

Constructing a database of asphaltenes: Quantum chemistry to contextualize single molecule experiments within ensemble properties

Gretel Stokes, Sydney Mazat, Benjamin Janesko*

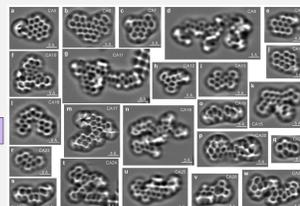
Texas Christian University, Fort Worth, Texas

Introduction

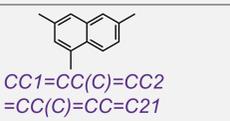
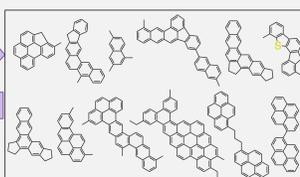
Asphaltene are the heaviest and most complex constituents of crude oil, with thousands of distinct species. Despite their structural diversity, which hampers complete understanding, modern studies have identified many asphaltene structures using atomic force microscopy. To address the challenge of storing and analyzing this information, we've created a database of 67 published asphaltene structures. Quantum chemistry calculations provide molecular properties such as weight, solubility, aromaticity, dipole moment, and HOMO-LUMO gap, which are stored in the database. Using this data, we generate graphs, like UV-visible absorbance spectra, to offer a comprehensive chemical description of asphaltene mixtures. Our computational predictions enhance understanding of individual asphaltene structures and their relationship to ensemble properties in crude oil.

Constructing the Database

1. Single molecule imaging studies elucidate asphaltene chemical structures from three datasets: SCT28, CA12, and P27, where the number denotes the quantity of molecules in each dataset



2. The published molecules are converted to their SMILES strings



3. Descriptors are generated from SMILES strings

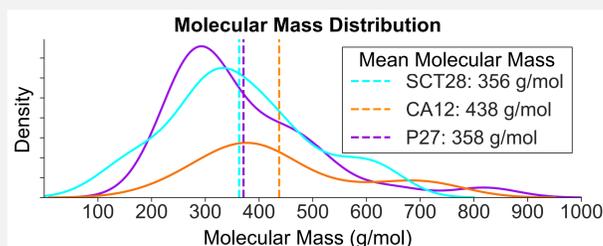
```
# create supac name, molecular weight, and molecular formula list from smiles
supac = []
for mol in mol_list:
    formula = []
    mol_weight = []
    mol_formula = []
    for mol in mol_list_SCT28:
        # smiles list
        compounds = pub.get_compounds(smiles, namespace='smiles')
        match = compounds[0]
        supac.append(match.supac_name)
        mol_weight.append(match.mol_weight)
        mol_formula.append(match.mol_formula)
    # smiles list
    for mol in mol_list_CA12:
        compounds = pub.get_compounds(smiles, namespace='smiles')
        match = compounds[0]
        supac.append(match.supac_name)
        mol_weight.append(match.mol_weight)
        mol_formula.append(match.mol_formula)
    # smiles list
    for mol in mol_list_P27:
        compounds = pub.get_compounds(smiles, namespace='smiles')
        match = compounds[0]
        supac.append(match.supac_name)
        mol_weight.append(match.mol_weight)
        mol_formula.append(match.mol_formula)
# Generate lists with number of N, H, C, number of rings,
# LogP values, DBE, Z, and dipole moment columns for each molecule.
number_of_H = []
number_of_C = []
number_of_N = []
rings = []
logP = []
DBE = []
dipole_moment = []
HOMO_LUMO_gap = []
for mol in mol_weight:
```

