

PetroPalette: The Petro-Informatics Chemical Structure Database

Sydney Mazat and Benjamin Janesko*

Department of Chemistry & Biochemistry, Texas Christian University



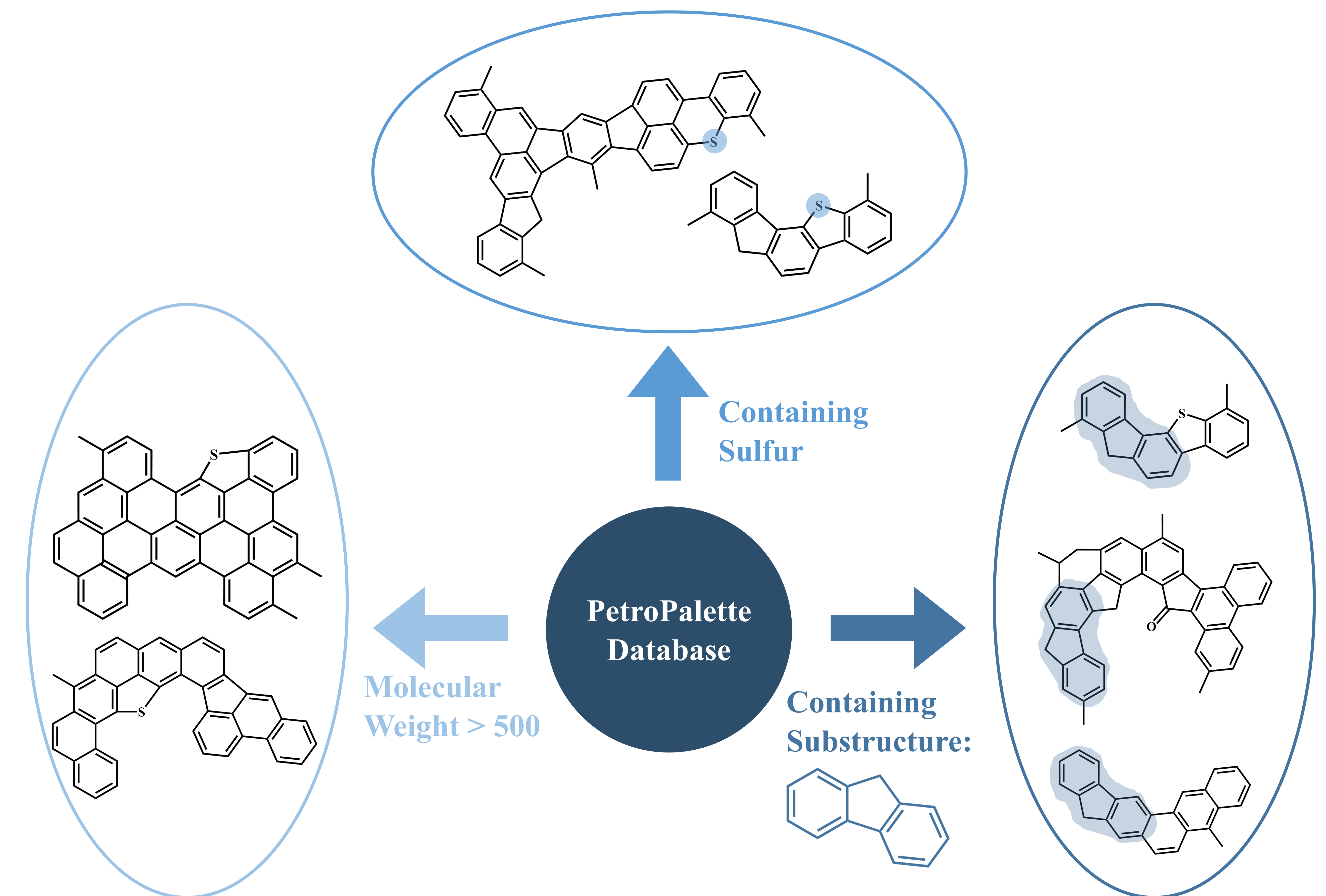
THIS PROJECT

Petroleum crude oil, unconventional crudes, and renewable bio-crudes are essential materials in our everyday lives. Crudes are highly complex chemical mixtures, estimated to contain between 100,000 and 100,000,000,000,000,000 unique molecules. Since 2015, single-molecule imaging has visualized hundreds of chemical structures, and historical literature has published thousands of proposed structures. This project builds an open database populated with published crude structures enabling data-driven analysis of these structures, and detailed workflows, allowing for easy future insertion of new molecules into the database.

