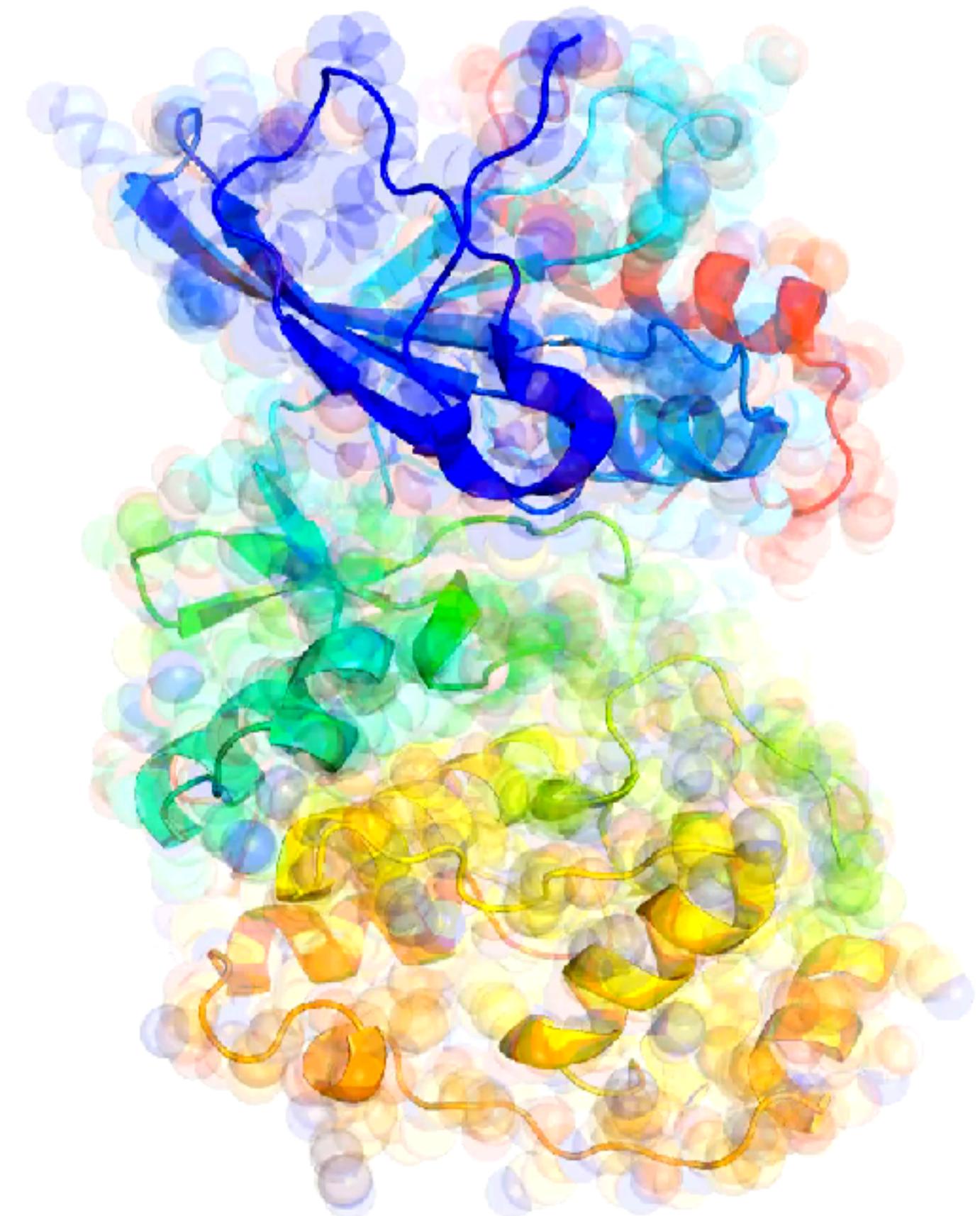


# Improving **molecular models** by generating **high-quality** quantum chemistry data

---

Chaya D. Stern, Jeffrey Wagner, Daniel G. A. Smith, Christopher Bayly, Andrea Rizzi, John D. Chodera  
RDKit User Group Meeting  
Sep 25, 2019

