



258th ACS National Meeting San Diego Aug 2019

Measuring R group similarity using medicinal chemistry data

Noel O'Boyle and Roger Sayle

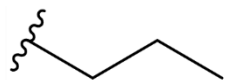
NextMove Software



INTRODUCTION

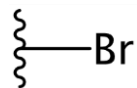


PREDICTING THE PAST

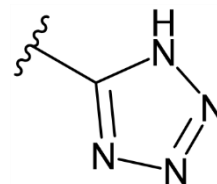


*is usually made
before*

?



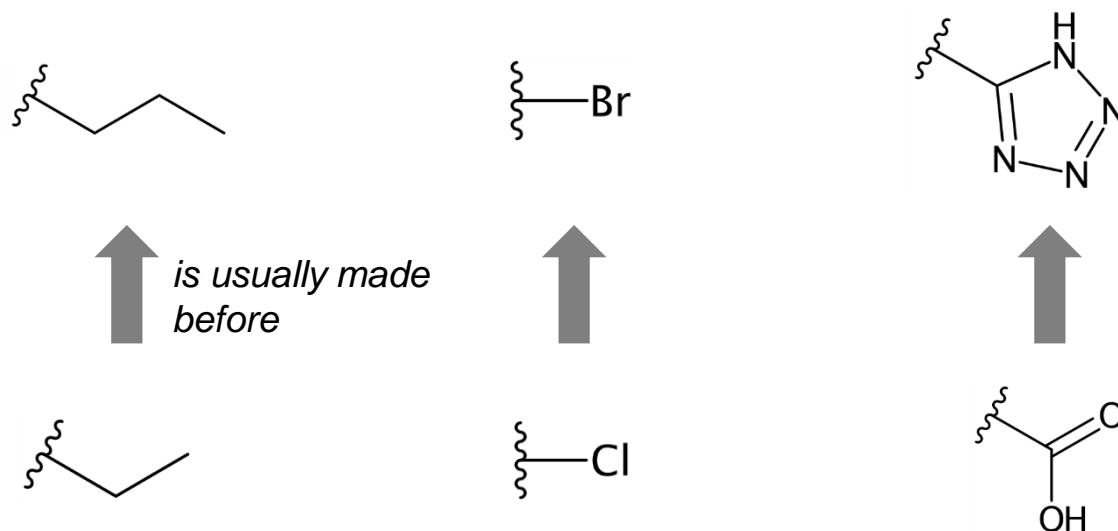
?



?



PREDICTING THE PAST



- Given a med chem project containing a molecule with a Br, it is likely that a Cl was tried at that position at an earlier stage
- This relationship between Cl and Br is a measure of similarity
 - A type of similarity relevant to medicinal chemistry projects
 - A type of similarity undetected by molecular fingerprints
- Can we find all pairs of R groups similarly related and exploit these relationships?



TIME SERIES DATA

- The **ideal dataset** is time series data from one or more large pharmaceutical companies over multiple projects
 - Use date of registration to track project progression
 - Caveats: compounds may have been registered in a previous project, parallel syntheses may obscure conceptual order
- For each project, identify all matched pairs
 - For matched pairs of Br/Cl, was the chlorine analog generally made first? (etc...)



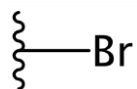
TIME SERIES DATA

- The **ideal dataset** is time series data from one or more large pharmaceutical companies over multiple projects
 - Use date of registration to track project progression
 - Caveats: compounds may have been registered in a previous project, parallel syntheses may obscure conceptual order
- For each project, identify all matched pairs
 - For matched pairs of Br/Cl, was the chlorine analog generally made first? (etc...)
- **Unfortunately....**
 - We don't have access to pharmaceutical data
 - ChEMBL/patents don't include time series information
 - Though I did try looking at papers published by the same author over time

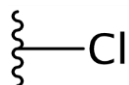


CO-OCCURRENCE INFORMATION

- But ChEMBL/patents **do** have a large number of independent datasets
 - Each is a subset of data from a medicinal chemistry project
 - Co-occurrence of R group analogs within the same subset is informative
- If Br is generally preceded by Cl, then:
 - Molecules with Br will often co-occur with the Cl analog
 - Analogs with Cl will be more frequent than analogs with Br
 - Since after making/testing Cl, Br may not be made



*is usually made
before*



FINDING
CO-OCCURRENCE
DATA



CHEMBL25

- Extract molecules associated with (non-duplicate) IC₅₀, EC₅₀ or K_i data, where *doc_type* is publication or patent
- Within each document, fragment to find single-cut matched series
- Normalise R groups to the most common form using a tautomer hash
 - (see my other talk*)
- Multiple co-occurrences within the same document only count as a single observation

* CIN 55: *Making a hash of it: The advantage of selectively leaving out structural information* (<https://nextmovesoftware.com/talks.html>)



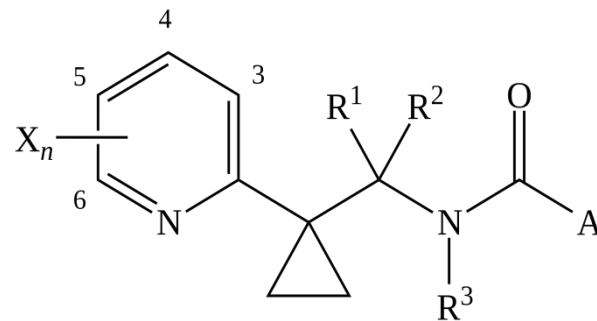
US PATENT APPLICATIONS

- LeadMine was used to extract 2.9M key molecules from pharmaceutically-related USPTO patent applications
 - Key molecules were those associated with bioactivity or present in R Group tables

TABLE 9

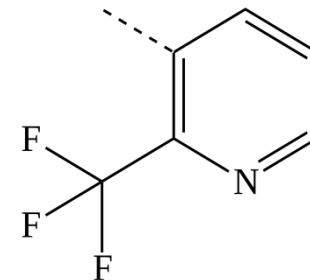
Inhibition of CETP Activity by Examples in Reconstituted Buffer Assay.