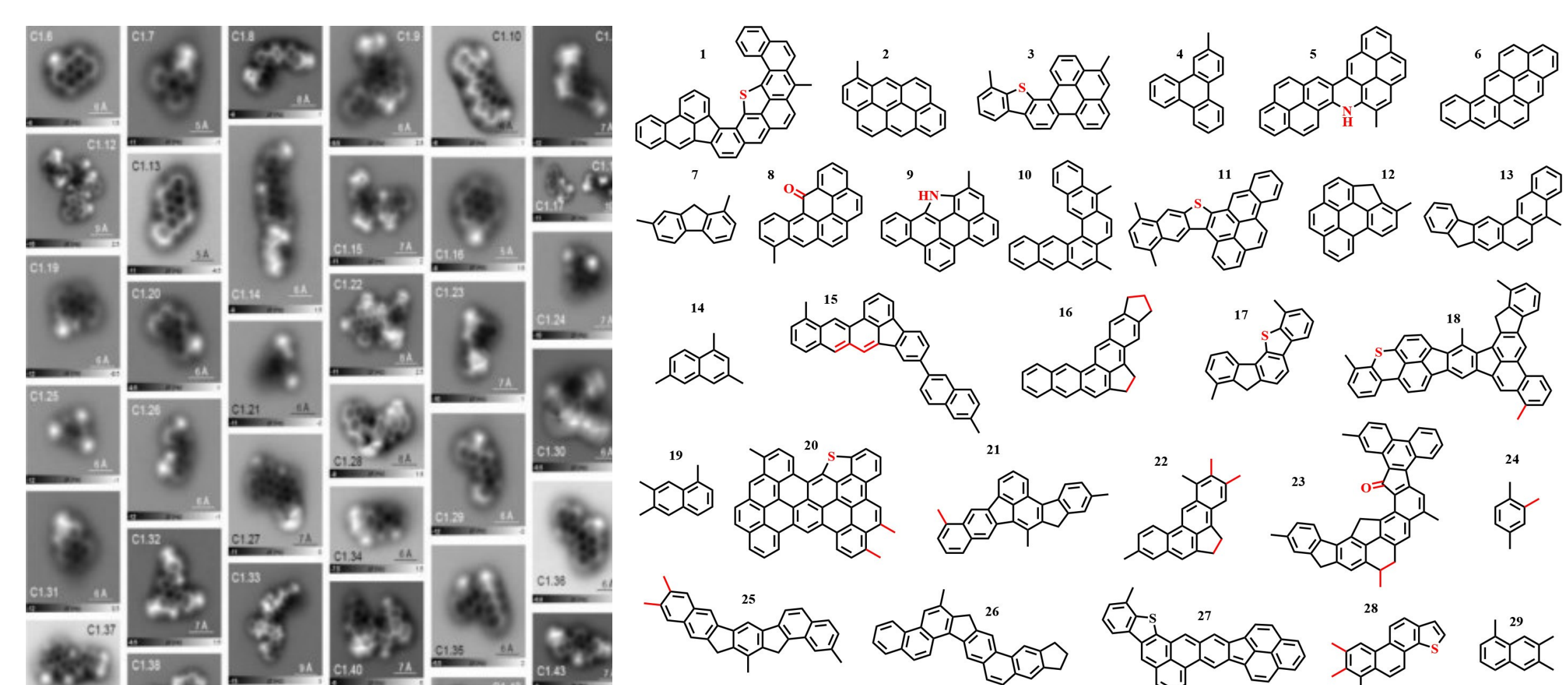
LONG TERM GOAL

This database can be used to make calculations and predict characteristics of molecules, such as viscosity, density, and reactivity, which are all critical in refinery plants, transportation, and usage of these fuels. Machine learning can be used to determine important characteristics of crudes molecules, leading to more refined and successful hypotheses. PetroPalette will eventually predict possible structures of petroleum crudes and their ensemble property distributions. The next step in this project is to add more real crudes into PetroPalette using the successful workflows, followed by machine learning to build predicted structures.