4. Properties are stored in an SQL database

smiles_string	IUPAC_name	molar_mass	molecular_weight	number_of_H	number_of_C	number_of_N	number_of_rings	LogP	DBE	Z	dipole_moment	HOMO_LUMO_gap
CC1=CC=CC=C1	1H-benzene	78.11	78.11	6	6	0	1	2.13	0	6	0.000	0.000
CC1=CC=CC=C1C	1H-naphthalene	128.17	128.17	10	10	0	2	4.15	1	10	0.000	0.000
CC1=CC=CC=C1C=C1	1H-fluorene	166.20	166.20	14	14	0	3	5.99	2	14	0.000	0.000
CC1=CC=CC=C1C=CC1	1H-acenaphthylene	152.15	152.15	12	12	0	3	5.00	2	12	0.000	0.000
CC1=CC=CC=C1C=CC=C1	1H-acenaphthene	166.20	166.20	14	14	0	3	5.99	2	14	0.000	0.000
CC1=CC=CC=C1C=CC=C1C	1H-benzofluorene	202.23	202.23	18	18	0	4	7.92	3	18	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=C1	1H-benzofluorene	216.25	216.25	20	20	0	4	8.94	3	20	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC1	1H-benzofluorene	230.27	230.27	22	22	0	4	9.96	3	22	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	244.29	244.29	24	24	0	4	10.98	3	24	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	258.31	258.31	26	26	0	4	12.00	3	26	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	272.33	272.33	28	28	0	4	13.02	3	28	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	286.35	286.35	30	30	0	4	14.04	3	30	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	300.37	300.37	32	32	0	4	15.06	3	32	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	314.39	314.39	34	34	0	4	16.08	3	34	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	328.41	328.41	36	36	0	4	17.10	3	36	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	342.43	342.43	38	38	0	4	18.12	3	38	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	356.45	356.45	40	40	0	4	19.14	3	40	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	370.47	370.47	42	42	0	4	20.16	3	42	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	384.49	384.49	44	44	0	4	21.18	3	44	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	398.51	398.51	46	46	0	4	22.20	3	46	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	412.53	412.53	48	48	0	4	23.22	3	48	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	426.55	426.55	50	50	0	4	24.24	3	50	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	440.57	440.57	52	52	0	4	25.26	3	52	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	454.59	454.59	54	54	0	4	26.28	3	54	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	468.61	468.61	56	56	0	4	27.30	3	56	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	482.63	482.63	58	58	0	4	28.32	3	58	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	496.65	496.65	60	60	0	4	29.34	3	60	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	510.67	510.67	62	62	0	4	30.36	3	62	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	524.69	524.69	64	64	0	4	31.38	3	64	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	538.71	538.71	66	66	0	4	32.40	3	66	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	552.73	552.73	68	68	0	4	33.42	3	68	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	566.75	566.75	70	70	0	4	34.44	3	70	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	580.77	580.77	72	72	0	4	35.46	3	72	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	594.79	594.79	74	74	0	4	36.48	3	74	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	608.81	608.81	76	76	0	4	37.50	3	76	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	622.83	622.83	78	78	0	4	38.52	3	78	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	636.85	636.85	80	80	0	4	39.54	3	80	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	650.87	650.87	82	82	0	4	40.56	3	82	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	664.89	664.89	84	84	0	4	41.58	3	84	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	678.91	678.91	86	86	0	4	42.60	3	86	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	692.93	692.93	88	88	0	4	43.62	3	88	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	706.95	706.95	90	90	0	4	44.64	3	90	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	720.97	720.97	92	92	0	4	45.66	3	92	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC1	1H-benzofluorene	734.99	734.99	94	94	0	4	46.68	3	94	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1	1H-benzofluorene	749.01	749.01	96	96	0	4	47.70	3	96	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C	1H-benzofluorene	763.03	763.03	98	98	0	4	48.72	3	98	0.000	0.000
CC1=CC=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=CC=C1C=C1	1H-benzofluorene	777.05	777.05	100	100	0	4	49.74	3	100	0.000	0.000

Exploiting the Database

Mass Spectra: Asphaltene exhibit molecular weights ranging from several hundred to several thousand g/mol. A plot that estimates the molecular mass distribution for each of the three datasets was generated, with masses recorded in the database being on the lower end of the spectrum.

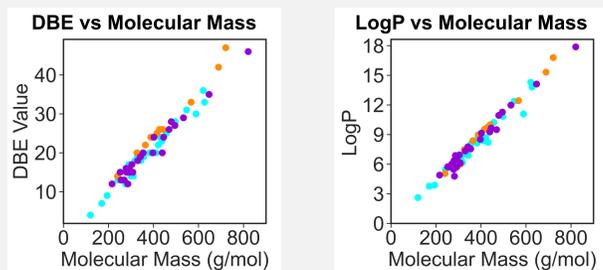


Double bond equivalents (DBE): a measure of unsaturation in organic molecules calculated by:

$$DBE = N_C - \frac{N_H}{2} + \frac{N_N}{2} + 1$$

For asphaltenes, DBE offers insight into their structural complexity and reactivity due to the presence of fused aromatic rings and aliphatic chains.

LogP: the logarithm of a molecule's octanol:water partition coefficient, which indicates relative solubility. Plotting logP against asphaltene molecular mass shows larger molecules are more hydrophobic, consistent with previous research. Predicted logP values span a wide range, signifying substantial variation in water solubility among asphaltenes.

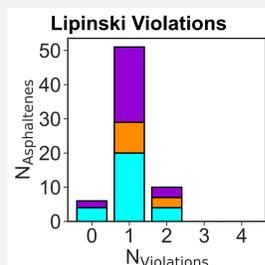


Dataset Legend
• SCT28 • CA12 • P27

Druglikeness: Lipinski's Rule of Five is widely used in pharmaceuticals to gauge oral availability. For a molecule to be considered orally available, it shouldn't violate more than one of these criteria:

1. No more than five hydrogen bond donors
2. No more than ten hydrogen bond acceptors
3. Molecular mass less than 500 daltons
4. LogP not greater than five

The number of violations of each asphaltene was computed. Most had one or fewer violations, suggesting they might be suitable drug candidates for oral delivery.



