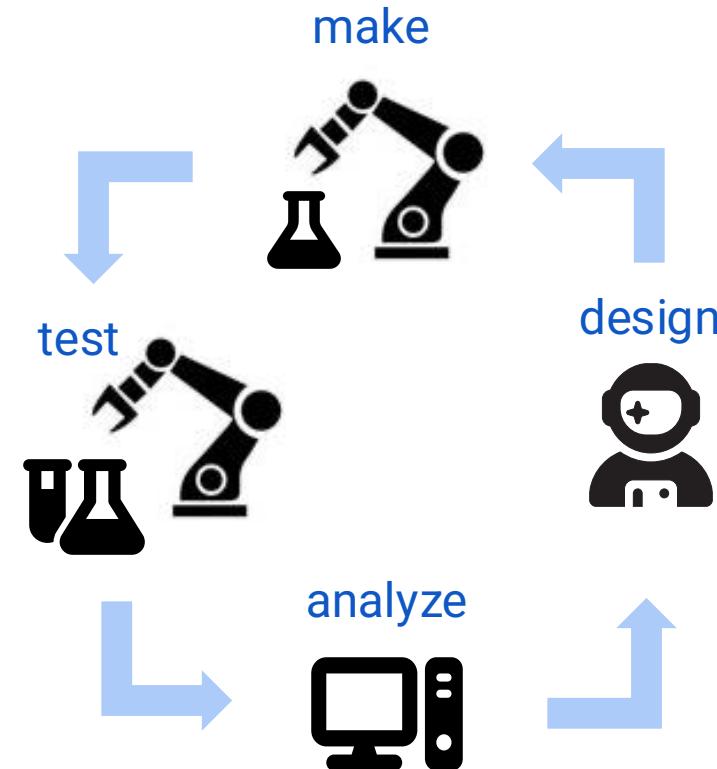
so – efficiency in chemical synthesis via improved precision and control offered by AMP platform scalability and flexibility

Design-make-test-analyze

conventional



automated



it's not all rainbows...

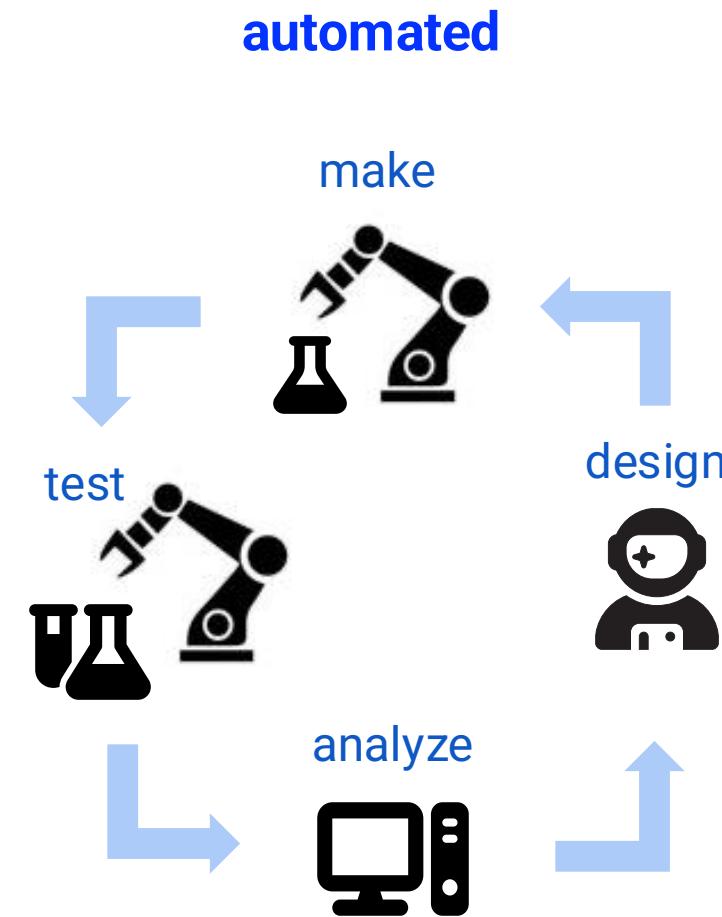
automation isn't automatic

Chem. Sci., 2021, 12, 15473

where are the bottlenecks?

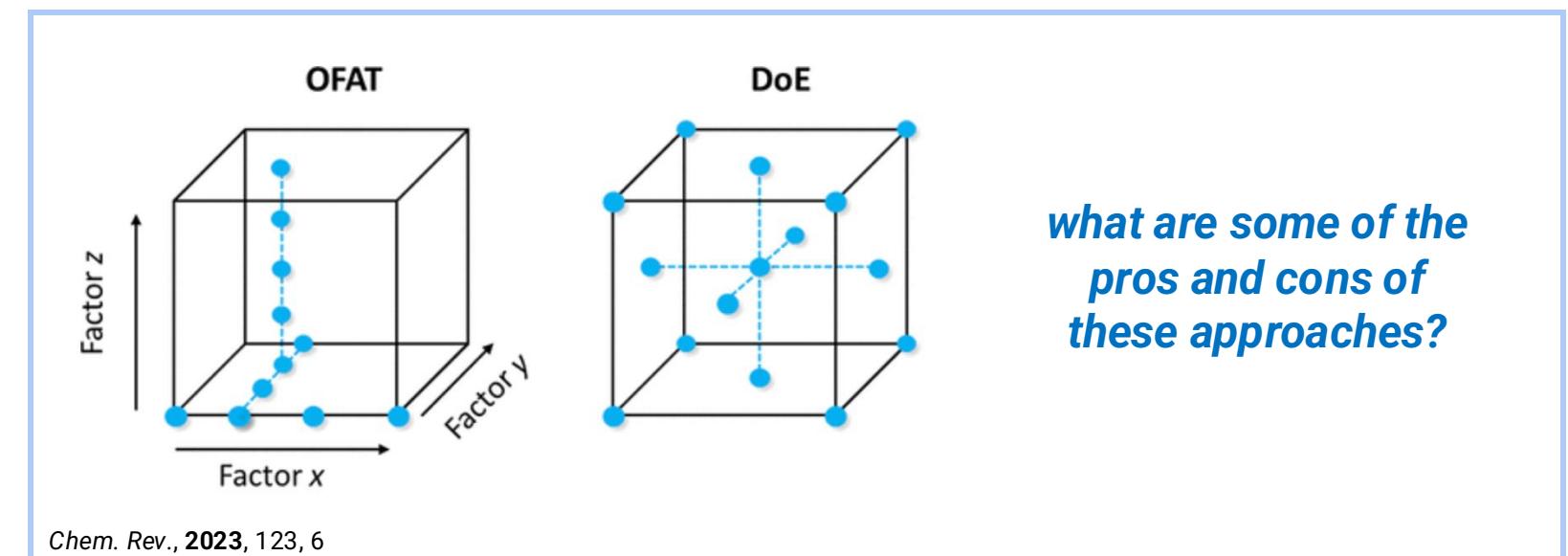
1. manual synthesis, characterization and target property measurements are time and resource intensive
2. **chemical intuition is largely limited to *small changes to either experimental conditions or structure***

Conventional design in chemistry

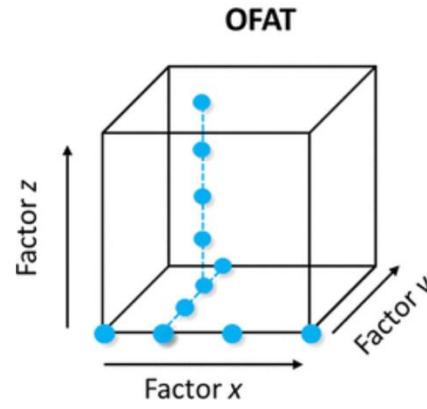


where are the bottlenecks?

1. manual synthesis, characterization and target property measurements are time and resource intensive
2. **chemical intuition is largely limited to small changes to either experimental conditions or structure**



OFAT & DOE – pros and cons

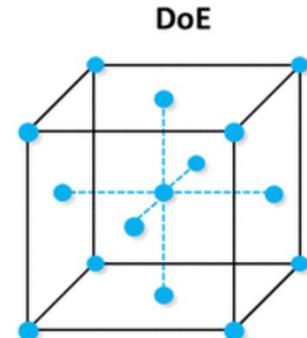


pros

widely taught
straightforward

cons

limited experimental space coverage
may miss optimal solution
fails to identify interactions
inefficient use of resources



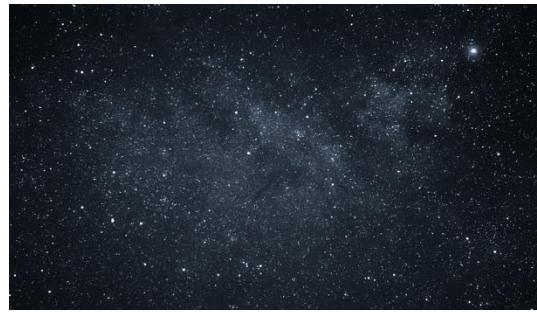
systematic – thorough coverage of experimental space

efficient – minimal resources required

minimum entry of experiments (~10)

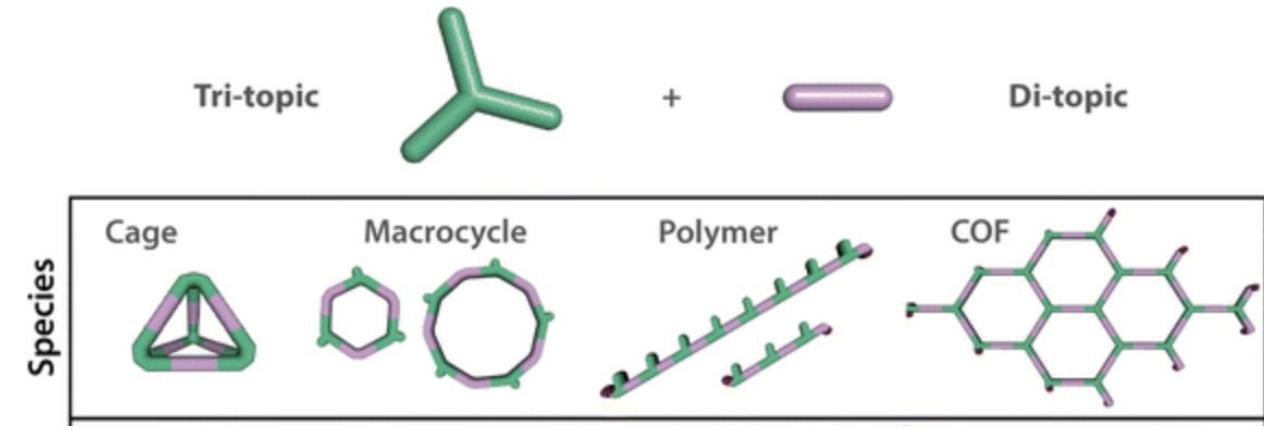
may need to perform experiments that researcher expects to fail

consider chemical space...