PETROPALETTE IN ACTION:

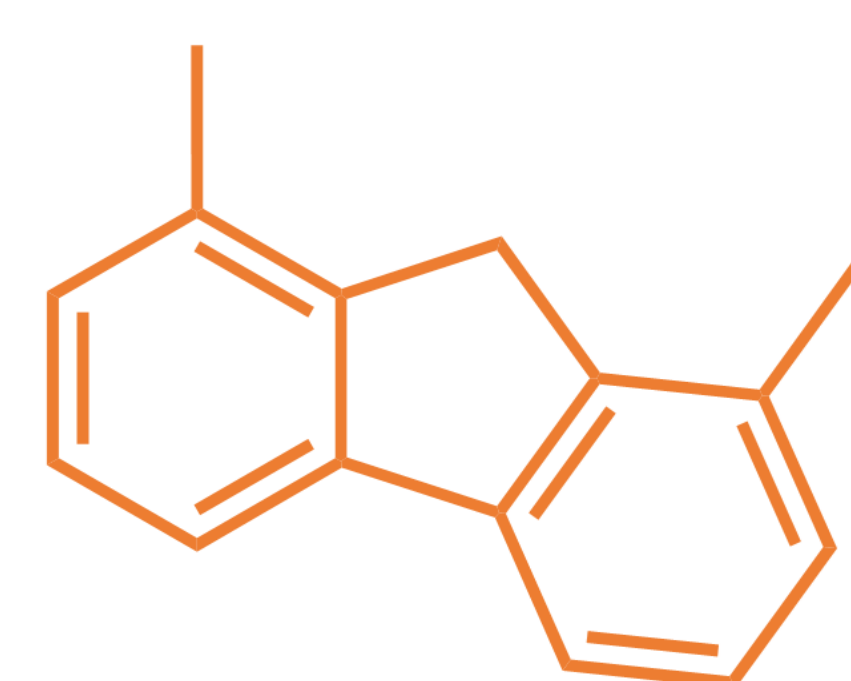


BUILDING PETROPALETTE:



1. Single Molecule Imaging produces chemical structures

2. Convert published structures into SMILES string



Cc1ccc2e(c1)Cc3c(C)cccc23

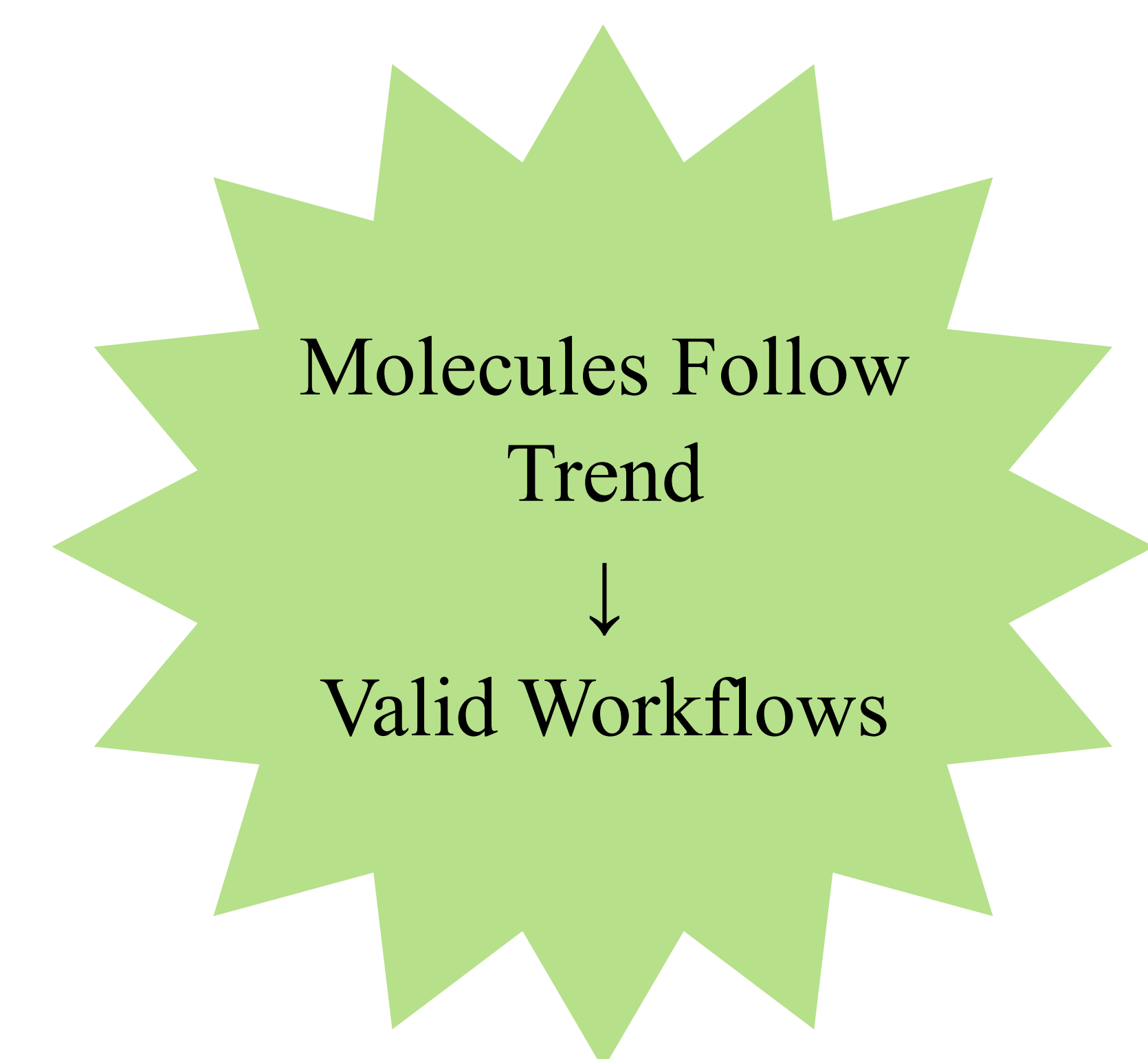
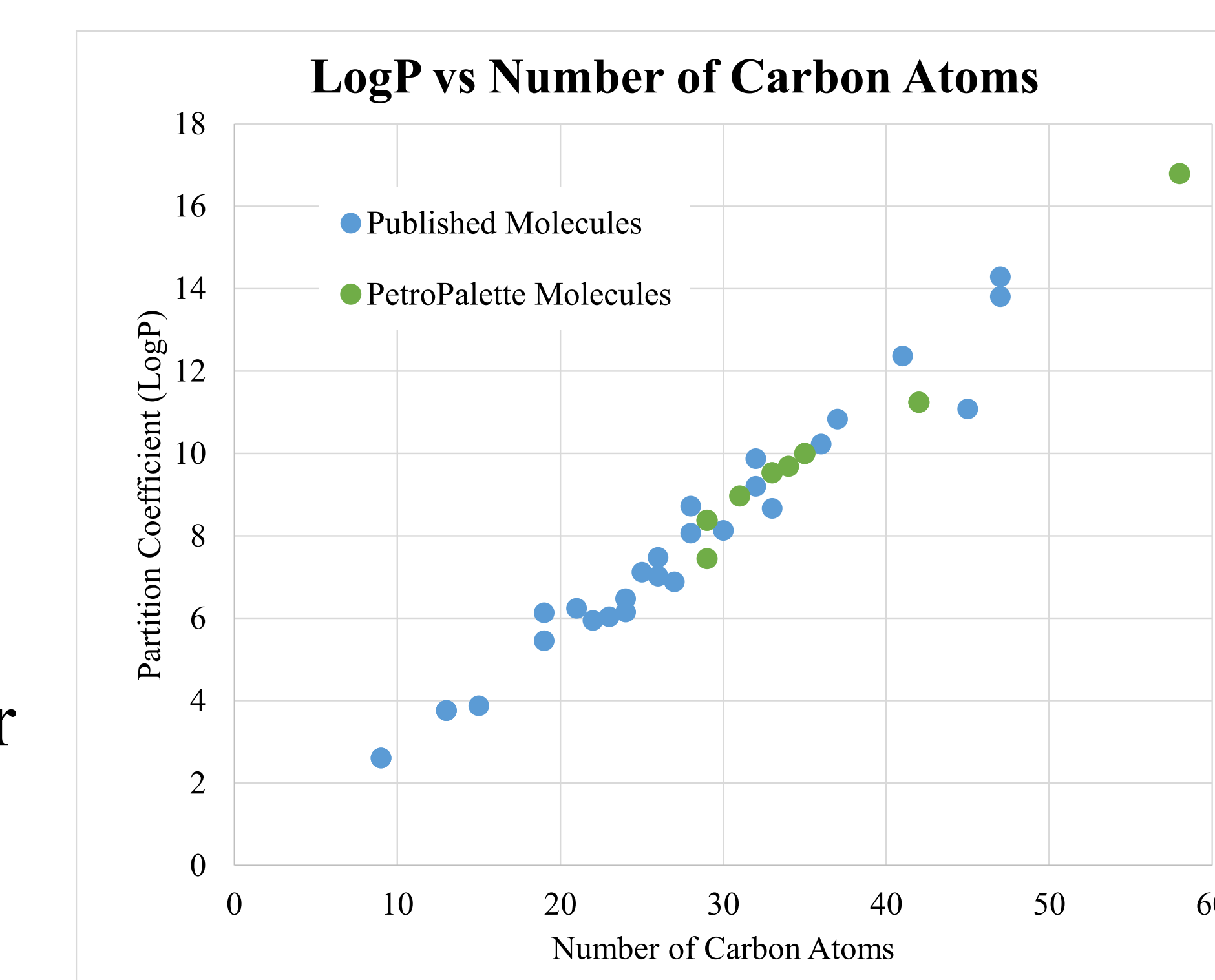
3. Produce descriptors from SMILES string

```
def descriptor_table(smiles, classification, doi):
    """Create lupac name, molecular weight, and molecular formula list from smiles"""
    lupac = []
    mw = []
    formula = []
    logp = []
    for sm in smiles:
        #lupac list
        compounds = pub.get_compounds(sm, namespace = 'smiles')
        match = compounds[0]
        lupac.append(match.lupac_name)
        molecular_weight_list =
        mol_wt = Chem.rdMolFromSmiles(sm)
        molwtmh = Chem.rdMolFromSmiles(sm)
        mass = descriptor.HolInt(molwtmh)
        mw.append(mass)
        #molecular formula list
        form = CalcMolFormula(molwtmh)
        formula.append(form)
        #logp list
        logp_value = crnp.HolLogP(molwtmh)
        logp.append(logp_value)
    """Create dataframe and add classification & doi source columns"""
    d = {'smiles_string': smiles,
        'lupac_name': lupac,
        'molecular_weight': mw,
        'molecular_formula': formula,
        'logp': logp}
    description = pd.DataFrame(data = d)
    description['classification'] = classification
    description['doi_source'] = doi
    return description
```