UV-Vis Spectral Modeling

The UV-Visible absorbance spectrum of the asphaltene was modeled using computed values of the HOMO-LUMO gap and transition dipole moment.

HOMO-LUMO gap (E_{HL}):

- Energy in eV is converted to corresponding wavelength
- Determines wavelength where absorption is the highest

Transition dipole moment (μ):

- Determines intensity of absorption peak.

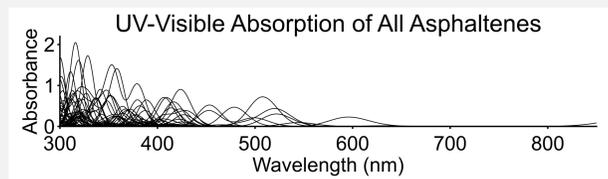
Thus, the absorption, A , is modeled by a Gaussian broadening function:

$$A(E) = \mu e^{-(E-E_{HL})^2/\sigma^2}$$

Where:

- E is the photon energy over which UV-visible spectra are plotted
- σ is a broadening factor, chosen to be 0.05, which determines how sharp the peaks appear.

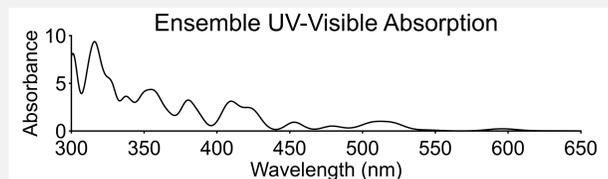
The following graph gives the plotted absorbance of each of the 67 asphaltene molecules:



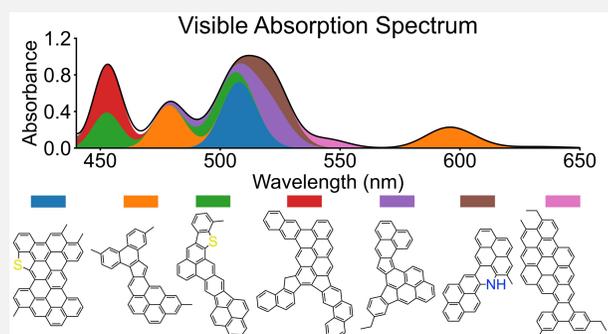
Summing the absorption and plotting against photon energy gives the following ensemble spectrum:

$$\sum_{n=1}^{67} A_n$$

Where A_n is the absorbance of each asphaltene



Most asphaltene absorb in the 300 nm – 400 nm range. Focusing on the visible spectrum, we illustrate the spectral absorption of the seven most significant asphaltene structures.



Conclusions

Our study illuminates the complexity of asphaltene in crude oil. Through computational analysis, we've cataloged 67 asphaltene structures in a SQL database and elucidated their predicted molecular properties. Our findings reveal insights into properties of asphaltene such as molecular weight distribution, unsaturation, solubility, and drug-likeness, with many adhering to Lipinski's Rule of Five. By modeling UV-visible absorbance spectra based on molecular characteristics, we enhance understanding of asphaltene behavior. Our interdisciplinary approach integrates experimental data and computational modeling, offering valuable tools for further research in petroleum science.

Future Directions

Asphaltene mixtures are structurally complex, containing thousands of distinct structures. The upcoming phase involves employing machine learning techniques to generate predictive structures based on the properties of real asphaltene stored in the database, advancing our understanding of crude oil composition and behavior.



Acknowledgements

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

References

1. Schuler, B.; Meyer, G.; Peña, D.; Mullins, O. C.; Gross, L. Unraveling the Molecular Structures of Asphaltene by Atomic Force Microscopy. *Journal of the American Chemical Society* **2015**, *137* (31), 9870-9876. DOI: 10.1021/jacs.5b04056.
2. Chen, P.; Metz, J. N.; Mennito, A. S.; Merchant, S.; Smith, S. E.; Siskin, M.; Rucker, S. P.; Dankworth, D. C.; Kushnerick, J. D.; Yao, N.; et al. Petroleum pitch: Exploring a 50-year structure puzzle with real-space molecular imaging. *Carbon* **2020**, *161*, 456-465. DOI: <https://doi.org/10.1016/j.carbon.2020.01.062>.
3. Zhang, Y.; Schuler, B.; Fatayer, S.; Gross, L.; Harper, M. R.; Kushnerick, J. D. Understanding the Effects of Sample Preparation on the Chemical Structures of Petroleum Imaged with Noncontact Atomic Force Microscopy. *Industrial & Engineering Chemistry Research* **2018**, *57* (46), 15935-15941. DOI: 10.1021/acs.iecr.8b03962.

Contact Information

gretel.jordan@tcu.edu
b.janesko@tcu.edu