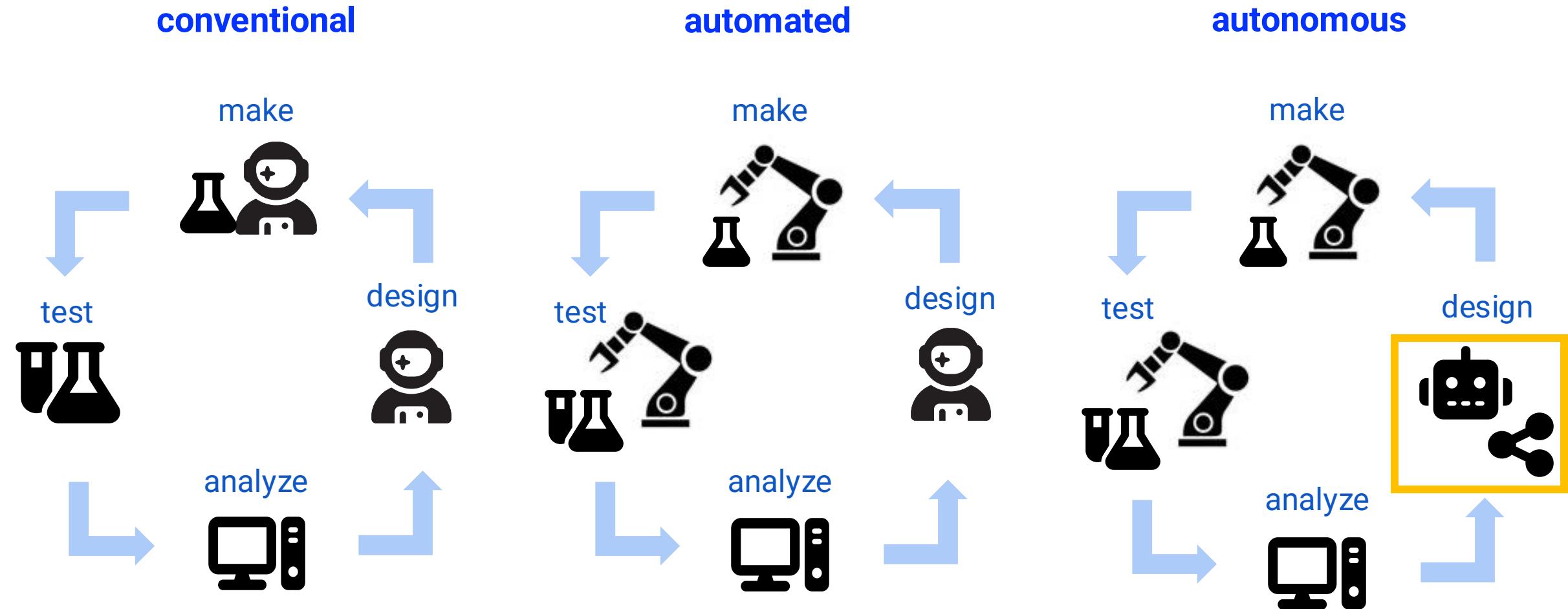


10^{60}

Scale

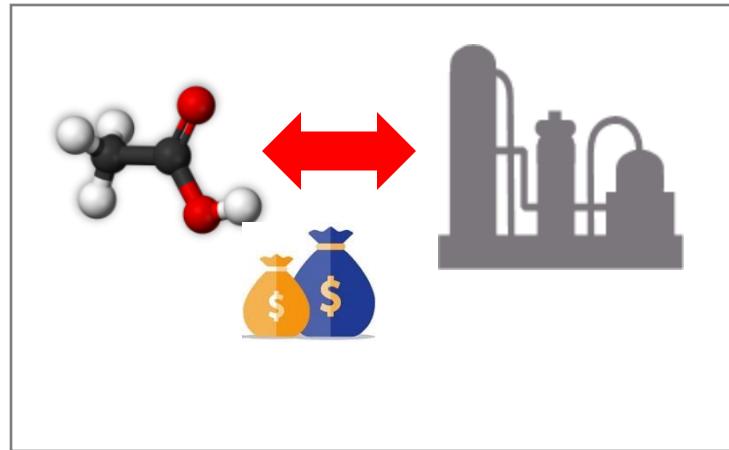


Design-make-test-analyze



Challenges

Challenge 1



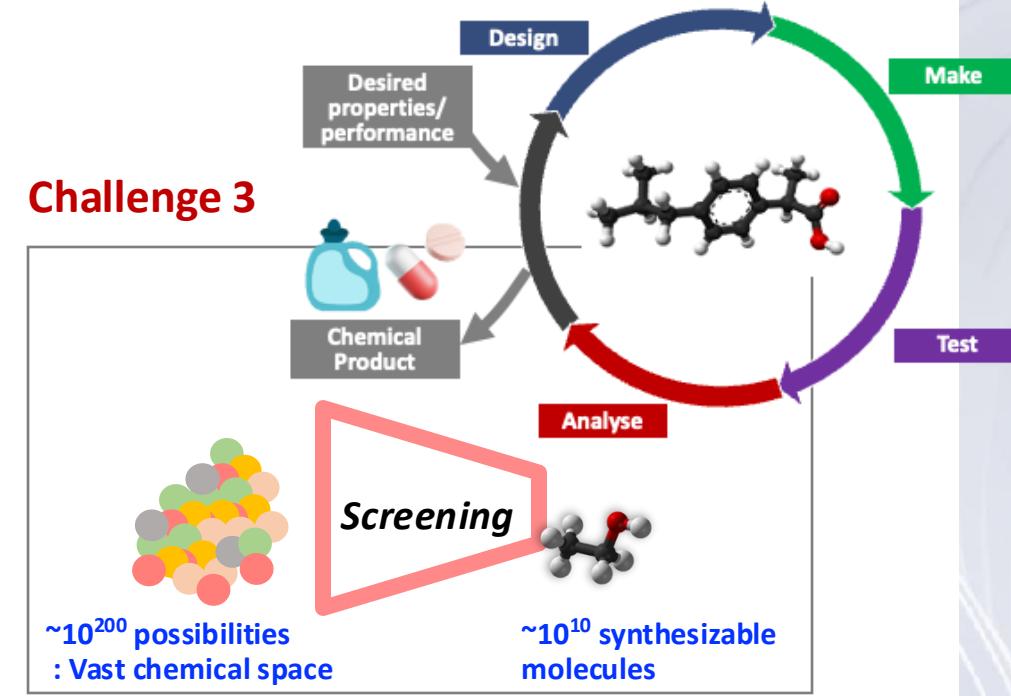
- Overall performance: molecular + system-level decisions
→ high numerical complexity

Challenge 2



- Unknown physical mechanism
- Uncertainty in property prediction models

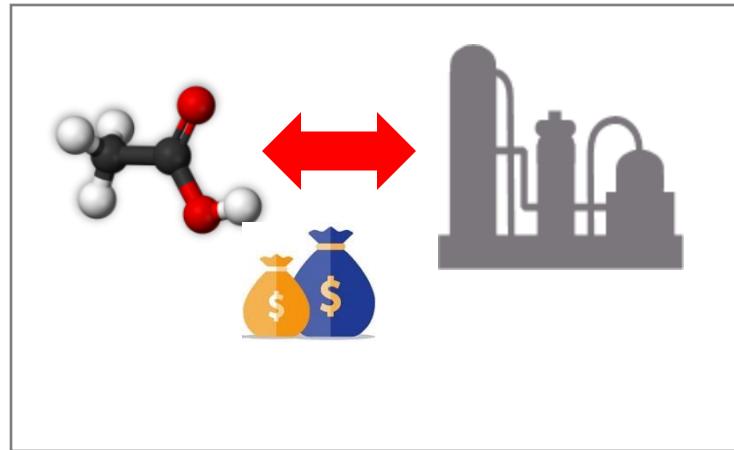
Challenge 3



- Developing a novel chemical involves more than just design
- Synthesizability of molecules
→ Non-smooth & nonlinear function

Challenges

Challenge 1



- Overall performance: molecular + system-level decisions
→ high numerical complexity

- Developing new techniques and tools for the design of better products and processes

Mechanistic Modelling



Statistical modelling
Machine learning techniques

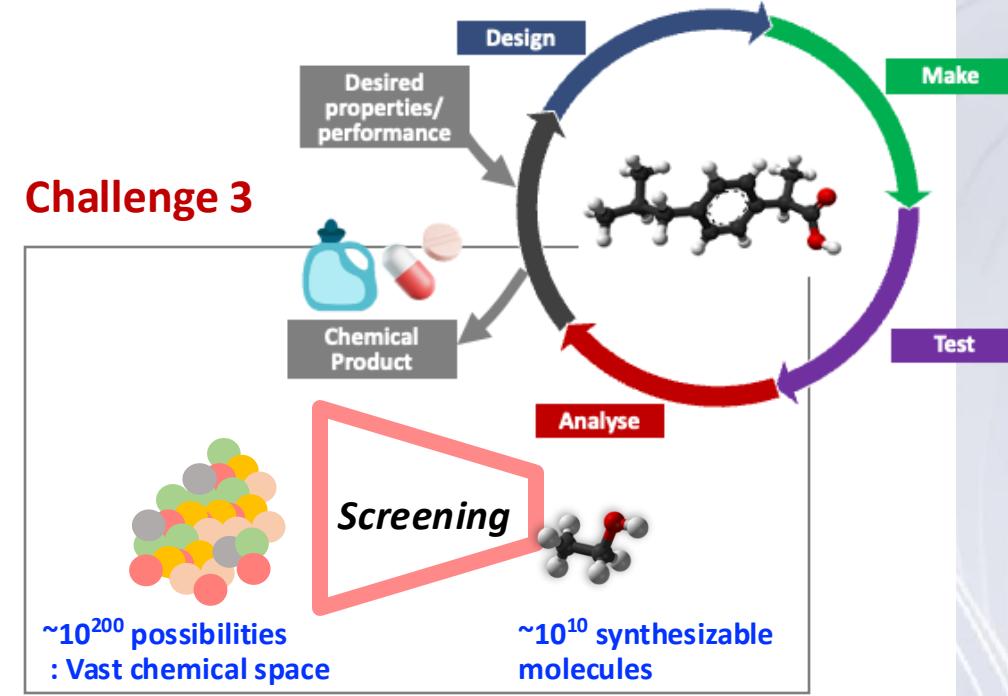
Challenge 2



- Unknown physical mechanism
- Uncertainty in property prediction models



Challenge 3



- Developing a novel chemical involves more than just design
- Synthesizability of molecules
→ Non-smooth & nonlinear function

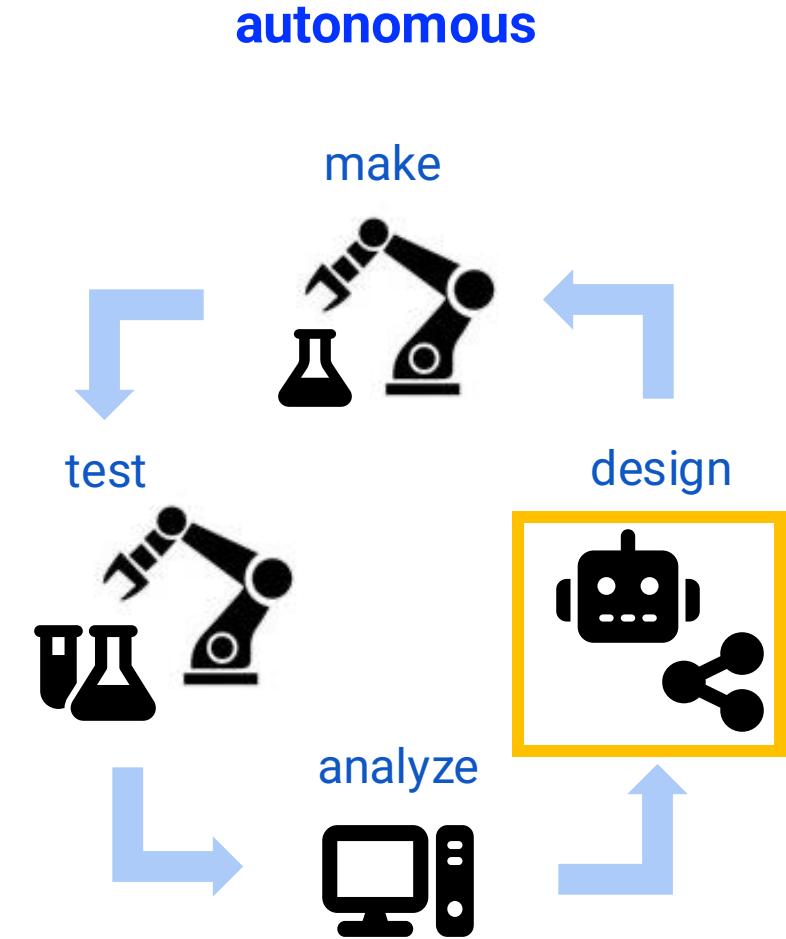
DataDrivenDesign-make-test-analyze

main considerations

1. we don't know what our objective function looks like
2. we are often operating under budgetary and resource constraints
3. we often don't have a lot of starting data, if any at all