Ex. No.	IC ₅₀ (μM)	Ex. No.	IC ₅₀ (μM)	Ex. No.	IC ₅₀ (μM)
249	0.020	419	0.19	425	0.34
244	0.029	230	0.20	514	0.34
634	0.032	248	0.20	237	0.35
221	0.034	266	0.20	399	0.35
229	0.034	378	0.20	645	0.35



No.	X _n	R ¹	R ²	R ³	A
-----	----------------	----------------	----------------	----------------	---

1	6-Cl	H	H	H	
---	------	---	---	---	--



US PATENT APPLICATIONS

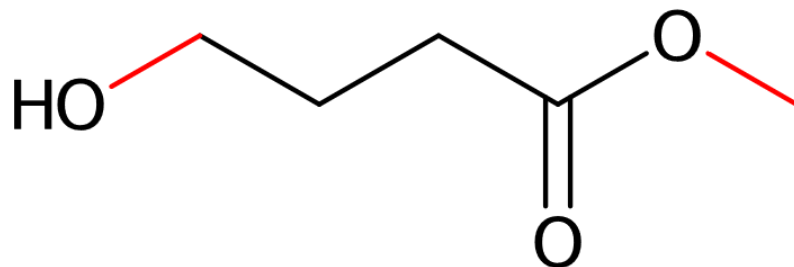
- LeadMine was used to extract 2.9M key molecules from pharmaceutically-related USPTO patent applications
 - Key molecules were those associated with bioactivity or present in R Group tables
- Patents were grouped into chemically-related patent families* (CRPFs)
- Molecules within the same patent were fragmented to find matched series
- Multiple co-occurrences within the same CRPF only count as a single observation

* <https://nextmovesoftware.com/blog/2017/07/04/chemically-related-patent-families/>



FRAGMENTATION SCHEME

- Intended to generate synthetically-relevant matched series*
- Acyclic single bonds broken:
 - If either end is in a ring OR
 - If the bond is between a non sp^2 -hybridised carbon atom and a non-carbon atom
- R group must be ≤ 20 heavy atoms

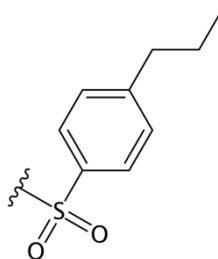
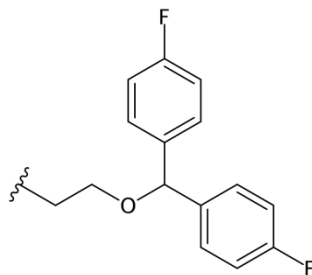
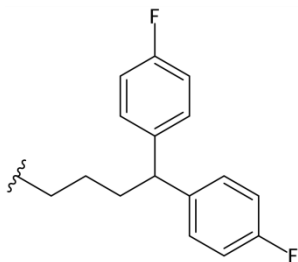
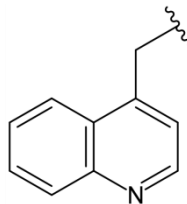
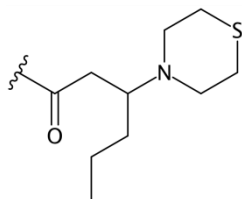
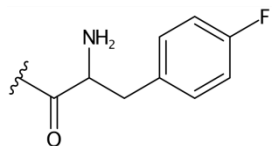
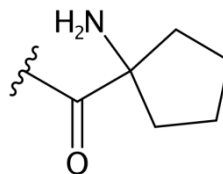
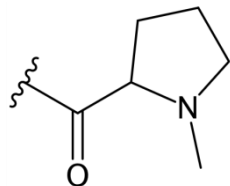
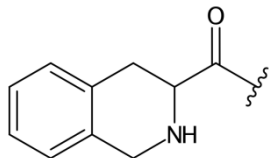
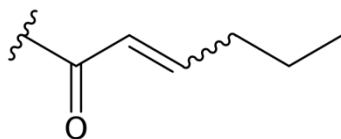


* As used in O'Boyle, Bostrom, Sayle, Gill. *Using Matched Molecular Series as a Predictive Tool To Optimize Biological Activity*. *J. Med. Chem.* **2014**, 57, 6.

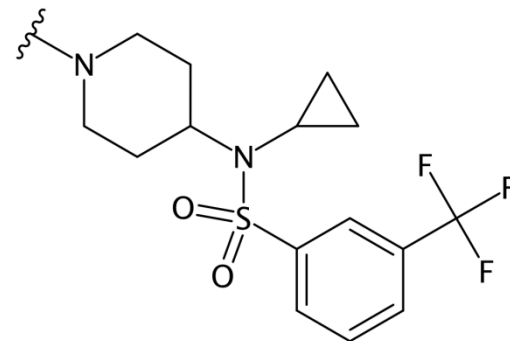


SINGLE-CUT MATCHED SERIES

R Groups

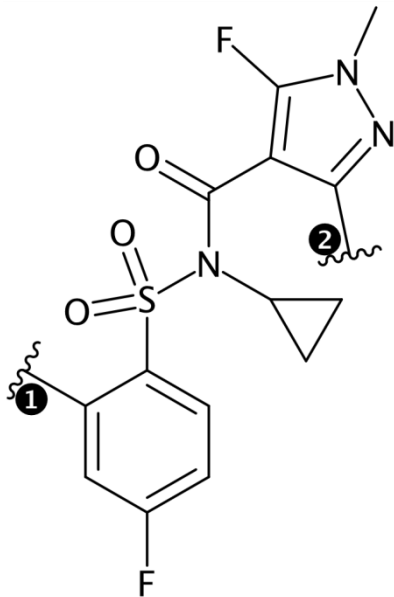


Scaffold



DOUBLE-CUT MATCHED SERIES

- To maximize recovery of R group co-occurrences, also infer co-occurrences from double-cut data

