A tour through Molecular Representations in Al-driven Drug Discovery



Laurianne David

David, L., Thakkar, A., Mercado, R. *et al.* Molecular representations in AI-driven drug discovery: a review and practical guide. *J Cheminform* **12**, 56 (2020). <u>https://doi.org/10.1186/s13321-020-00460-5</u>

The Field of Molecular Representation



- Molecular representations have been of interest since the 19th century and many challenges arose to create a functional representation
- Challenges are related to the diversity of chemical structures and their complexity
- Empirical formula of Alanine: C₃H₇NO₂
- This formula can also match sarcosine and lactamide.
- A chemical representation is any encoding of a chemical compound
- Linear notations are referred to as notations



Graph Representations for Small Molecules



- In a molecular graph representation, the atoms and bonds of a molecule are represented as a set of nodes and edges.
- How the atoms are connected can be represented in an adjacency matrix.
- The node order used in a matrix representation is determined by a graph traversal algorithm.





- Matrix representations of graphs are node-order dependent.
- If consistency is required, depth-first search and breadth-first search are applicable.
- If necessity of noisier data: random search (e.g. deep-learning)

Graph Representations for Small Molecules



Graph Representations: Related Formats



Linear Notations of Small Molecules



- Matrix representations:
 - require a large amount of disk space
 - Are not well adapted to basic cheminformatics analysis
- Molecules are often represented as strings of characters encoding the Ctab.

D-Alanine (implicit H)				
Molfile	612 bytes			
SMILES	15 bytes			
InChI	59 bytes			

• Linear notations are compact, easy to manipulate.

Storyline of Linear Notations

Alchemy: Compounds and elements named based on properties (e.g. aqua fortis = nitric acid, sweet oil of vitriol = diethyl ether)

20th century : Systematic nomenclature of organic chemistry described in the IUPAC Color Books and used for literature, patents, governments legislation.

1949-1961: International standard for electronic chemical notations





- 5. Ability to generate a unique chemical nomenclature
- 6. Compatibility with accepted practices of inorganic chemical nomenclature
- 7. Uniqueness
- 8. Generation of an unambiguous and useful enumeration pattern
- 9. Ease of manipulation by machine methods
- 10. Exhibition of associations
- 11. Ability to deal with partial indeterminates

- Non-ambiguous (i.e. notation will regenerate only the original compound)
- Chosen notation by IUPAC in 1961: Dyson cyphering
 - Drawbacks: contained many arbitrary rules / could not be handled on standard typewriters or ordinary punched-card machines

SMALL MOLECULES REACTIONS MACROMOLECULES REPRESENTATIONS
--





Punched-card machine

Typewriter



- 7. Uniqueness
- 8. Generation of an unambiguous and useful enumeration pattern
- 9. Ease of manipulation by machine methods
- 10. Exhibition of associations
- 11. Ability to deal with partial indeterminates

- Drawbacks: could not be handled on standard typewriters or ordinary punched-card machines / contained many arbitrary rules
- Most used notation by the community: Wiswesser Line Notation (WLN)

SMALL MOLECULES	REACTIONS	МА	CROMOL		GRAPHICAL EPRESENTATIONS			
	E	Bromine atom	F	Fluorine atom				
	G	Chlorine atom	н	Hydrogen atom				
	1	lodine atom	Q	Hydroxyl group, -OH				
	R	Benzene ring	S	Sulfur atom				
	U	Double bond	UU	Triple bond				
	V	Carbonyl, -C(=O)-						
	С	Unbranched carbon multiply bonded to non-carbon atom						
	к	Nitrogen atom bonded to more than three other atoms						
	L	First symbol of a ca	rbocyclic ring	notation				
<u>Wiswesser Line Notation</u>	M	Imino or imido -NH- group Nitrogen atom, hydrogen free, bonded to fewer than 4 atoms						
	N							
	0	Oxygen atom, hydrogen-free						
	т	First symbol of a heterocyclic ring notation						
	W	Non-linear dioxo group, as in -NO ₂ or -SO ₂ -						
	X	Carbon attached to four atoms other than hydrogen						
	Y	Carbon attached to	three atoms	other then hydrogen				
	Z	Amino and amido NH ₂ group						
	<digit></digit>	Digits "1" to "9" denote unbranched alkyl chains						
	&	Sidechain terminator or, after a space, a component separator 4						
		WLN fo	or Alanine	: QVYZ				