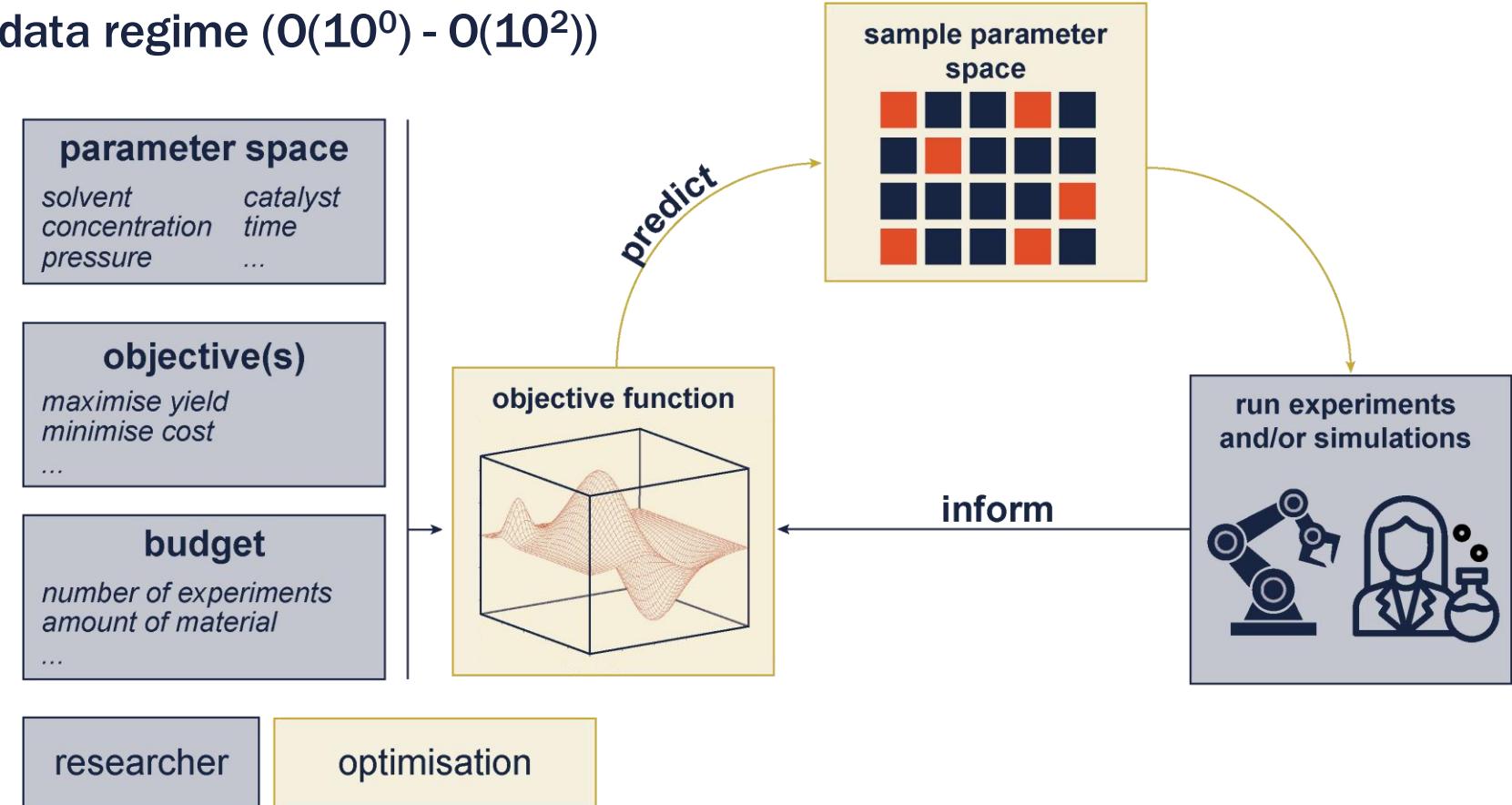
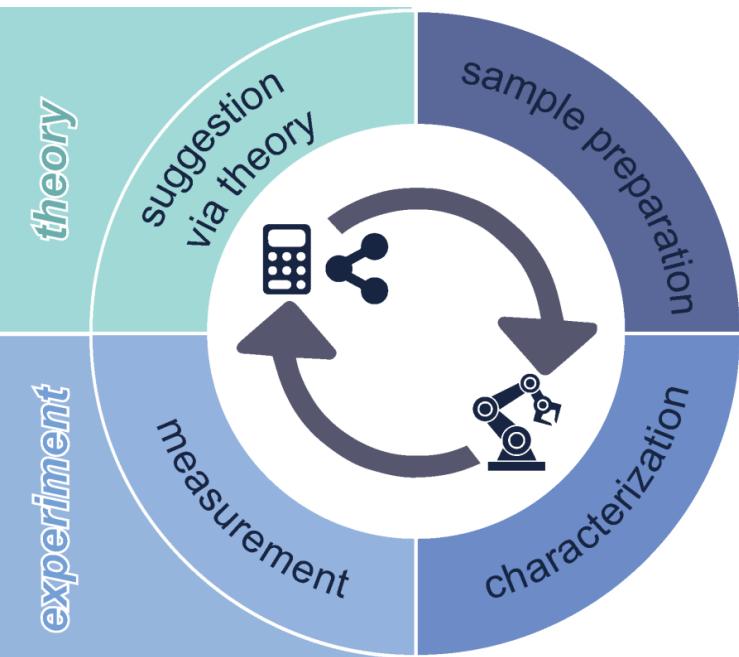
what does our solution need?

1. we need to be able to sample
2. resource- and budget-aware
3. need to work well in the sparse/low data regime

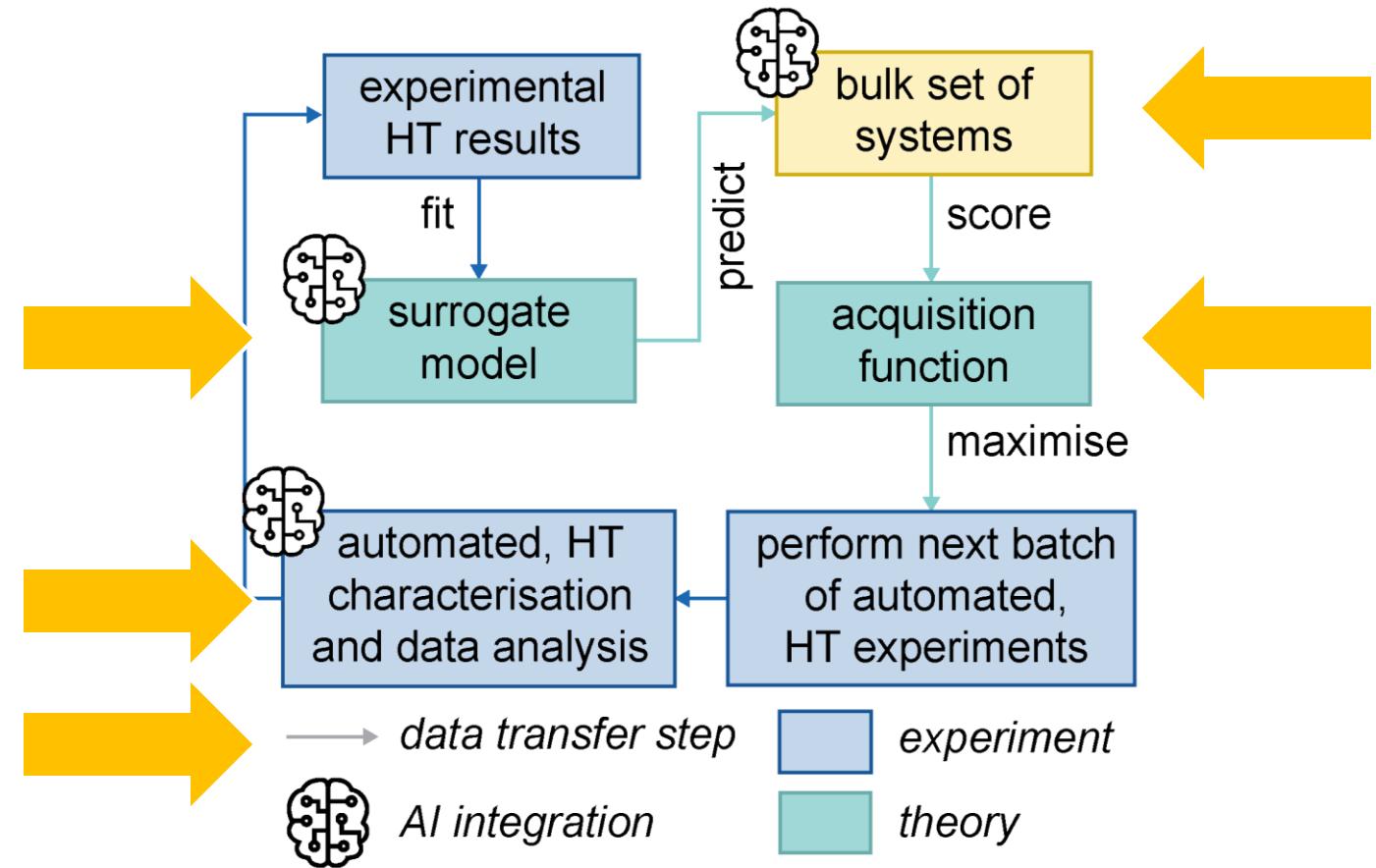
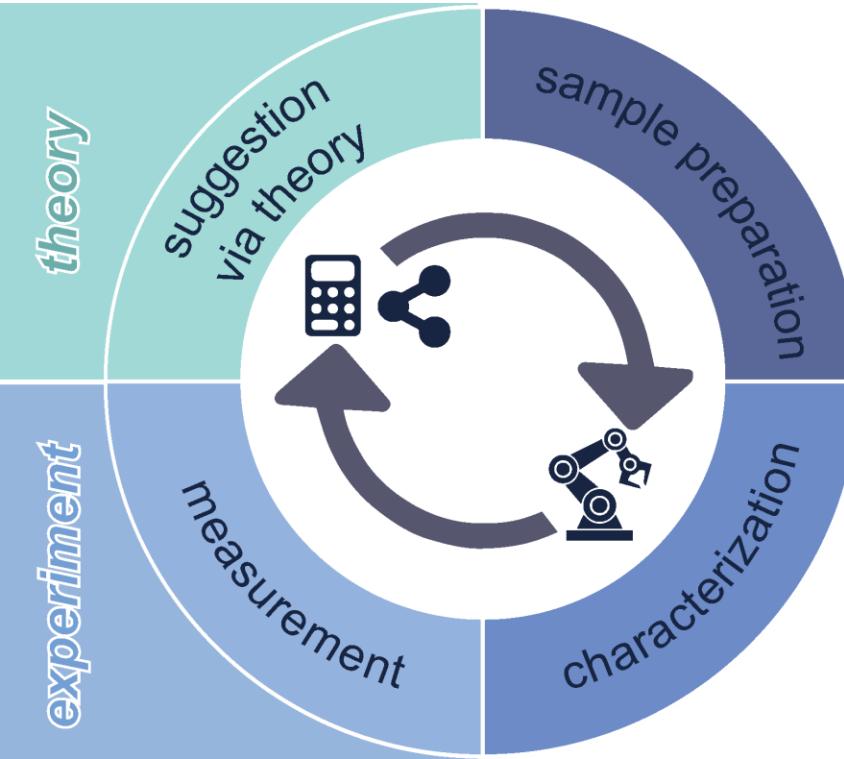