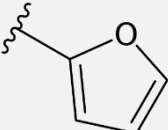
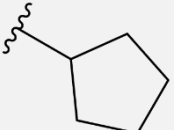
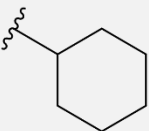
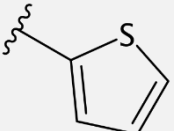
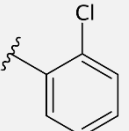
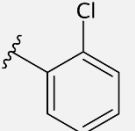
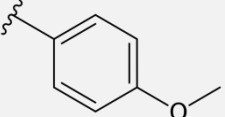
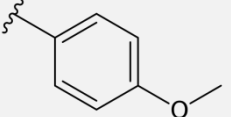


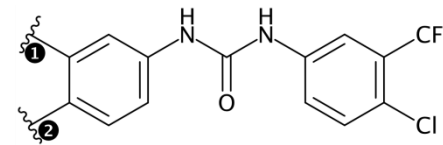
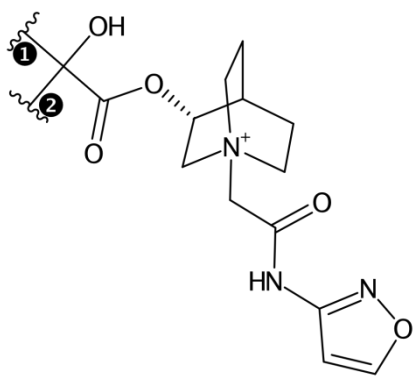
R1	R2

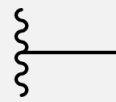
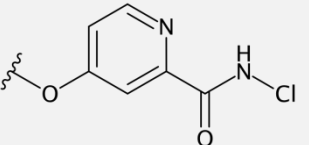
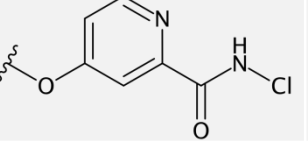
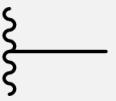


DOUBLE-CUT MATCHED PAIRS

- To maximize recovery of R group co-occurrences, also infer co-occurrences from double-cut data
- Carefully...
 - Where R1 and R2 positions are symmetry-related, discard if $R1 \neq R2$
 - If any R1 appears at R2 (or v.v.), discard the whole series

R1	R2
	
	
	
	

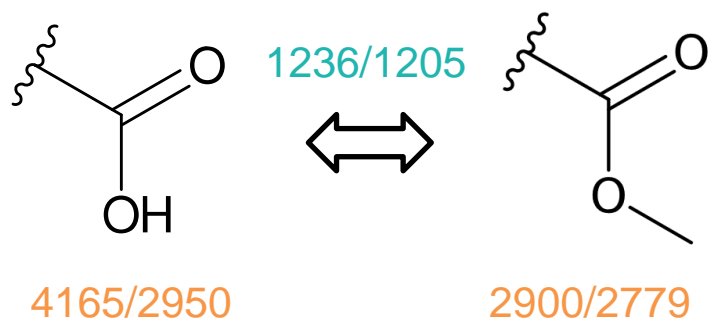


R1	R2
	
	

