

## Supporting Information

# **RDChiral: An RDKit Wrapper for Handling Stereochemistry in Retrosynthetic Template Extraction and Application**

Connor W. Coley, William H. Green\*, and Klavs F. Jensen\*

*Department of Chemical Engineering, MIT, Cambridge, MA 02139*

E-mail: whgreen@mit.edu; kfjensen@mit.edu

## Code

All code can be found at <https://github.com/connorcoley/rdchiral> in addition to Jupyter notebooks containing all of the examples included in this manuscript.

# Additional Results

## Special Functional Group List

Table S1: List of special functional groups currently used for template extraction to ensure important motifs are included in the template. Atom IDs refers to the indices of the atoms in the SMARTS pattern that will trigger inclusion of that motif. For example, only the atom matching the wildcard in C#C-[\*] triggers inclusion of that motif; the alkynyl carbons will not trigger inclusion of their neighbors. This list is not meant to be exhaustive or definitive and is easily modified in the code.

SMARTS	Atom IDs	Description
<chem>[OH0,SH0]=C[O,C1,I,Br,F]</chem>	All	carboxylic acid / halogen
<chem>[OH0,SH0]=CN</chem>	All	amide/sulfamide
<chem>S(O)(O)[C1]</chem>	All	sulfonyl chloride
<chem>B(O)O</chem>	All	boronic acid/ester
<chem>[Si](C)(C)C</chem>	0	trialkyl silane
<chem>[Si](OC)(OC)(OC)</chem>	0	trialkoxo silane, default to methyl
<chem>[N;HO;\$ (N-[#6]) ;D2]-,=[N;D2]-,=[N;D1]</chem>	All	azide
<chem>O=C1N([Br,I,F,C1])C(=O)CC1</chem>	All	NBS-like reagent
<chem>Cc1ccc(S(=O)(=O)O)cc1</chem>	All	Tosyl
<chem>CC(C)(C)OC(=O)[N]</chem>	7	NBoc
<chem>[CH3][CH0]([CH3])([CH3])O</chem>	4	trimethylmethoxy
<chem>[C,N]=[C,N]</chem>	All	alkene/imine
<chem>[C,N]#[C,N]</chem>	All	alkyne/nitrile
<chem>C=C-[*]</chem>	2	adj to alkene
<chem>C#C-[*]</chem>	2	adj to alkyne
<chem>O=C-[*]</chem>	2	adj to carbonyl
<chem>O=C([CH3])-[*]</chem>	3	adj to methyl ketone
<chem>O=C([O,N])-[*]</chem>	3	adj to carboxylic acid/amide/ester
<chem>ClS(Cl)=O</chem>	All	thionyl chloride
<chem>[Mg,Li,Zn,Sn][Br,C1,I,F]</chem>	All	grinard/metal (non-disassociated)
<chem>S(O)(O)</chem>	All	SO2 group
<chem>N~N</chem>	All	incl. diazo
<chem>[!#6;R]@[#6;R]</chem>	1	adj to heteroatom in ring
<chem>[a!c]:a:a</chem>	2	two-steps from aromatic heteroatom
<chem>[B,C](F)(F)F</chem>	0	CF3, BF3 should have the F3 included

## Template Extraction

The ring closing reaction in [Figure S1A](#), as defined by the reaction SMILES, requires the trans cyclobutane in the precursor. The ketone in [Figure S1B](#) is prepared from the chiral alcohol, so the corresponding template specifies a chiral alcohol precursor. [Figure S1C](#) defines a reaction

between an alkyl iodide and the  $\beta$  position of an  $\alpha,\beta$ -unsaturated ketone. **Figure S1D** defines the preparation of a bromohydrin from a cis epoxide. And, finally, **Figure S1E** describes the stereoinversion of a chiral mesylate with an azide via an  $S_N2$  reaction.