Bayesian optimization for chemistry

- ✓ resource- and budget-aware
- ✓ balances exploitation and exploration
- ✓ sparse data regime ($O(10^0) - O(10^2)$)



Bayesian optimization for chemistry



chemistry-specific considerations

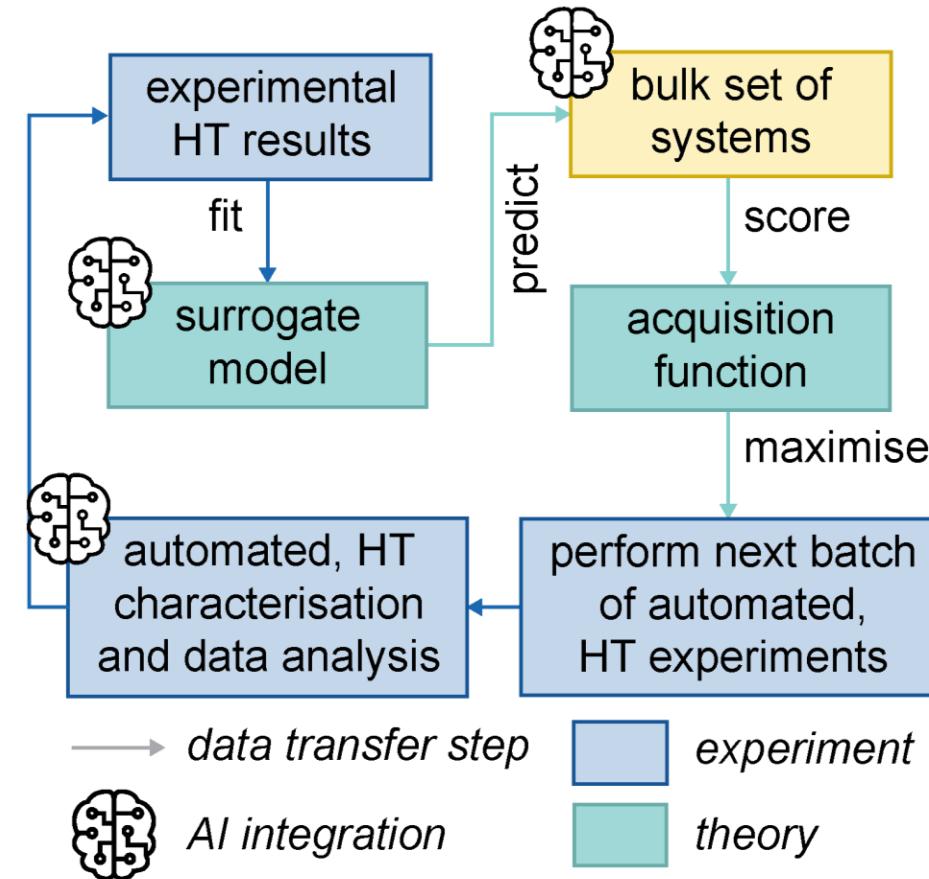
chemistry-specific considerations and problem formulations

objective

- maximise, minimise, combination?
- new material
- improved performance
- experimental conditions

information streams

- experiment, theory, combination?
- timescales
- batch?



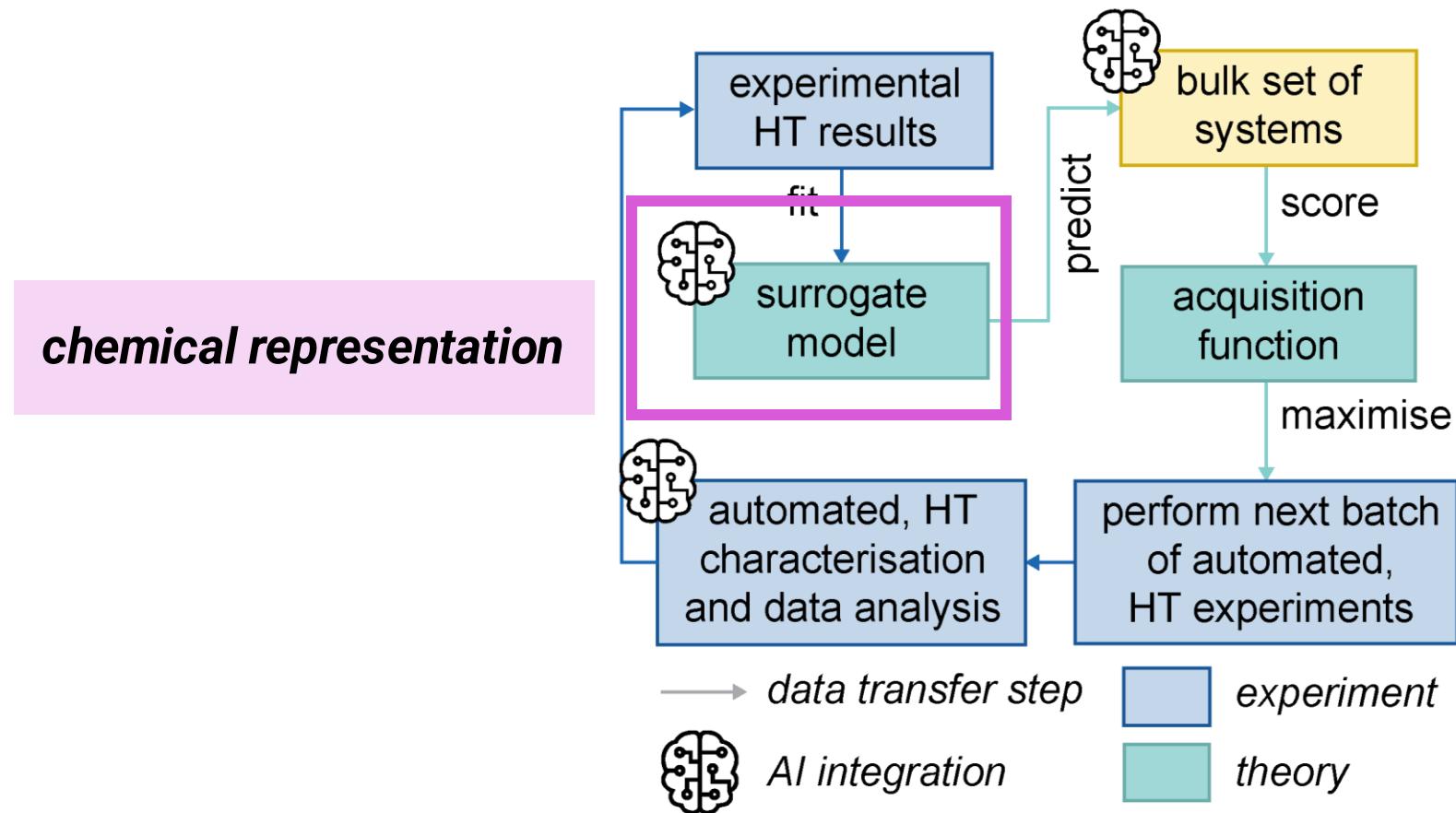
design space

- integer, continuous, categorical, chemical?
- chemical representation / encoding

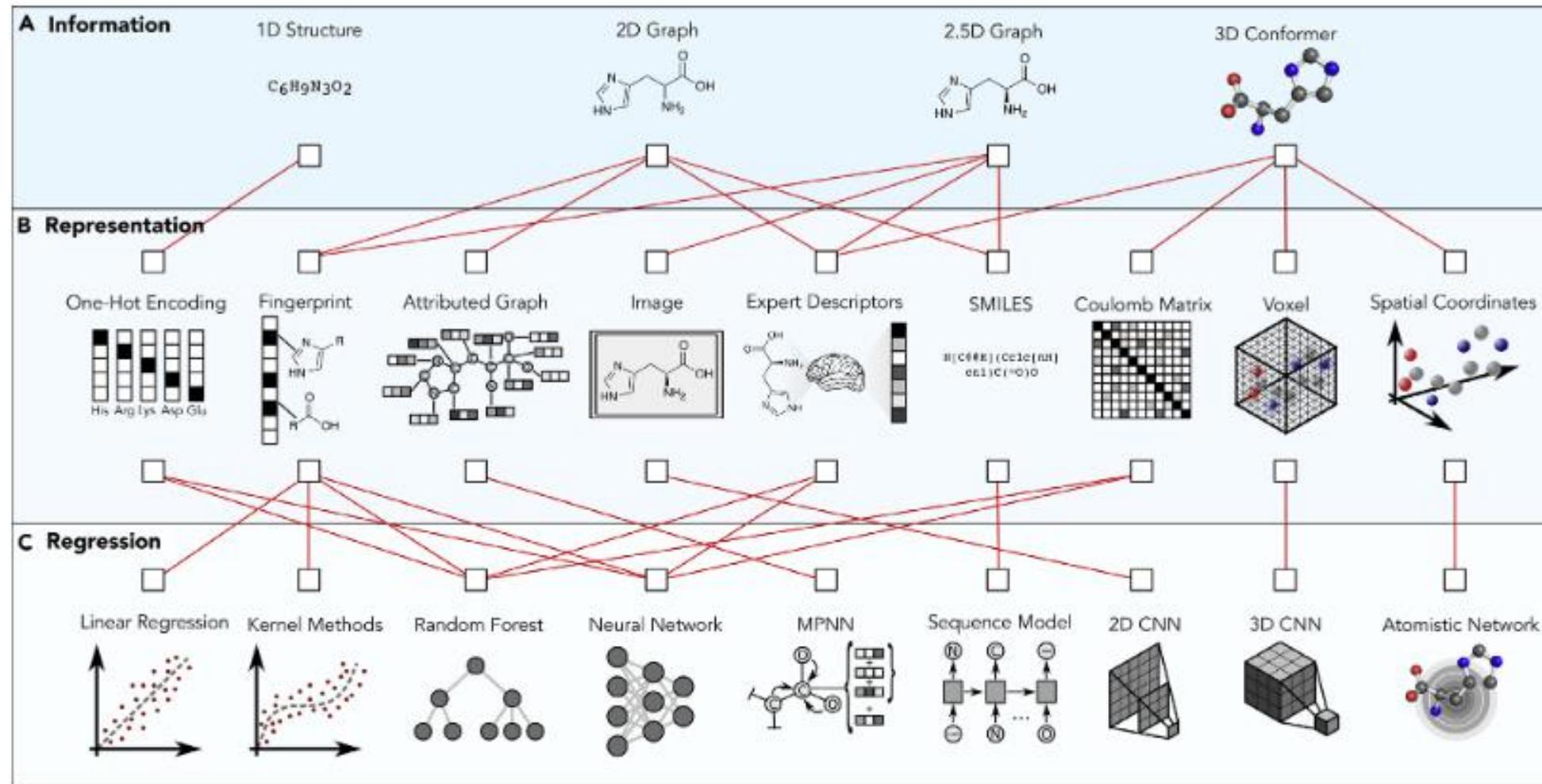
stopping criteria

- precursor amounts
- number of experiments
- total time

chemical representations & surrogate model selection



Overview of chemical representations



Chem., 2020, 6, 1204

One-Hot Encoding

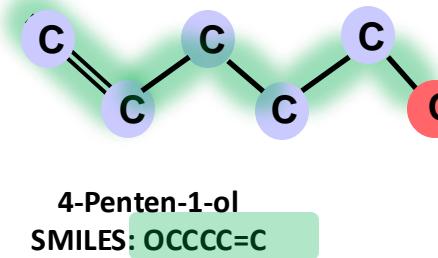
- The conversion of categorical information into a format that may be fed into machine learning.