DATA-DRIVEN R_GROUP SIMILARITY



FREQUENCY AND CO-OCCURRENCE OF R GROUPS

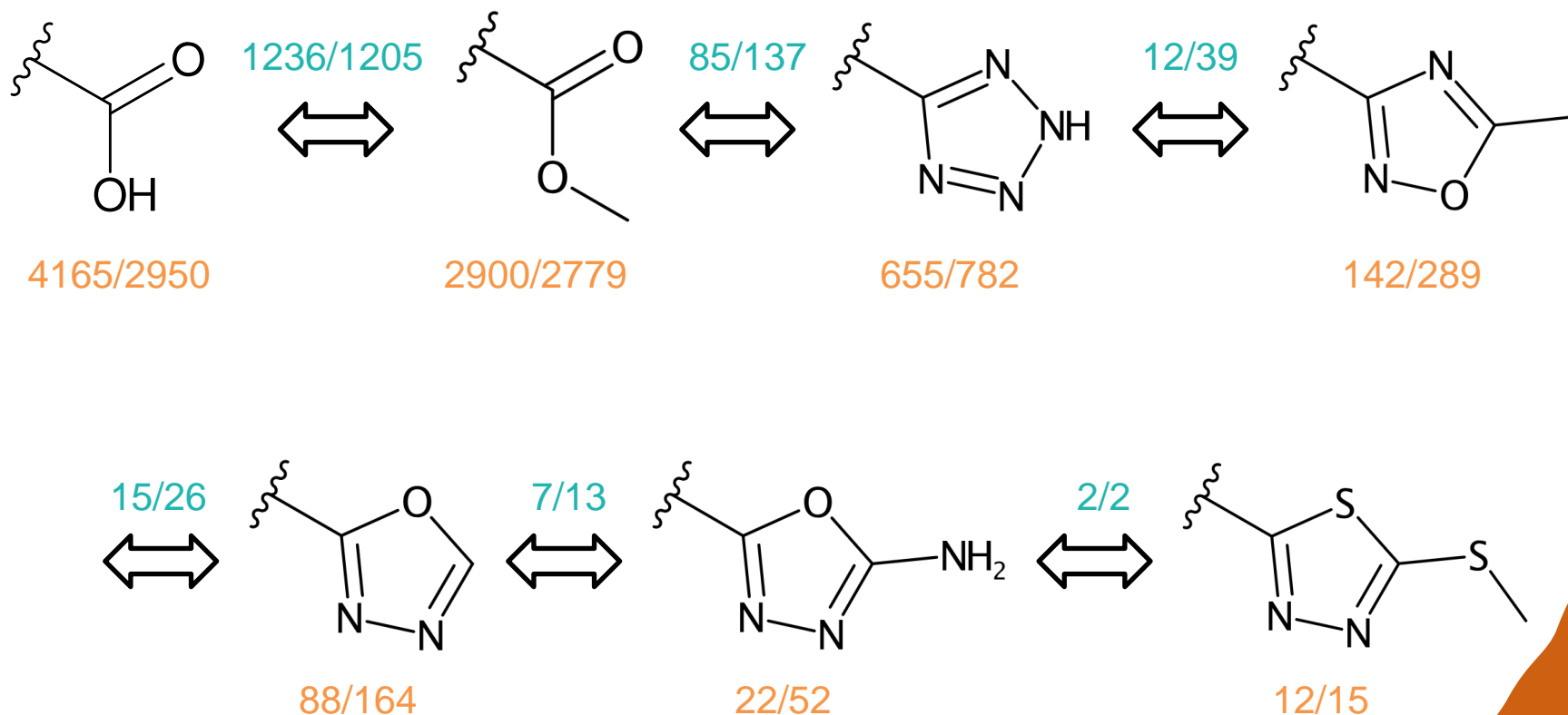


Color key: *frequency*
co-occurrence

Dataset key: ChEMBL/USPTO



FREQUENCY AND CO-OCCURRENCE OF R GROUPS

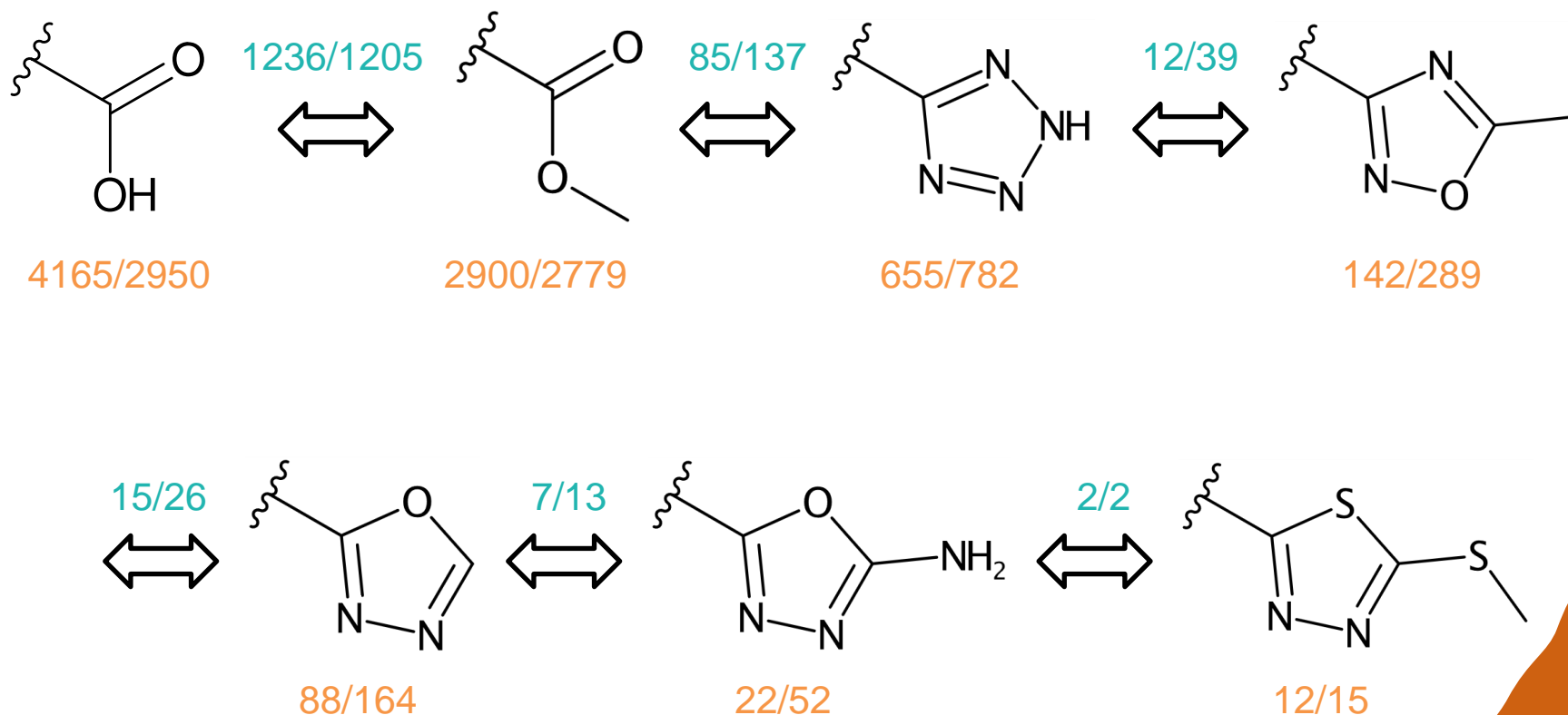


Color key: *frequency*
co-occurrence

Dataset key: ChEMBL/USPTO



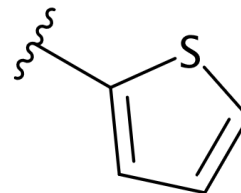
FREQUENCY AND CO-OCCURRENCE OF R GROUPS



Decreasing frequency \approx likely to be synthesised afterwards (if at all)