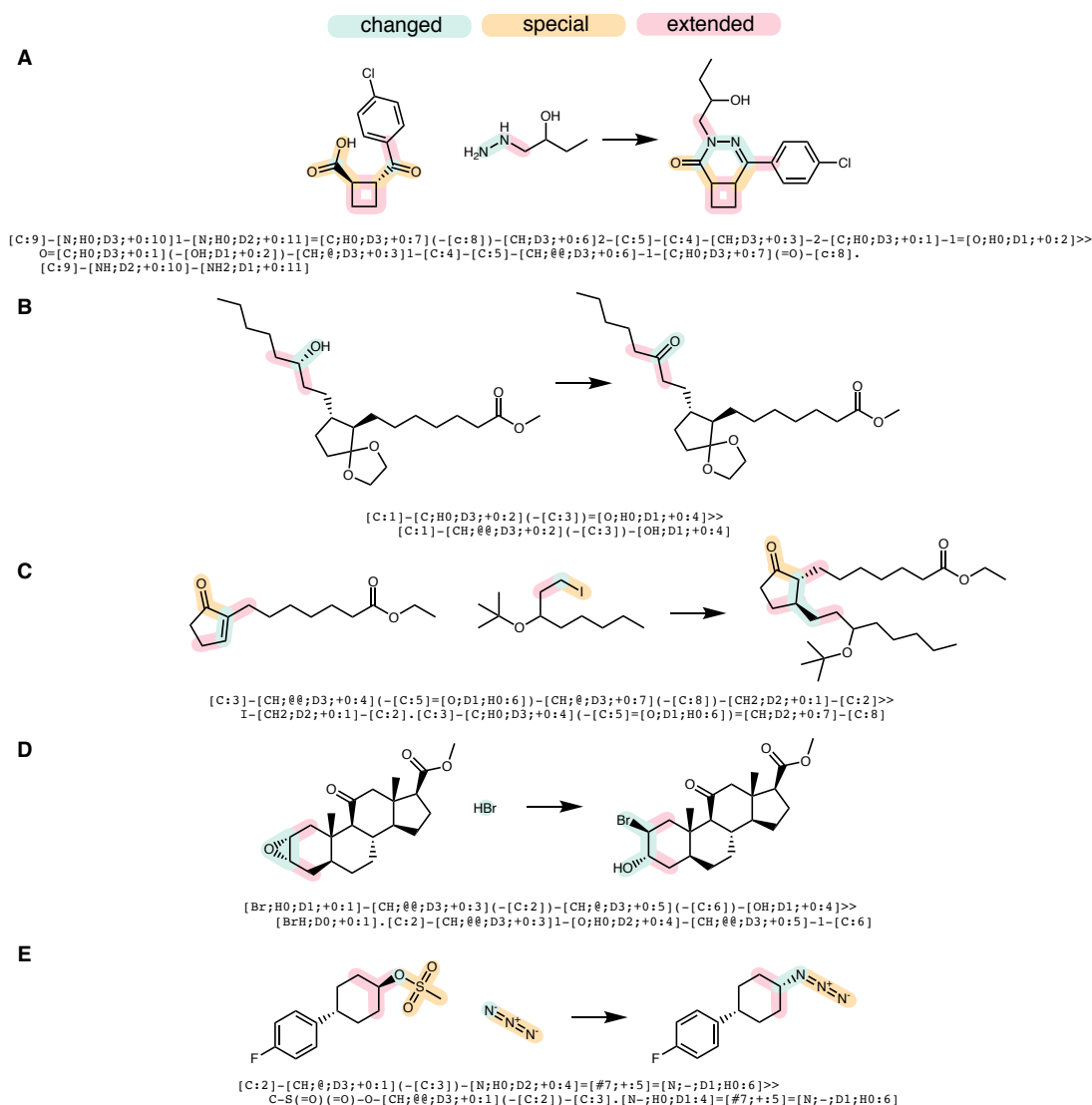


Figure S1: Additional examples of reactions from the USPTO and the retrosynthetic templates we extract from them, where stereochemistry must be considered. Atom mapping is omitted for brevity, but is unambiguous in all cases shown. Spectator molecules are also omitted, as they do not contribute heavy atoms to the product and are not included in the resulting reaction template.

## Data Quality Issues

The workflow relies on the availability of atom-mapped reaction SMILES strings. There are cases where the quality of examples can lead to nonsensical reaction templates. Two such examples are shown in [Figure S2](#). The first, shown in [Figure S2A](#), is a result of poor atom mapping. The reaction is a substitution of an alkyl chloride using sodium cyanide to prepare the alkyl cyanide. However, the solvent dimethylsulfoxide (DMSO) is labeled as contributing to the heterocycle in the product molecule. To the algorithm, all mapped atoms aside from atom 3 appear to undergo a change in connectivity, which causes the resulting template to include the entire product molecule.

The second example, shown in [Figure S2B](#), appears to be an erroneous reaction example entirely. It is likely that the unlabeled molecule was meant to be the reactant and the labeled reactant was meant to be the product of an oxidation. However, given the entry and its atom mapping, the apparent reaction is a shortening of the propyl side chain to ethyl. The corresponding template describes a reaction where any alkyl chain of at least two carbons can be synthesized from the same compound with that chain extended by a single carbon.

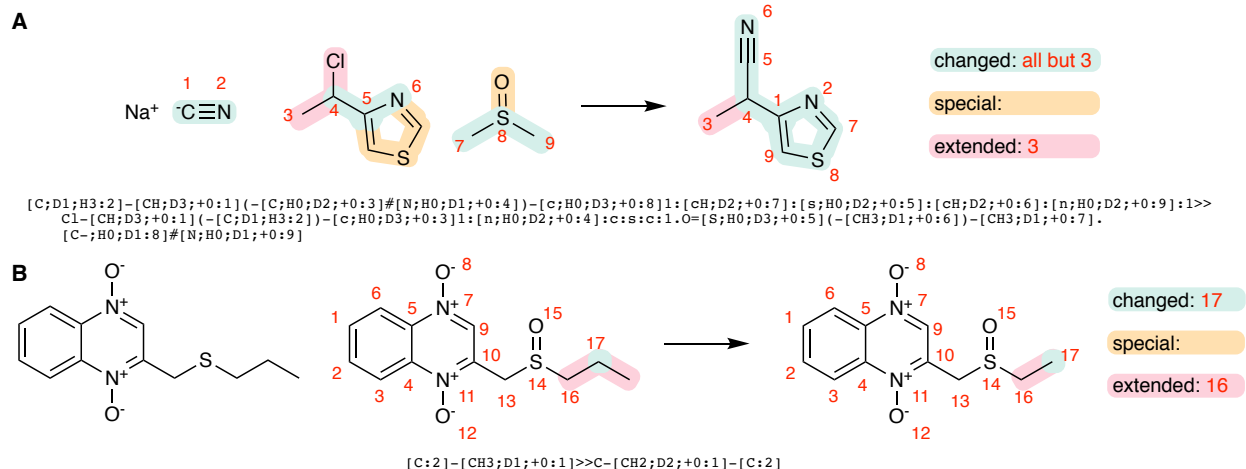


Figure S2: Reactions where the extracted template SMARTS is not chemically meaningful. (A) A case of poor atom mapping, where what should be a simple substitution reaction results in a template that encompasses the entire product structure. (B) A case of poor data quality, where the resulting template suggests that any alkyl chain of length  $n \geq 2$  can be prepared from an alkyl chain of length  $n + 1$ .