Example 1: Simple

A grid showing fruit categories and their one-hot encoding. The categories are represented by icons: lemon, cherries, watermelon, grapes, and another lemon. The first row shows the icons, with the second icon (cherries) highlighted by a red box. Below the icons is a 4x5 grid of binary values:

1	0	1	0	0
0	0	0	0	0
0	0	0	0	1
0	1	0	0	0

Example 2: molecule



A large matrix illustrating one-hot encoding for molecules. The columns represent different elements: C, O, N, O, ..., =, C, The rows are labeled with element symbols: C, N, O, ..., =. The matrix entries are binary values (0 or 1), indicating the presence or absence of each element in a molecule. For example, the first row (C) has a 1 in the first column and 0s elsewhere, while the last row (=) has a 1 in the eighth column and 0s elsewhere.

0	C	C	C	C	=	C	...	
C	0	1	1	1	1	0	1	0
N	0	0	0	0	0	0	0	0
O	1	0	0	0	0	0	0	0
...	0	0	0	0	0	0	0	0
=	0	0	0	0	0	1	0	0

pros

- resulting representation prevents ordinal misinterpretation & prevents ranking issues
- works best for low-dimensional classes
 - (e.g. small number of categories)

cons

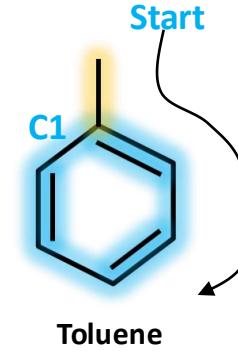
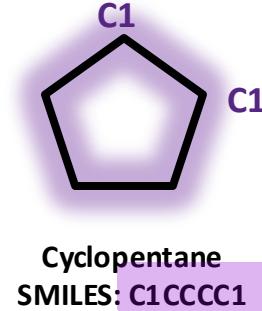
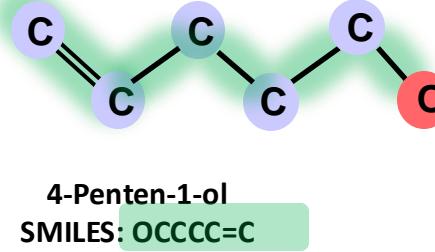
- curse of dimensionality
- may increase data sparsity

what are some of the pros & cons?

Text-based (or String-based) Representation

□ SMILES

- Simplified Molecular-Input Line-Entry system
- It is a line notation for describing the structure of chemical species
- A series of characters that represents atoms and bonds.
→ Normally, hydrogen atoms are not explicitly shown.



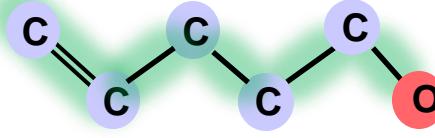
- Some vocab:
 - = : double bond, #: triple bond
 - () : branches
 - Digits: ring

Do it yourself !

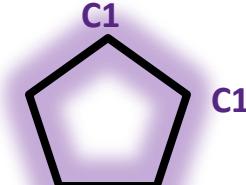
Text-based (or String-based) Representation

□ SMILES

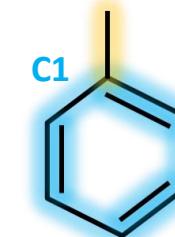
- Simplified Molecular-Input Line-Entry system
- It is a line notation for describing the structure of chemical species
- A series of characters that represents atoms and bonds.
→ Normally, hydrogen atoms are not explicitly shown.



4-Penten-1-ol
SMILES: OCCCC=C



Cyclopentane
SMILES: C1CCCC1



Toluene
SMILES: CC1=CC=CC=C1 (Pubchem)
Cc1ccccc1
c1(C)ccccc1

- Some vocab:
 - = : double bond, #: triple bond
 - () : branches
 - Digits: ring

What can you notice here?

Text-based (or String-based) Representation

□ SMILES

- It is necessary to select a “unique SMILES” for a molecule
→ “Canonicalization”

similar molecules can have very different
SMILES representation → fail to capture
structural similarity!

Canonical SMILES

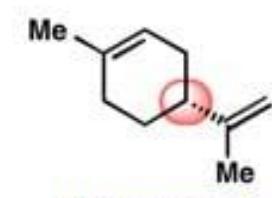
PubChem Cyclopentane (Compound)

2.1.3 InChIKey

RGSFGYAAUTVSQA-UHFFFAOYSA-N

Computed by InChI 1.0.6 (PubChem release 2021.10.14)

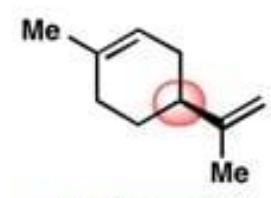
▶ PubChem



(+)-limonene



orange smell



(-)-limonene



lemon smell

2.1.4 Canonical SMILES

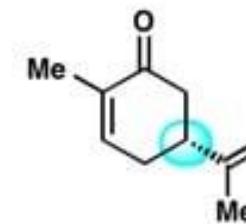
C1CCCC1

Computed by OEChem 2.3.0 (PubChem release 2021.10.14)

▶ PubChem

2.2 Molecular Formula

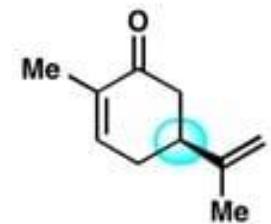
C₁₀H₁₆



(-)-carvone



spearmint smell



(+)-carvone

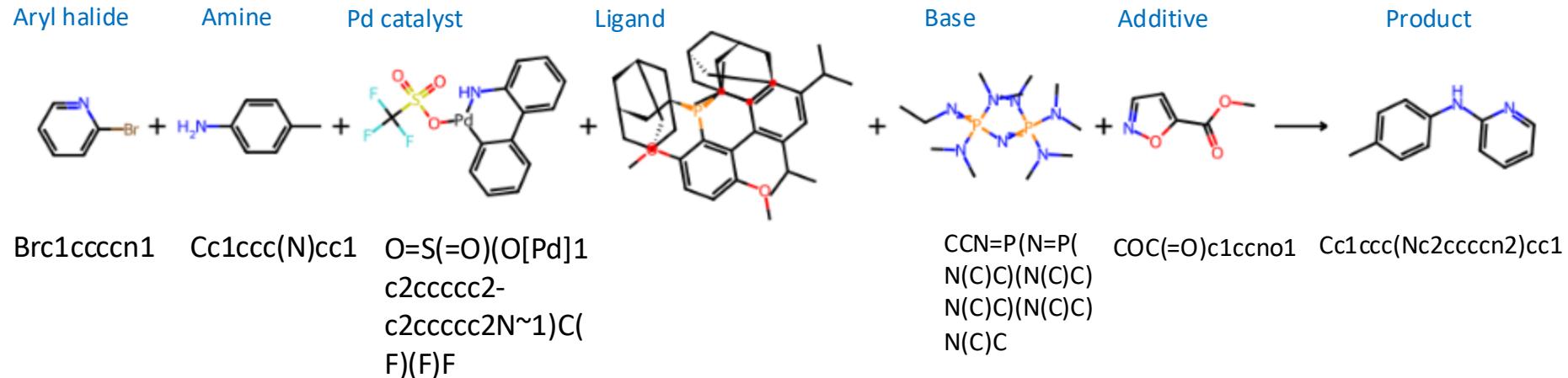


dill/caraway smell

Text-based (or String-based) Representation

□ Reaction Representation

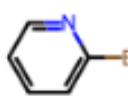
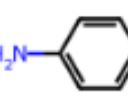
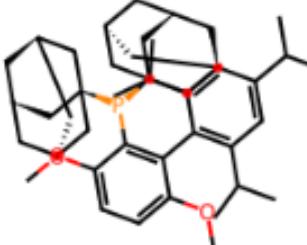
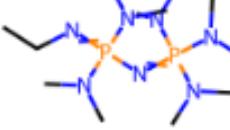
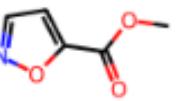
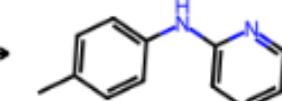
- The SMILES format used for describing molecules has been extended to so-called Reaction SMILES
- Each molecule in the reactants, agents, and products is represented by a SMILES string, and disconnected structures are separated by a period; this includes the individual molecules, ions and ligands, which are listed in an arbitrary manner. Reactants, agents, and products are separated by either the ‘>’ or ‘>>’ symbol



Text-based (or String-based) Representation

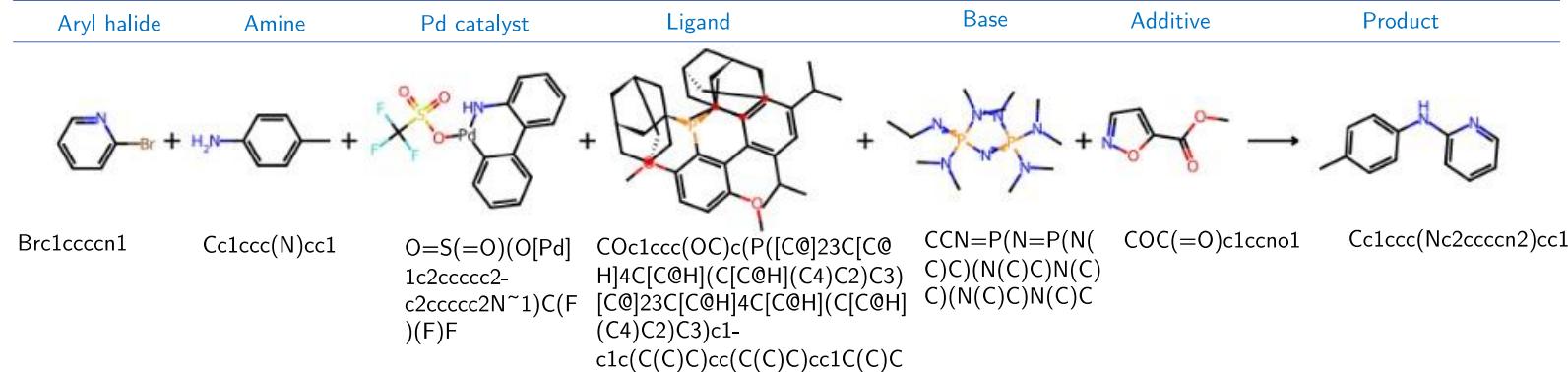
□ Reaction Representation

- The SMILES format used for describing molecules has been extended to so-called Reaction SMILES
- Each molecule in the reactants, agents, and products is represented by a SMILES string, and disconnected structures are separated by a period; this includes the individual molecules, ions and ligands, which are listed in an arbitrary manner. Reactants, agents, and products are separated by either the ‘>’ or ‘>>’ symbol