Find the 8 R groups most similar to



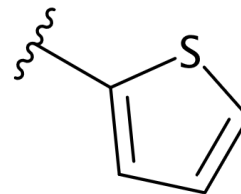
More frequent

- We need to balance high co-occurrence against high frequency
 - Everything has a high co-occurrence with methyl, for example
- Rank R groups based on co-occurrence divided by frequency

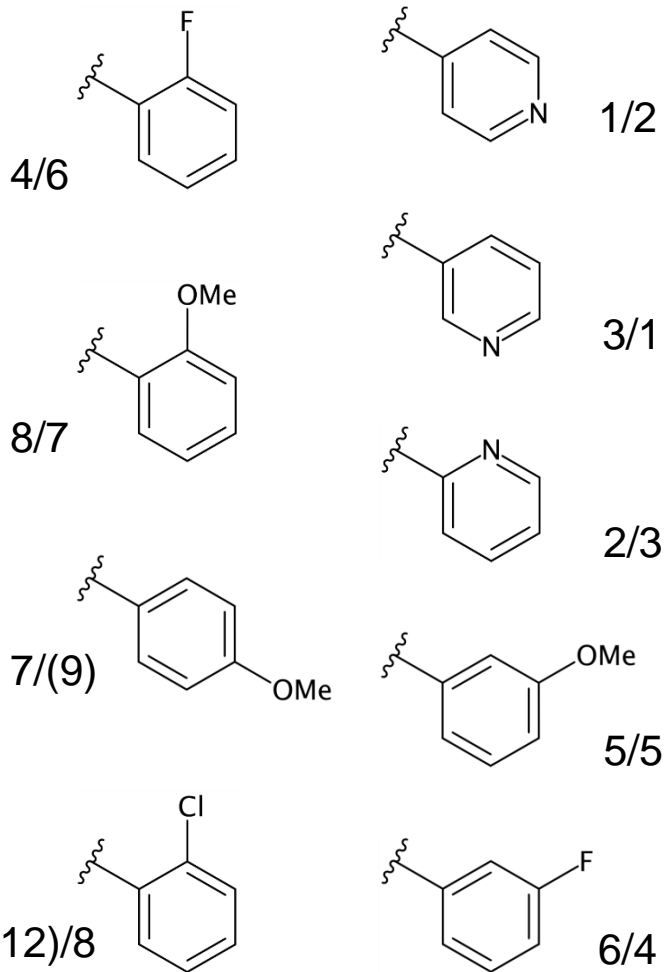
Less frequent

- Rank R groups based on co-occurrence

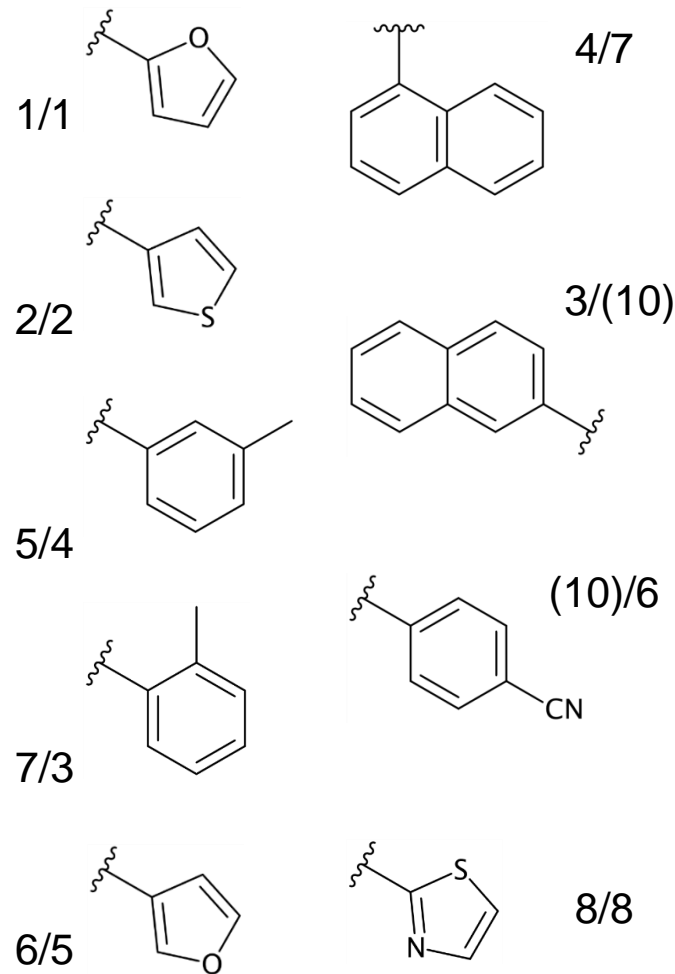
Find the 8 R groups most similar to



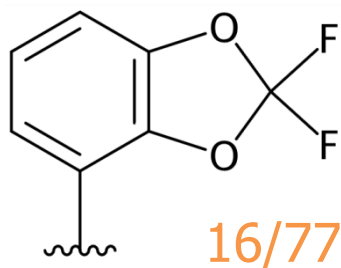
More frequent



Less frequent

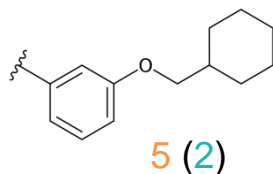
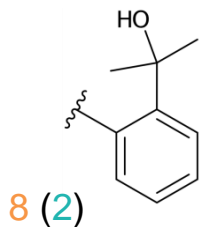
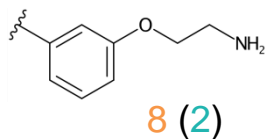
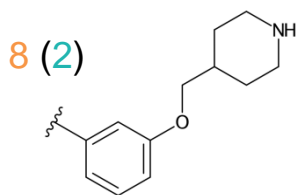
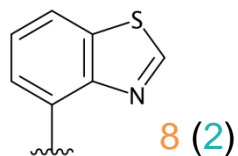
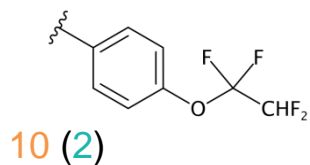


Find R groups similar to

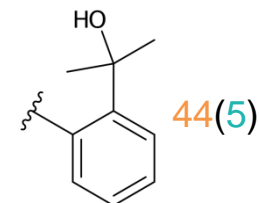
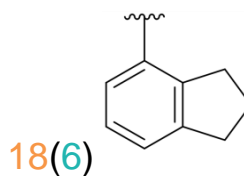
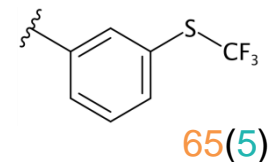
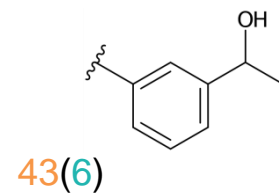
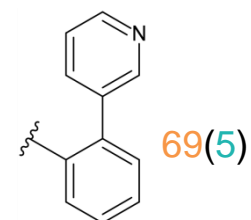
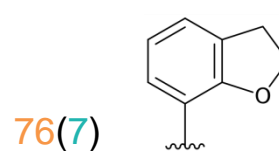


Color key: *frequency*
co-occurrence
Dataset key: ChEMBL/USPTO

Less frequent (ChEMBL)

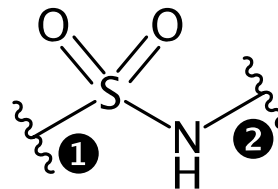