4. Insert structure and its descriptors into PetroPalette

smiles_string	IUPAC_name	molar_mass	molecular_formula	log_p	DQE	classification	doi_source
Cc1ccc2ccccc2c1	18-methyl-41-thiaund...	547	C41H22S	12.3673	31	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962
Cc1ccc2ccccc2c1	294	C23H18	6.04052	15	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	387	C28H18S	8.72194	20	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	242	C19H14	5.45662	13	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	415	C33H17N	9.2054	25	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	326	C28H14	7.4814	20	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	194	C19H14	3.87464	9	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	318	C24H14O	6.15222	18	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	329	C24H14N	7.11792	19	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	356	C28H20	8.06544	19	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	437	C32H20S	9.87514	23	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	
Cc1ccc2ccccc2c1	278	C22H14	5.94982	16	steam-cracked tar	https://doi.org/10.1021/acs.iecr.8b03962	

This figure compares the Partition Coefficient value of published crudes¹ (blue) to that of a fully automated analysis of other crudes² (green) using the workflows and PetroPalette. These values following the linear trend help to validate the workflows.



Acknowledgement: This work is supported by the American Chemical Society Petroleum Research Fund

1. *Ind. Eng. Chem. Res.* 2018, 57, 46, 15935–15941

2. *Energy & Fuels* 2022, 36, 16, 8714–8724

Contact: s.mazat@tcu.edu