NH₂

Simplified Molecular Input Line Entry System (SMILES)



- SMILES were developed by Weininger et al. in 1988
- SMILES are:
 - Non-unique and unambiguous
 - Obtained by assigning a number to each atom in the molecule, then traversing the molecular graph using that order
 - Rdkit version uses depth-first search
- SMILES can encode for stereochemistry
- SMILES can be generated canonically
- Canonical implementations vary between companies and research teams
- Randomized SMILES can be generated by changing the atom numbering in the structure:
 - Such SMILES can be employed for data augmentation in QSAR modelling or molecular generative modelling



CC(=0)Oc1cccc1C(=0)O

clcc(c(ccl)C(0)=0)OC(C)=0

Simplified Molecular Input Line Entry System (SMILES)



Molecular Descriptors



- Molecular descriptors such as structural keys and hashed fingerprints are notations encoding physicochemical, structural, topological, and/or electronical properties of a compound
- Descriptors are unique and ambiguous
- Structural keys are bit strings, encoding for the absence (0) and presence (1) of a specific chemical group (e.g. MACCS Keys)
- Chemical fingerprints are vectors containing indexed elements encoding for physicochemical or structural properties (e.g. Extended Connectivity Fingerprints ECFP)

smartsPatts = {

```
1: ('?', 0), # ISOTOPE
```

```
#2:('[#104,#105,#106,#107,#106,#109,#110,#111,#112]',0), # atomic num >103 Not complete
```

2: ('[#104]', 0), # limit the above def'n since the RDKit only accepts up to #104

```
3: ('[#32,#33,#34,#50,#51,#52,#82,#83,#84]', 0), # Group IVa,Va,VIa Rows 4-6
```

- 4: ('[Ac,Th,Pa,U,Np,Pu,Am,Cm,Bk,Cf,Es,Fm,Md,No,Lr]', 0), # actinide
- 5: ('[Sc,Ti,Y,Zr,Hf]', 0), # Group IIIB,IVB (Sc...)
- 6: ('[La,Ce,Pr,Nd,Pm,Sm,Eu,Gd,Tb,Dy,Ho,Er,Tm,Yb,Lu]', 0), # Lanthanide
- 7: ('[V,Cr,Mn,Nb,Mo,Tc,Ta,W,Re]', 0), # Group VB,VIB,VIIB
- 8: ('[!#6;!#1]1~*~*~*~1', 0), # QAAA@1
- 9: ('[Fe,Co,Ni,Ru,Rh,Pd,Os,Ir,Pt]', 0), # Group VIII (Fe...)
- 10: ('[Be,Mg,Ca,Sr,Ba,Ra]', 0), # Group IIa (Alkaline earth)
- 11: ('*1~*~*~1', 0), # 4M Ring
- 12: ('[Cu,Zn,Ag,Cd,Au,Hg]', 0), # Group IB,IIB (Cu..)
- 13: ('[#8]~[#7](~[#6])~[#6]', 0), # ON(C)C
- 14: ('[#16]-[#16]', 0), # S-S

Some keys from MACCS Keys (<u>rdkit/MACCSkeys.py at master ·</u> <u>rdkit/rdkit · GitHub</u>)

International Chemistry Identifier (InChI)



- InChI is an open-source canonical notation introduced in 2006 by NIST
- InChI are composed of multiple layers such as the *Main, Charge, Stereochemical* and *Isotopic layers* (non exhaustive lit)
- InChI might not be decodable back to the molecular graph of origin
- A hashed version of InChI, InChIKey is used for library searching
- There exists a universal SMILES system which relies on the InChI representation





Representation for Chemical Reactions



- Approximately 141 million reactions have been recorded since 1840
- Chemical reactions represent the interconversion of one set of molecules into another related set, under a set of specific conditions
- Graphical representations of reactions are not easily machine-readable
- Several notations and representations are available to solve this problem