Less frequent (USPTO)



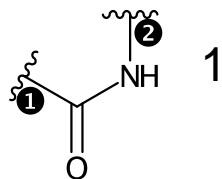
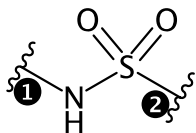
plus 200 more, 56 with
co-occurrence > 2

Find the 8 linkers most similar to



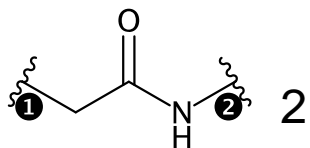
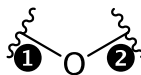
More frequent

5



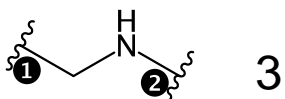
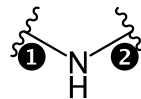
1

6



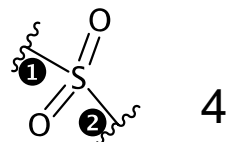
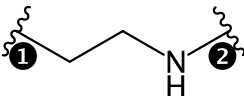
2

7



3

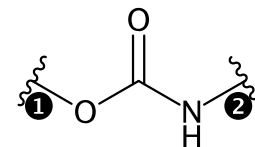
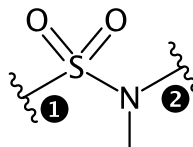
8



4

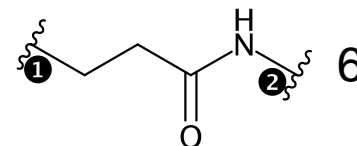
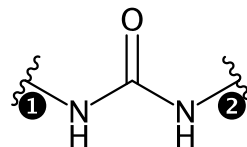
Less frequent

1



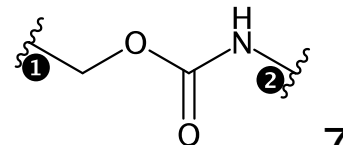
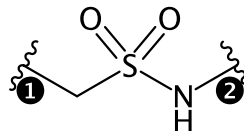
5

2



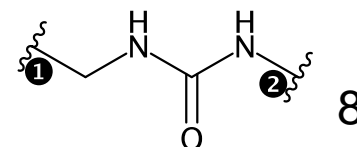
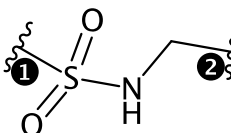
6

3



7

4



8

APPLICATIONS

- Ensure that gaps are covered
 - List R groups that should have been tested already
- Idea generation
 - List R groups that are commonly tested afterwards
- Enumerate
 - Generate molecules that are similar in R group space (and implicitly likely to be synthesisable)
- Similarity search...