Aryl halide	Amine	Pd catalyst	Ligand	Base	Additive	Product
 Brcc1ccccc1	 Cc1ccc(N)cc1	 O=S(=O)(O[Pd]1c2cccc2-c2cccc2N~1)C(F)(F)F		 CCN=P(N=P(N(C)C)(N(C)C)N(C)C)(N(C)C)N(C)C	 COC(=O)c1ccno1	 Cc1ccc(Nc2ccccc2)cc1

Text-based (or String-based) Representation

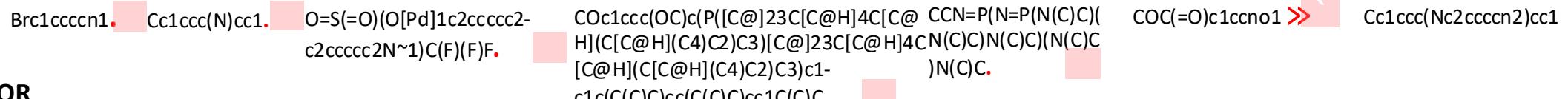
Reaction Representation



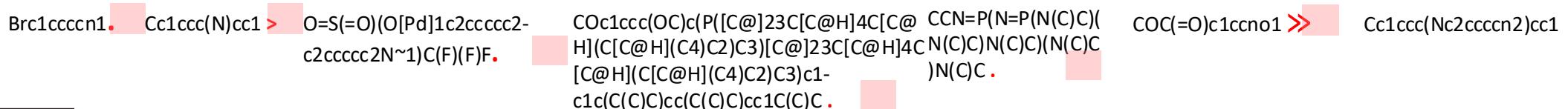
- Each reactant, solvent, additive, catalyst is separated by a “.”
- Sometime, for the substrate which do not contribute atoms to the product is separated by >
- A product is separated by >> (precursors >> products) OR (reactants > substrates >> products)

Dose the order of precursor matter?
Does the SMILES have physical meaning?

Reaction SMILES:



OR

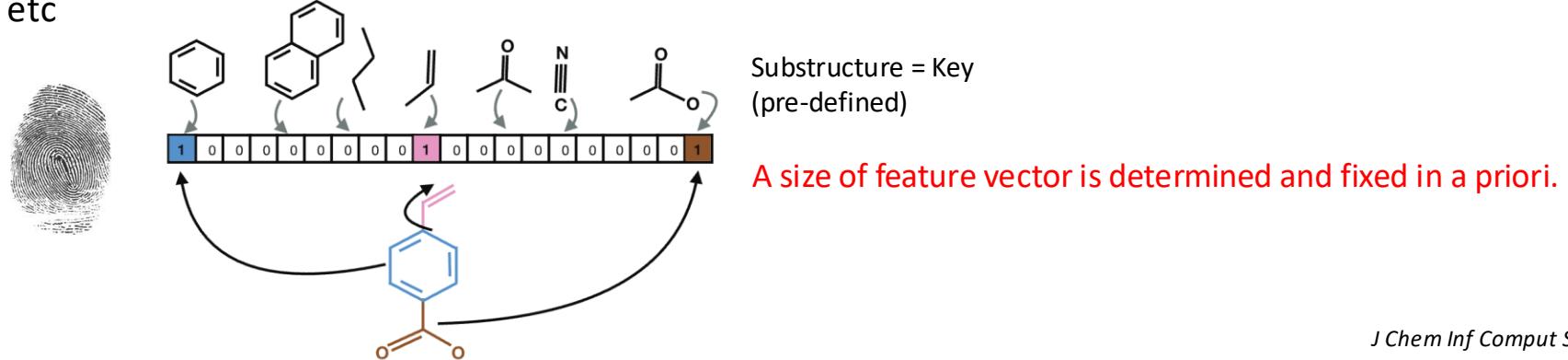


Molecular Fingerprint

- A vectorised representation of molecules capturing precise details of atomic configurations within a molecule.

□ 1) Structural keys

- In structural keys, the structure of a molecule is encoded into a binary bit string (that is, a sequence of 0's and 1's), each bit of which corresponds to a “pre-defined” structural feature (e.g., substructure or fragment).
 - It encode the presence (1) or absence (0) of certain substructures in a compound (one-hot-encoding)
 - It is important to understand that structural keys cannot encode structural features that are not pre-defined in the fragment library.
 - Examples are the **MACCS keys** and PubChem Fingerprints.
- **MACCSkeys** (Molecular ACCess System)
 - The **predefined keys** are implemented in popular open-source cheminformatics software packages, including **RDKit**, **OpenBabel**, **CDK**, etc



J Chem Inf Comput Sci., 2002, 42, 1273

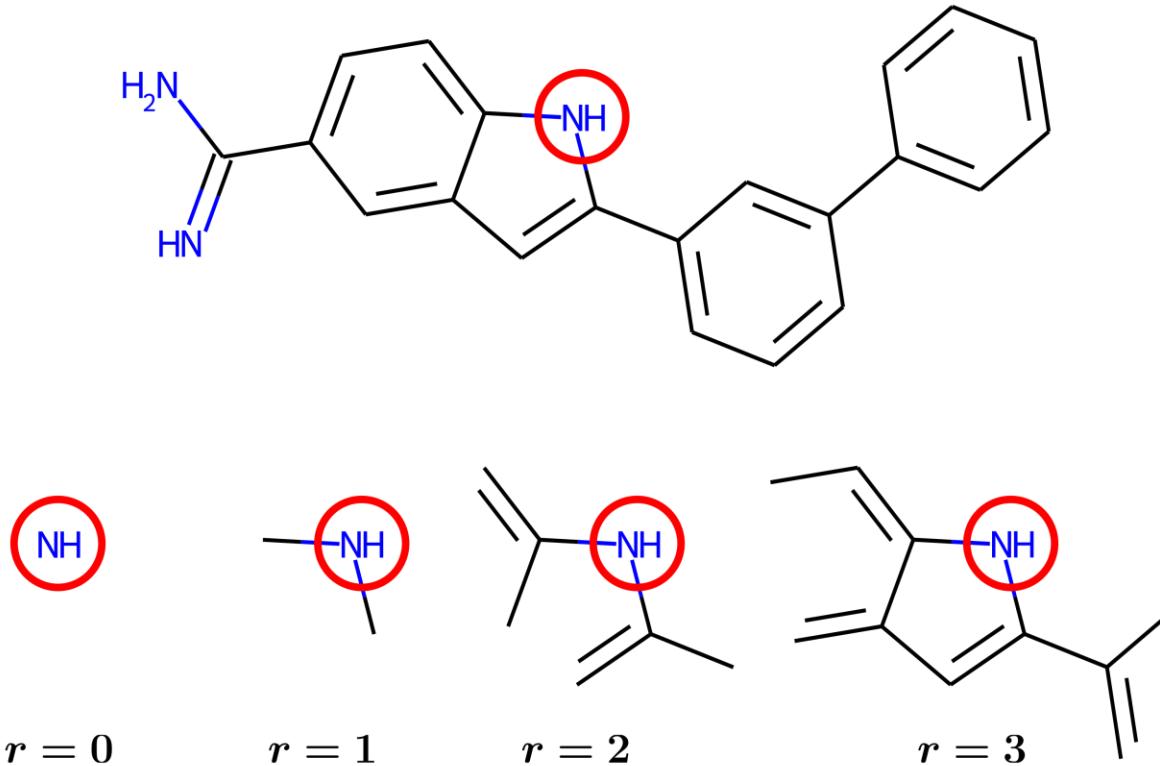
Molecular Fingerprint

□ 2) Hashed keys

- An alternative to structural keys.
- Hashed fingerprints **do not require** a pre-defined fragment library.
- Instead, they are **generated by enumerating through the molecule all possible fragments that are not bigger than a certain size** and then **converting these fragments into numeric values using a “hash” function**.
- These numeric values can be used to indicate bit positions in the hashed fingerprints.
- Hashing them into values within a fixed range inevitably results in “bit collisions”, in which different fragments are converted into the same numeric value (and the same bit position). Because of this, there is no one-to-one correspondence between fragments and fingerprint bits (contrary to structural keys).
- Hashed fingerprints **may be further classified into** topological or **path-based fingerprints** and **circular fingerprints**, according to the way by which the fragments are enumerated.
- Let's focus on circular fingerprint – Morgan fingerprint

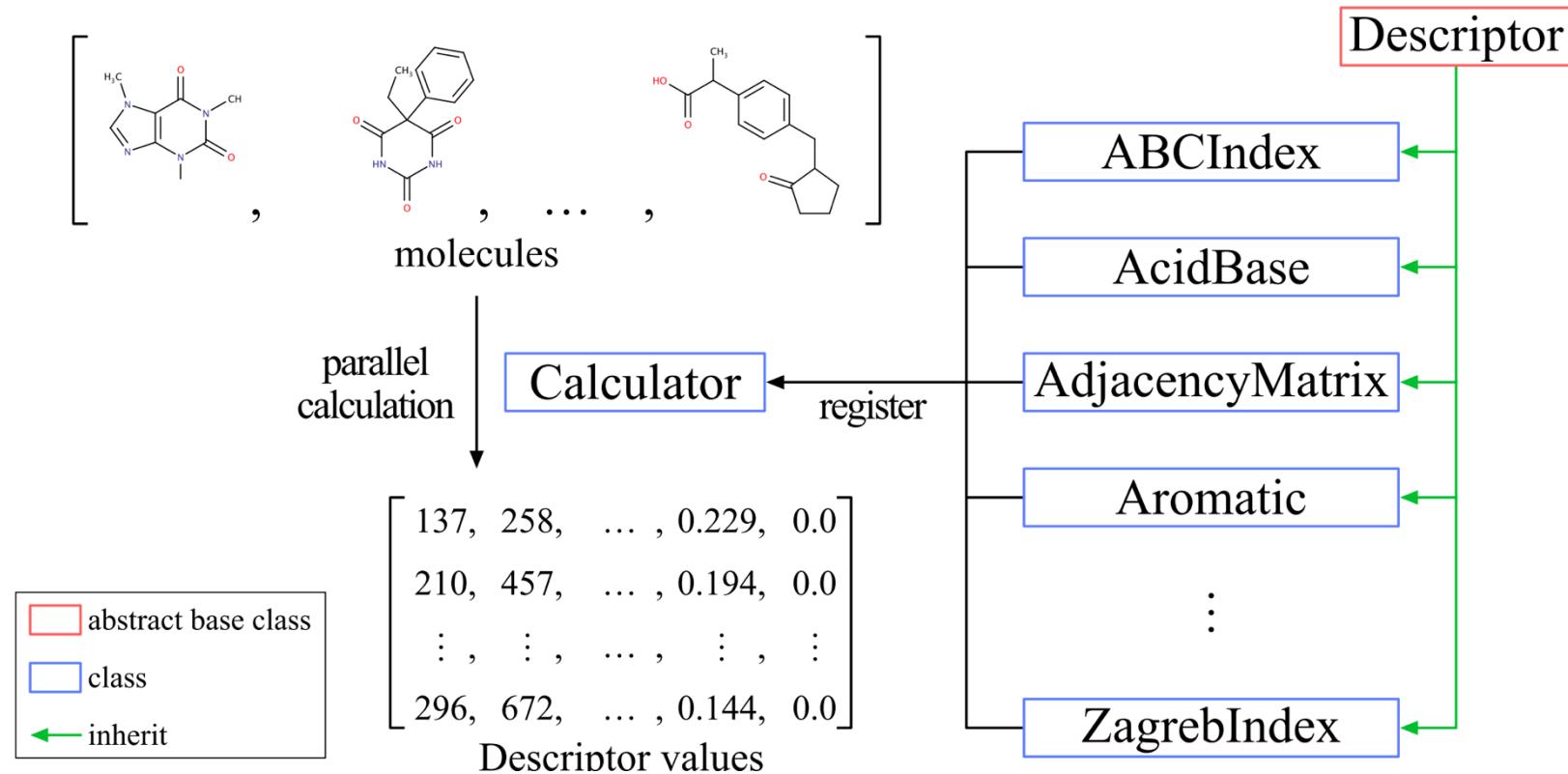
Morgan fingerprints

- family of fingerprints based on the Morgan algorithm
- bits correspond to the circular environments of each atom in a molecule