Notations (non-exhaustive)	Details
Reaction SMILES	Encode reactants, agent and product Storage of reactions conditions or reaction center not supported
SMIRKS	Define generic reactions transformations Can describe reaction centre, enumerate libraries
RInChI	Direction of reaction can be described Enables identification of duplicate reactions ProcAuxInfo extension stores metadata (yield, temperature,)

Representation for Chemical Reactions



- A method to represent reactions was developed by Saller et. al. (InfoChem CLASSIFY)
- The steps are the following:
 - 1. Identify and extract the reaction centre (i.e. set of atoms that have changed their number of implicit hydrogens, valency, number of electrons, etc ...)
 - 2. Atom hash codes are calculated for all atoms belonging to the reaction centre. A unique representation of the reaction centre is obtained



Representation for Macromolecules



Amino-acid based structures



- The Hierarchical Editing Language for Macromolecules (HELM) was developed by Pfizer under the auspices of the Pistoia Alliance.
- HELM can be used to describe peptides, antibodies, polymers, nucleotides, ...
- HELM is widely adopted by the community





Amino-acid based structures



Key Macromolecules



Graphical Representations



- Published data often have molecules given as 2D depictions, which makes it complicated to mine the data
- Optical Chemical Recognition (OCR) systems are required to translate a depiction into representation or notation.

Graphical Representations



Graphical Representations



Conclusions

SMALL MOLECULES REACTIONS MACROMOLECULES	GRAPHIC REPRESENTATI	CAL TIONS	
	Notation system	Representation	
	Generic name	D-Alanine	
 Molecular representations must describe: different types of structures (small molecules, peptides, polymers,) 	IUPAC name	(2R)-2-aminopropanoic acid (English) acide (2R)-2- aminopropanoïque (French) (2R)-2-аминопропановая (Bussian)	
 with different properties (stereochemistry, valence,) 	WLN ^a	QVYZ [46]	
	SMILES	C[C@H](C(=O)O)N	
 Precise representations for specific structures are needed to optimize the process of AI-driven discovery 	InChl	InChI = 1S/C3H7NO2/c1- 2(4)3(5)6/h2H,4H2,1H3,(H,5,6) InChI = 1S/C3H7NO2/c1- 2(4)3(5)6/h2H,4H2,1H3,(H,5,6)/t2-/m1/s1	
• Our aim was to provide a general overview of these representations,	InChl Key	QNAYBMKLOCPYGJ- UWTATZPHSA-N	
some which are well known in our field of AI-driven drug discovery,	HELM	PEPTIDE1{[dA]}	
and some which are specialized	Three-letter symbol	D-Ala	
	Protein Line Notation (PLN)	H-{d}A-OH	

Acknowledgments



- Rocío Mercado (<u>rociomer@mit.edu</u>)
- Amol Thakkar (<u>tha@zurich.ibm.com</u>)
- Ola Engkvist (AZ)
- Hongming Chen
- Noé Sturm
- Thierry Kogel
- Lionel Colliandre
- All the BigChem Fellows and PIs

CO-AUTHORS

https://jcheminf.biomedcentral.com/a rticles/10.1186/s13321-020-00460-5



The project leading to this presentation received funding from the European Union's Horizon 2020 research and innovation program under the Marie Skłodowska-Curie grant agreement No. 676434, "Big Data in Chemistry" ("BIGCHEM", <u>http://bigchem.eu</u>). This presentation reflects only the authors' view, and neither the European Commission nor the Research Executive Agency are responsible for any use that may be made of the information it contains.

Bibliography

- <u>Molecular representations in AI-driven drug discovery: a review and practical guide | Journal of Cheminformatics</u> | <u>Full Text (biomedcentral.com)</u>
- rdkit/MACCSkeys.py at master · rdkit/rdkit · GitHub
- Polymers as drugs—Advances in therapeutic applications of polymer binding agents Connor 2017 Journal of Polymer Science Part A: Polymer Chemistry - Wiley Online Library
- HELM Project Pistoia Alliance
- <u>PepSeA: Peptide Sequence Alignment and Visualization Tools to Enable Lead Optimization | Journal of Chemical Information and Modeling (acs.org)</u>
- <u>BigSMILES Olsen Research Group (mit.edu)</u>
- <u>Nanome</u>