# Molecular dynamics simulations provide atomic detail to dynamics that is usually difficult to observe

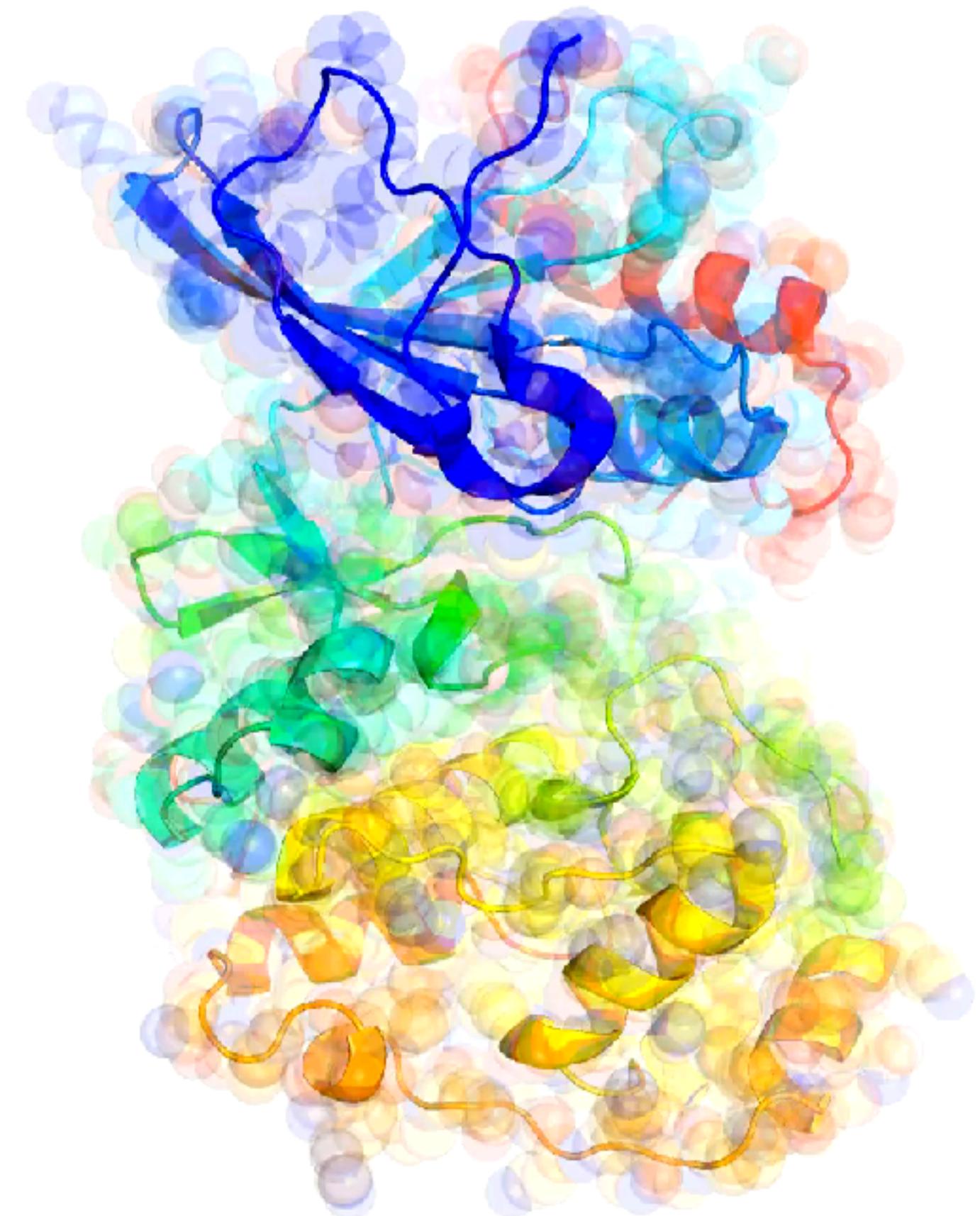


But more importantly, we can use MD simulations to make quantitative predictions.

**Rafal  
Wiewiora**



# Molecular dynamics simulations provide atomic detail to dynamics that is usually difficult to observe



But more importantly, we can use MD simulations to make quantitative predictions.