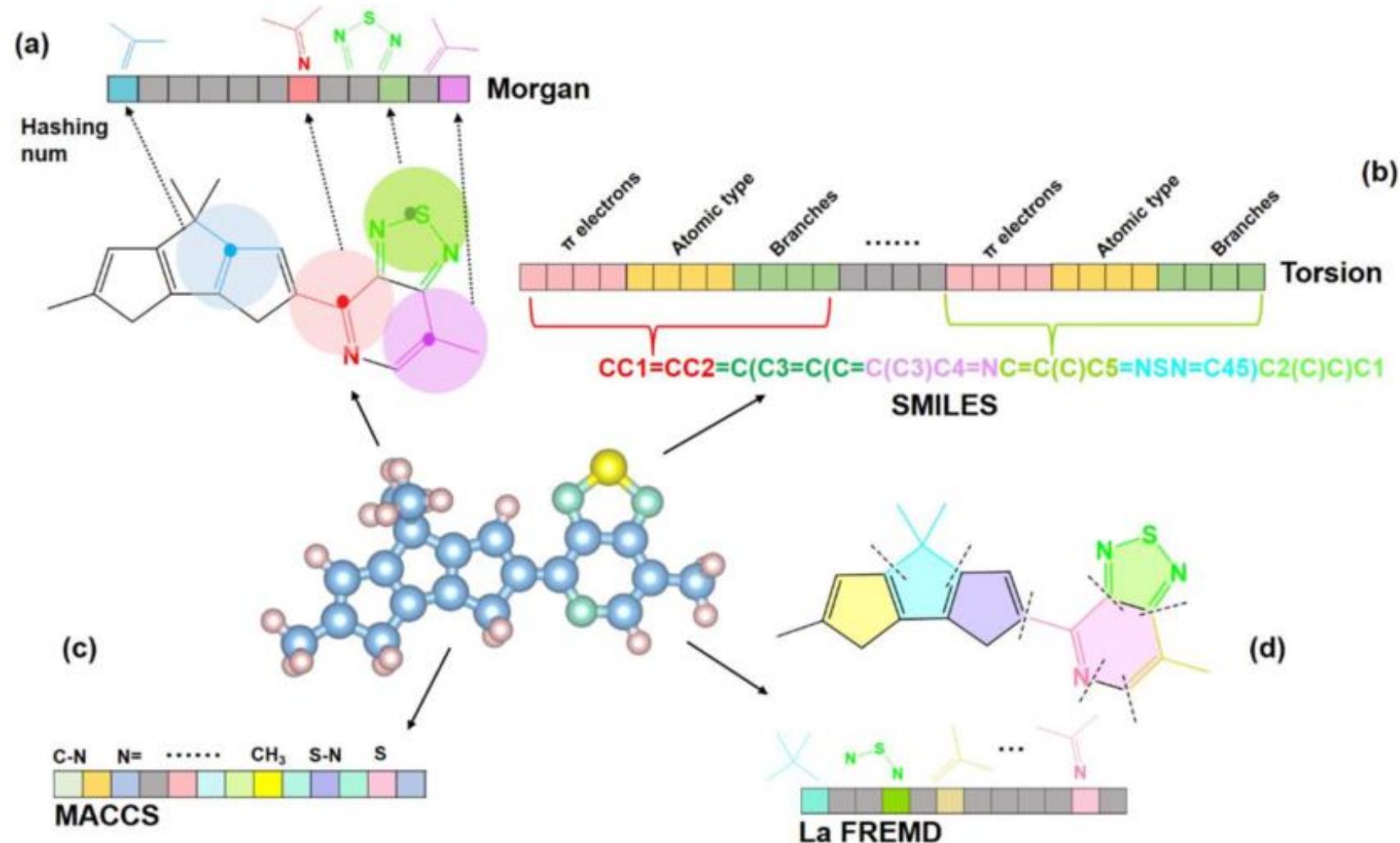


Mordred descriptors

a numerical vector of several molecular properties known to be important in structure-property relationships

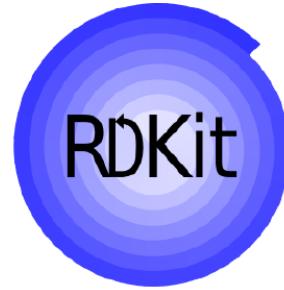


representation comparison



ACS Appl Mater Inter., 2023, 15, 17

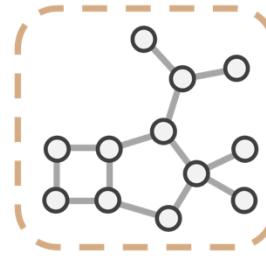
packages to support chemical representations



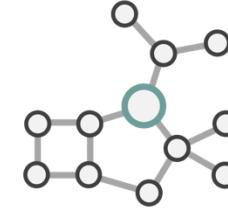
MORFEUS

molecular **features** for machine learning

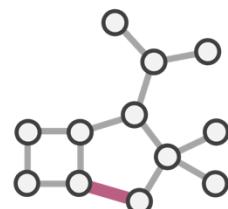
MORDRED



Graph level
e.g. total energy
of a molecule



Node level
e.g. oxidation state
of an atom



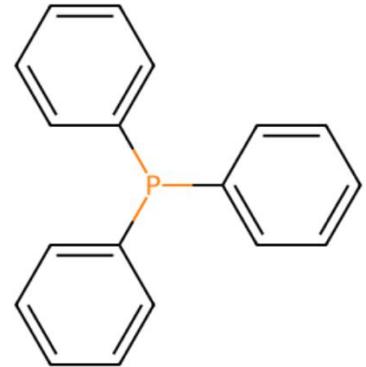
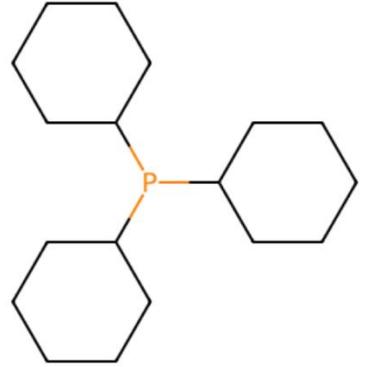
Edge level
e.g. strength of
a bond

023333321
33111222244
331112222333
22 022222343
2301222225 0233333224112222344
23322211344344443111243333330
012135333344444443245
53333444445556664
52123444455556665
43 1444455556672
34234455556662
0345555666530
01110

D Scribe

- learned embeddings (contrastive learning)
- custom fingerprints (MOFid)

what does it look like in practice?



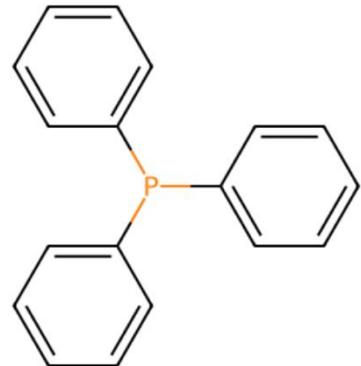
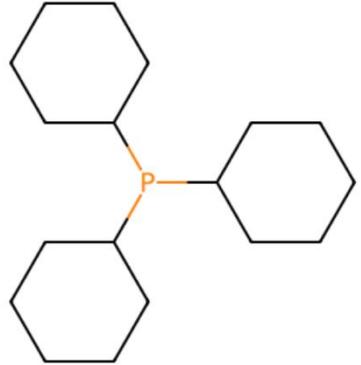
```
from rdkit import Chem
from rdkit.Chem import Draw

# first, generate the rdkit molecule object from the smiles
mol_L6 = Chem.MolFromSmiles('P(C1=CC=CC=C1)(C2=CC=CC=C2)C3=CC=CC=C3')
Draw.MolToImage(mol_L6)

mol_L5 = Chem.MolFromSmiles('C3CCC(P(C1CCCCC1)C2CCCCC2)CC3')
Draw.MolToImage(mol_L5)
```

snappify.com

what does it look like in practice?



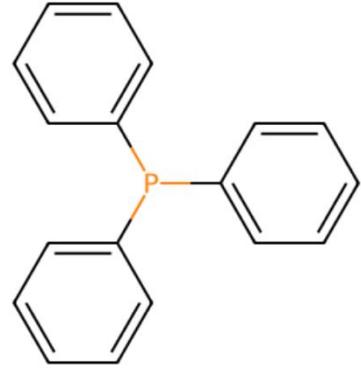
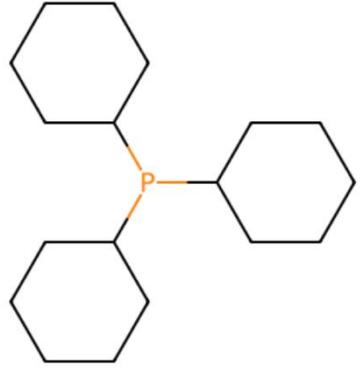
```
from rdkit import Chem
from rdkit.Chem import Draw, AllChem

# first, generate the rdkit molecule object from the smiles
smiles_list = [
    'C3CCC(P(c1ccccc1)c2ccccc2)CC3',
    'P(c1=CC=CC=C1)(c2=CC=CC=C2)c3=CC=CC=C3',
]
rdmols = [Chem.MolFromSmiles(smiles) for smiles in smiles_list]

from rdkit.Chem import AllChem
fpgen = AllChem.GetMorganGenerator(radius=3)
morgan_fingerprints = [
    fpgen.GetSparseCountFingerprint(mol) for mol in rdmols
]
```

snappyf.com

what does it look like in practice?



```
from mordred import Calculator, descriptors

# create descriptor calculator with all descriptors
calc_2D = Calculator(descriptors, ignore_3D=True)
calc_3D = Calculator(descriptors, ignore_3D=False)

df_2D = calc_2D.pandas(rdmols)
df_3D = calc_3D.pandas(rdmols)
```

snappyf.com

BO implementations

Implementing BO for chemistry

Table 2 A collection of open-source Python software libraries for Bayesian optimisation. Surrogate models include Gaussian process (GP), random forest (RF) and tree of Parzen estimators (TPE)