SIMILARITY SEARCH



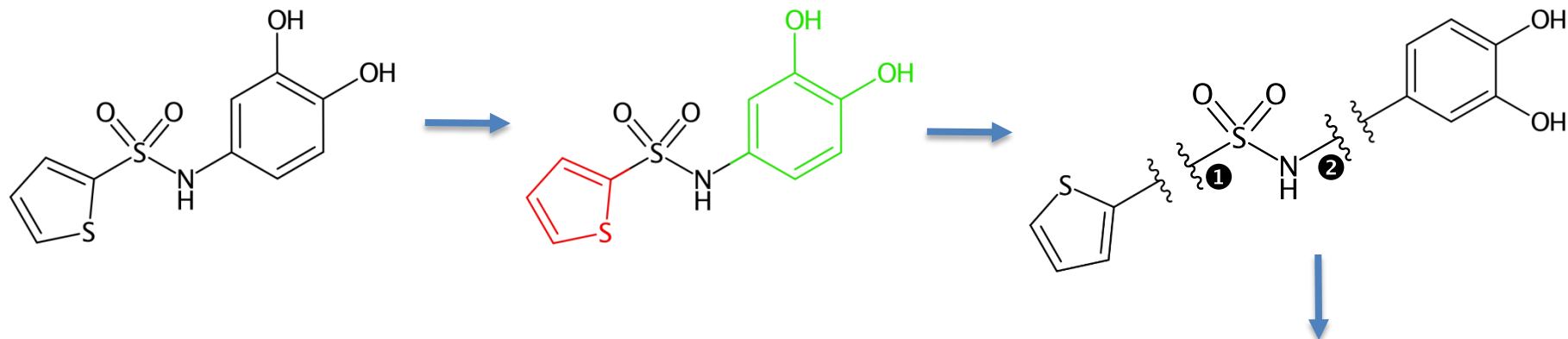
SIMILARITY SEARCH

- Are there molecules available for purchase that are similar in R group/linker space?
 - Worth purchasing and testing to complement in-house efforts
- Does ChEMBL contain bioactivity information for molecules that are similar in R group/linker space?
 - May affect patentability, or provide knowledge of off-target effects
- Are there exemplified compounds in patents* that are similar in R group/linker space?
 - May affect patentability, or provide knowledge of off-target effects

* For example, as extracted by LeadMine



FRAGMENTATION TO REDUCED GRAPHS

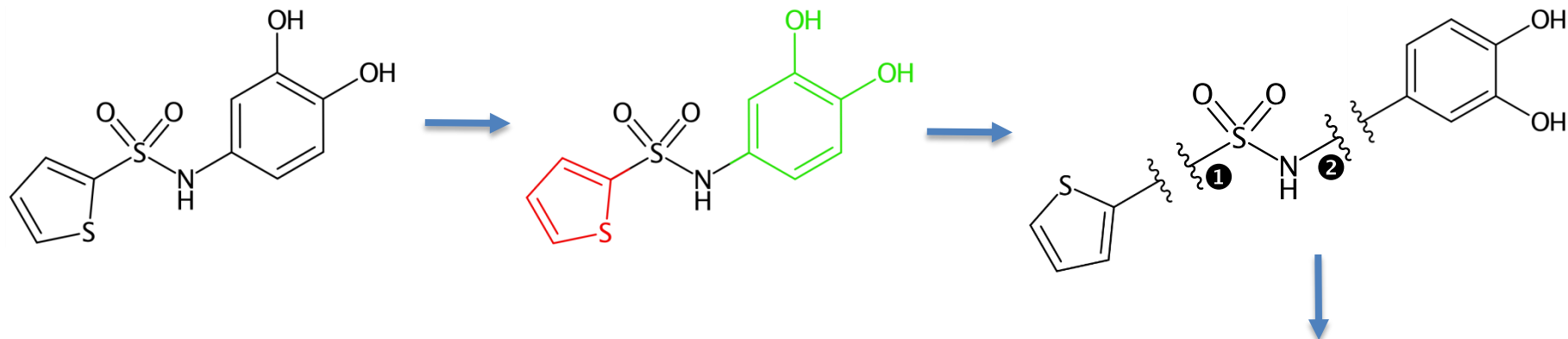


- Note that these reduced graphs:
 - Preserve the substitution position on rings
 - Preserve which end of the linker is attached to which R group

Birchall, K. and Gillet, V.J. (2011) Reduced graphs and their applications in chemoinformatics. In: Bajorath, J., (ed.) Chemoinformatics and Computational Chemical Biology. Methods in Molecular Biology, 672 . Humana Press , 197-212.
(http://eprints.whiterose.ac.uk/78616/9/Gillet_reduced-graphs-revised.pdf)



FRAGMENTATION TO REDUCED GRAPHS

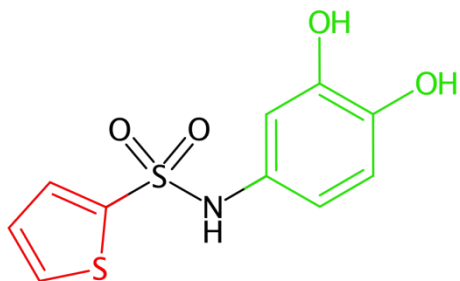


- Note that these reduced graphs:
 - Preserve the substitution position on rings
 - Preserve which end of the linker is attached to which R group
- Reduced graphs (RGs) are stored as SMILES strings
 - Isotopes indicate the node type: 1 for R group, 2 for linker, 10/11 indicate linker locants
 - Map indices indicate the specific identity of the nodes (via a lookup, e.g. R group 53 is 2-thiophenyl)

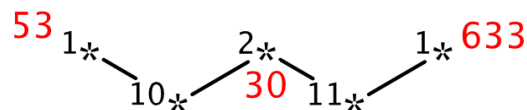
```
[1*:53][10*][2*:30][11*][1*:633]
```



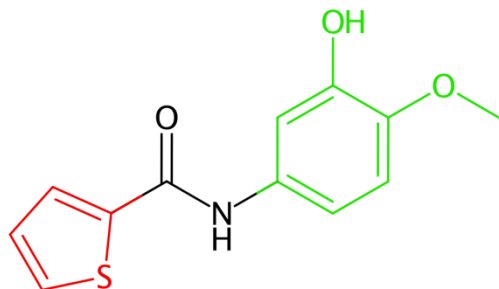
SIMILARITY SEARCH USING RGS



[1*:53] [10*] [2*:30] [11*] [1*:633]



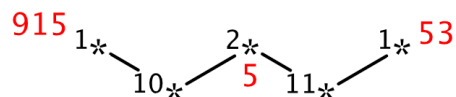
- Convert a database (e.g. ChEMBL) to reduced graphs
- For a given query:
 - Convert to reduced graph
 - Interpret as SMARTS and search the database
 - Score the hits based on distance to the query, a function of whether/where the matching groups appear on the query's list of nearest groups



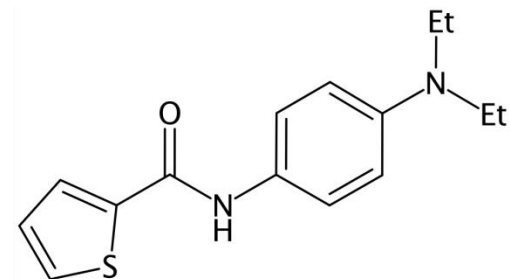
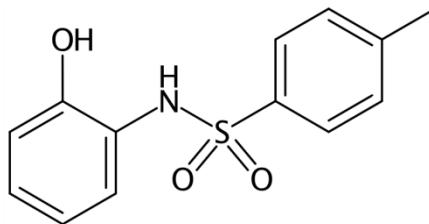
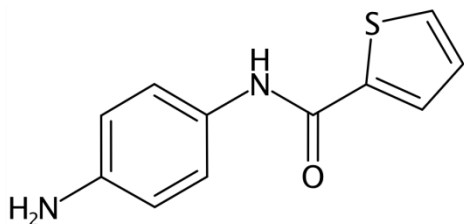
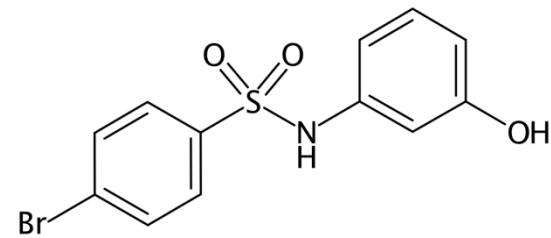
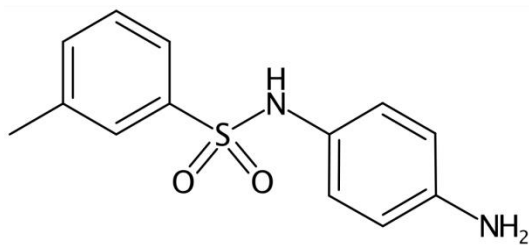
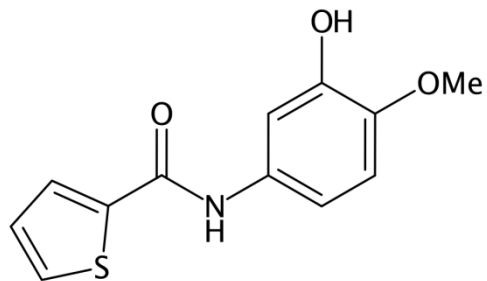
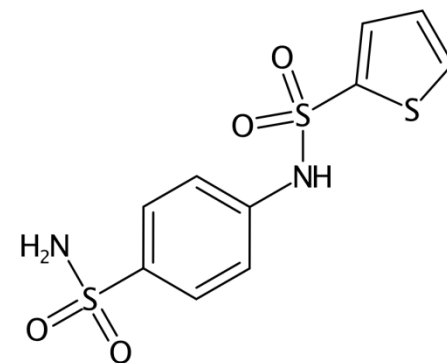
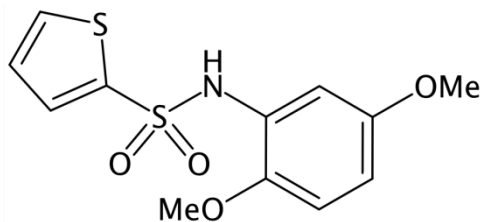
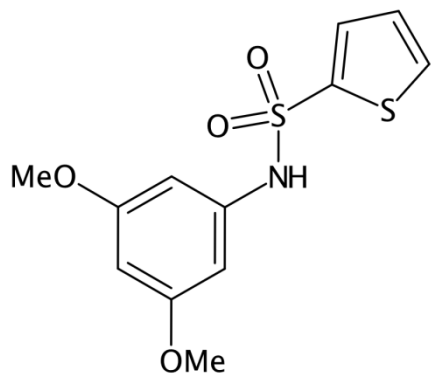
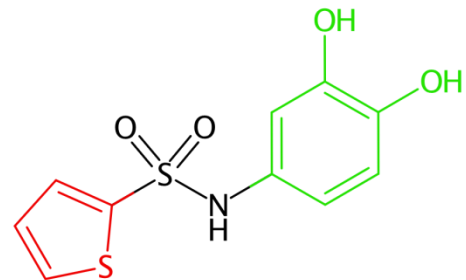
R group **53** is matched (does not contribute to distance)

R group **915** appears at rank 2 in 633's nearest (less frequent) R groups

Linker **5** appears at rank 1 in 30's nearest (more frequent) R groups



Search ChEMBL for molecules similar to



CONCLUSIONS



IN SUMMARY

- Co-occurrence of R groups in matched series from the same paper or patent can be used as an implicit measure of R group similarity
 - For less frequent R groups, patents provide more data
- Relative frequencies of R groups approximate the time dimension
 - Which R groups are typically made *after* versus *before*
- This approach can be used for idea generation, filling gaps, enumeration, and similarity search