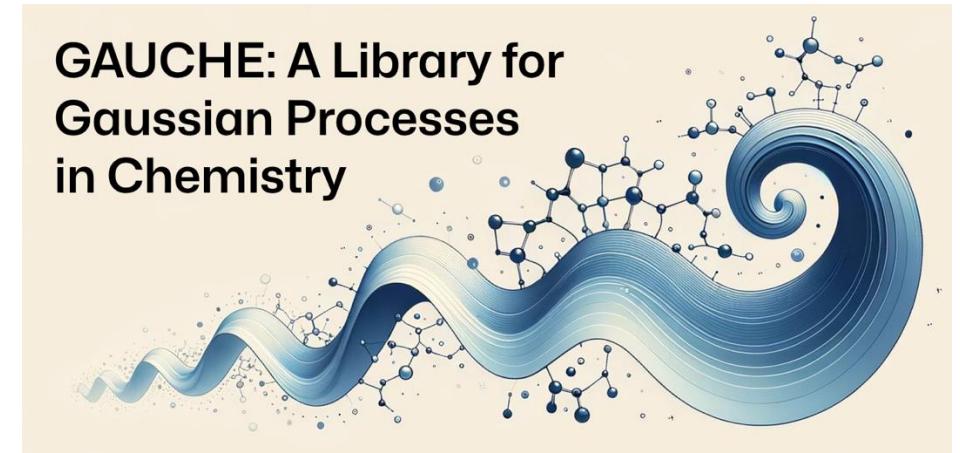
Package	Models	Features	License	Ref.
General purpose				
Ax ^a	GP, others	Modular framework built on BoTorch	MIT	
Bayesianopt ^b	GP	Parallel optimisation	MIT	29
BayesOpt ^c	GP	Single objective	MIT	30
BoTorch ^d	GP, others	Multi-objective optimisation	MIT	31
COMBO ^e	GP	Multi-objective optimisation	MIT	32
Dragonfly ^f	GP	Multi-fidelity optimisation	Apache	33
GPyOpt ^g	GP	Parallel optimisation	BSD	34
Hyperopt ^h	TPE	Serial/parallel optimisation	BSD	35
NEXTorch ⁱ	GP, others	Modular framework built on BoTorch	MIT	36
Optuna ^j	RF	Hyperparameter tuning	MIT	37
Skopt ^k	RF, GP	Batch optimisation	BSD	38
SMAC3 ^l	GP, RF	Hyperparameter tuning	BSD	39
GPax ^m	GP	Multi-task/fidelity	MIT	40 and 41
Physical science domain				
Atlas ⁿ	GP	Mixed-parameter optimisation for self-driving labs	MIT	42
BOSS ^o	GP	Crystal structure optimisation	Apache	43
Edbo ^p	GP	Tailored chemical synthesis descriptors	MIT	5
GAUCHE ^q	GP	Tailored molecular representations	MIT	44
NUBO ^r	GP	Transparent BO to personalise problem	BSD	45
Olympus ^s	GP, TPE, BNN	Benchmarking and noisy optimisation	MIT	46
Phoenics ^t	BNN	Bayesian kernel density estimation	Apache	47
Summit ^u	GP, RF	Multi-task optimisation for chemical reactions	MIT	48

chemistry-specific BO packages and tools

- What makes something chemistry-specific?
 - Chemical/Reaction encoding – one-hot encoding, Mordred, ECFP
- Packages that can readily implement this
 - GAUCHE
 - BayBE
 - BoFire

but what if I don't want to code?

open-source GUIs for BO!



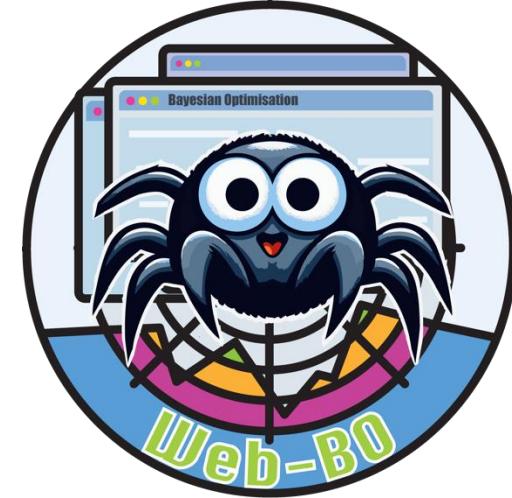
chemistry-specific BO packages and tools

- What makes something chemistry-specific?
 - Chemical/Reaction encoding – one-hot encoding, Mordred, ECFP
- Packages that can readily implement this
 - GAUCHE
 - BayBE
 - BOFire

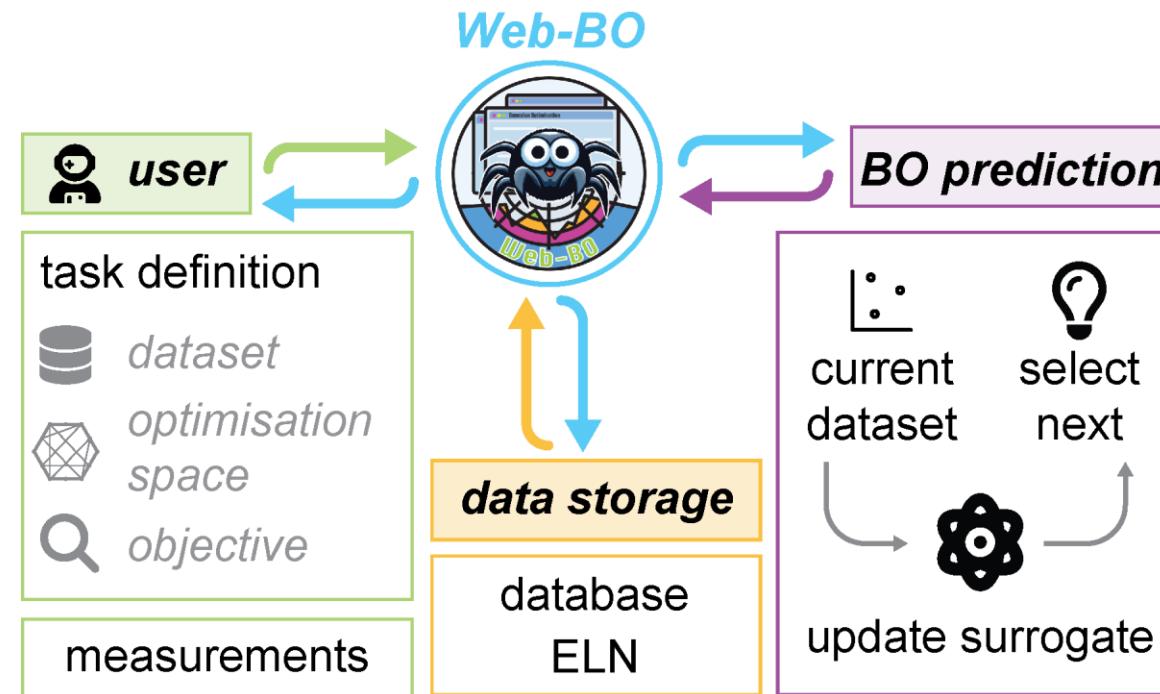


but what if I don't want to code?

open-source GUIs for BO!



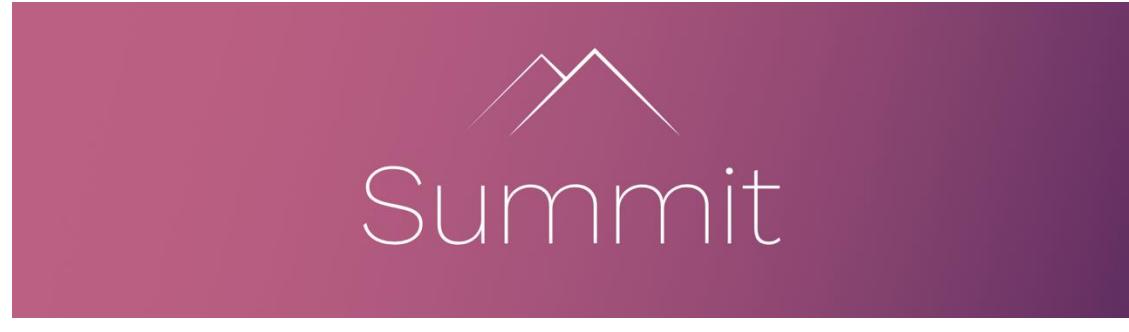
Improving accessibility of data-driven optimisation for chemical tasks via a graphical user interface



Interactive Web-BO task: Let's compare representations for reaction optimization

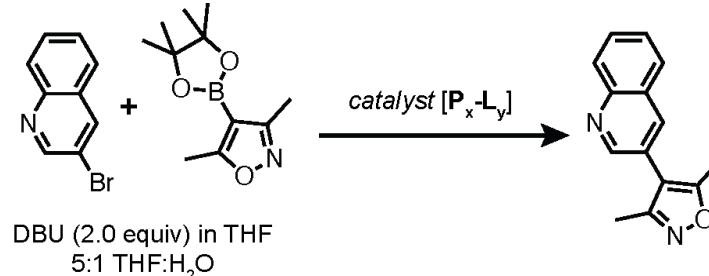
Emulators as tools for benchmarking

- ❑ Allow us to explore varying algorithm formulations without performing experiments
- ❑ Packages that offer emulators supported/maintained for recent versions of PyTorch
 - SUMMIT
 - EDBO
- ❑ Example cases
 - Suzuki-Miyura Cross Coupling optimisation



Emulators as tools for benchmarking

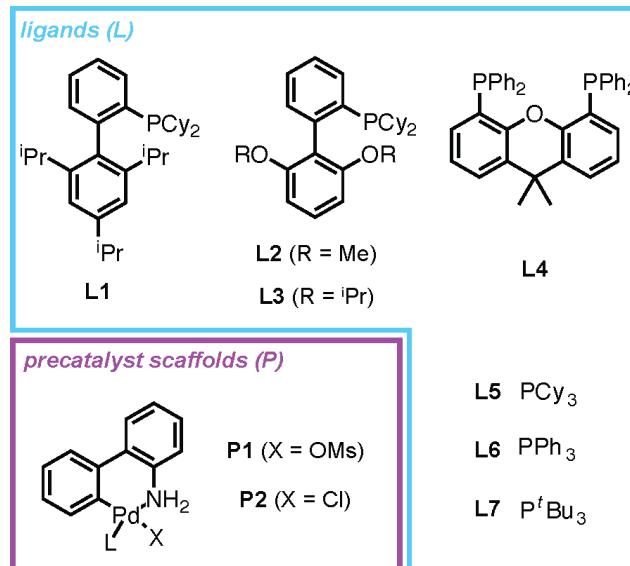
optimizing the coupling of 3-bromoquinoline with 3,5-dimethylisoxazole-4-boronic acid pinacol ester in the presence of 1,8-diazobicyclo[5.4.0]undec-7-ene (DBU) and THF/water



"experiments" are supported by an emulator developed by Lapkin, et al.

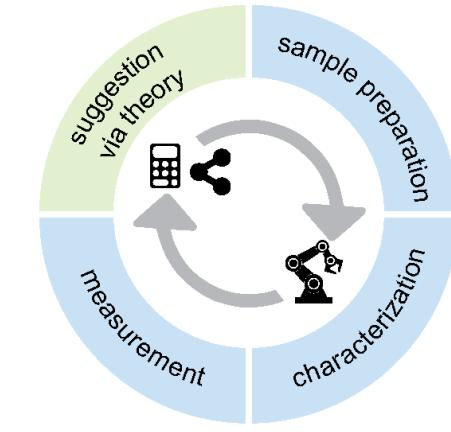
Summit

parameter	parameter type	parameter space
catalyst [P _x -L _y]	categorical	[P1-L1], [P2-L1], [P1-L2], [P1-L3] [P1-L4], [P1-L5], [P1-L6], [P1-L7]
catalyst loading	continuous	(0.5 - 2.0 %)
temperature	continuous	(30 - 110 °C)
residence time	continuous	(1 - 10 min.)



let's see how different chemical representations compare for our catalysts

Let's code!



suprashare



tutorial notebooks

break