**Rafal  
Wiewiora**

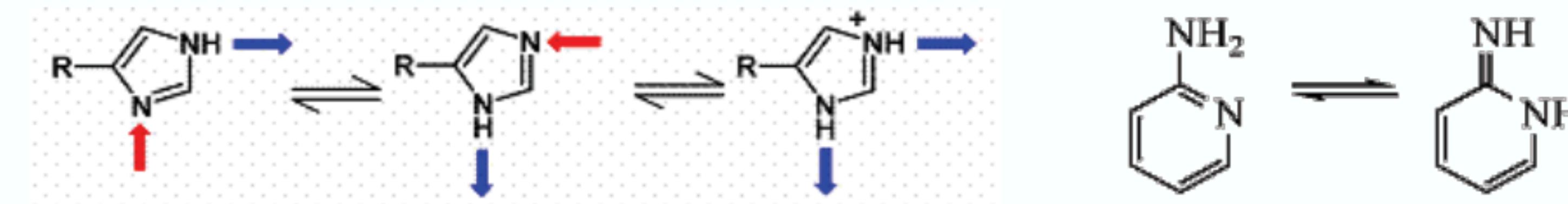


# Predictions fail for three main reasons.

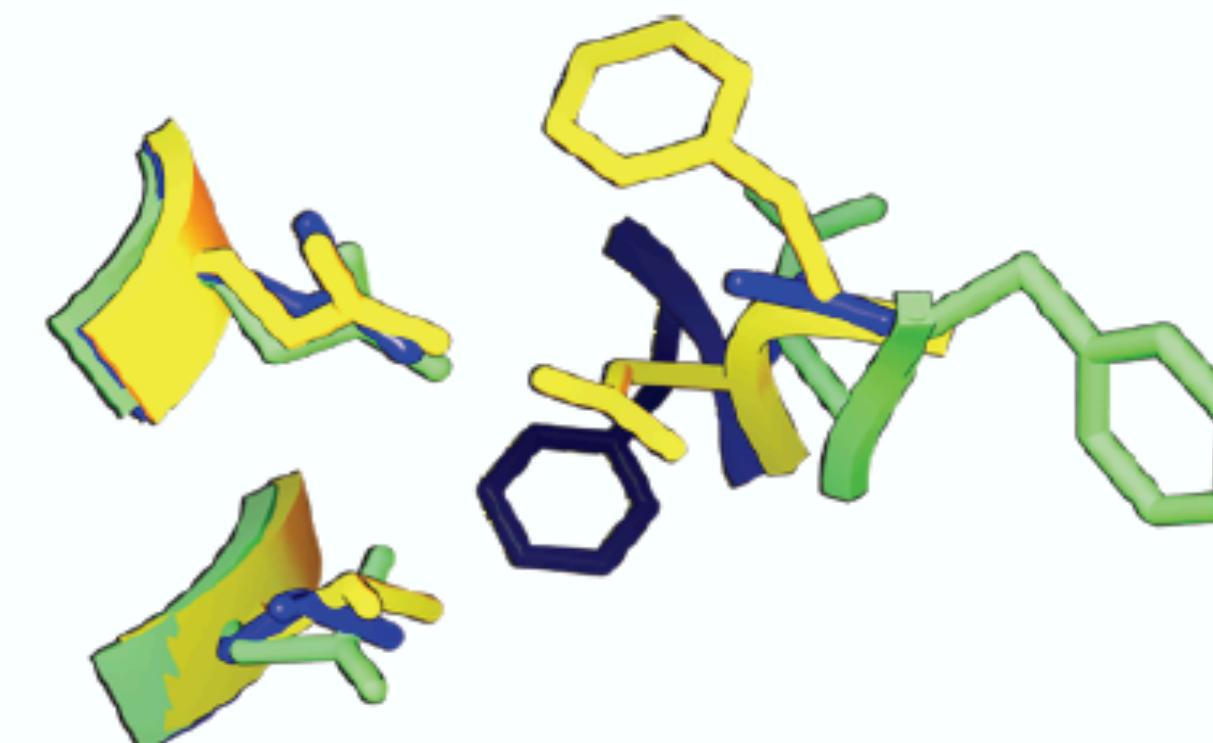
1. The **forcefield** does a poor job of modeling the physics of our system

$$V(\mathbf{q}) = \sum_{\text{bonds}} K_r(r - r_{eq})^2 + \sum_{\text{angles}} K_\theta(\theta - \theta_{eq})^2 + \sum_{\text{dihedrals}} \frac{V_n}{2}[1 + \cos(n\phi - \gamma)] + \sum_{i < j} \left[ \frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^6} + \frac{q_i q_j}{\epsilon R_{ij}} \right]$$

2. We're missing some **essential chemical** in our simulations  
(e.g. protonation states, tautomers, covalent association)



3. We haven't **sampled** all of the relevant conformations

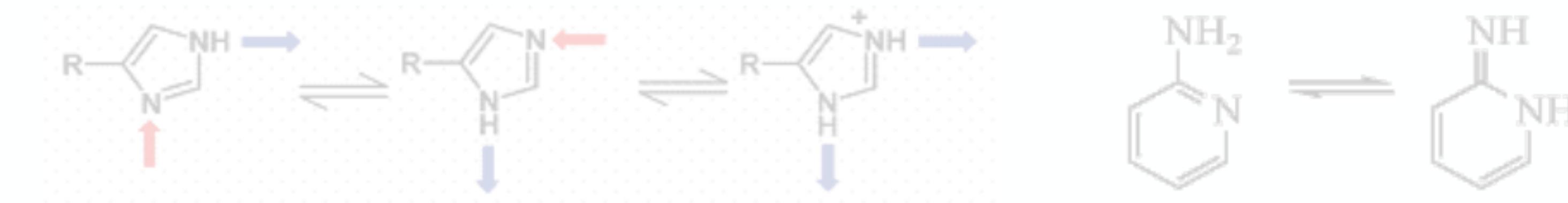


# Predictions fail for three main reasons.

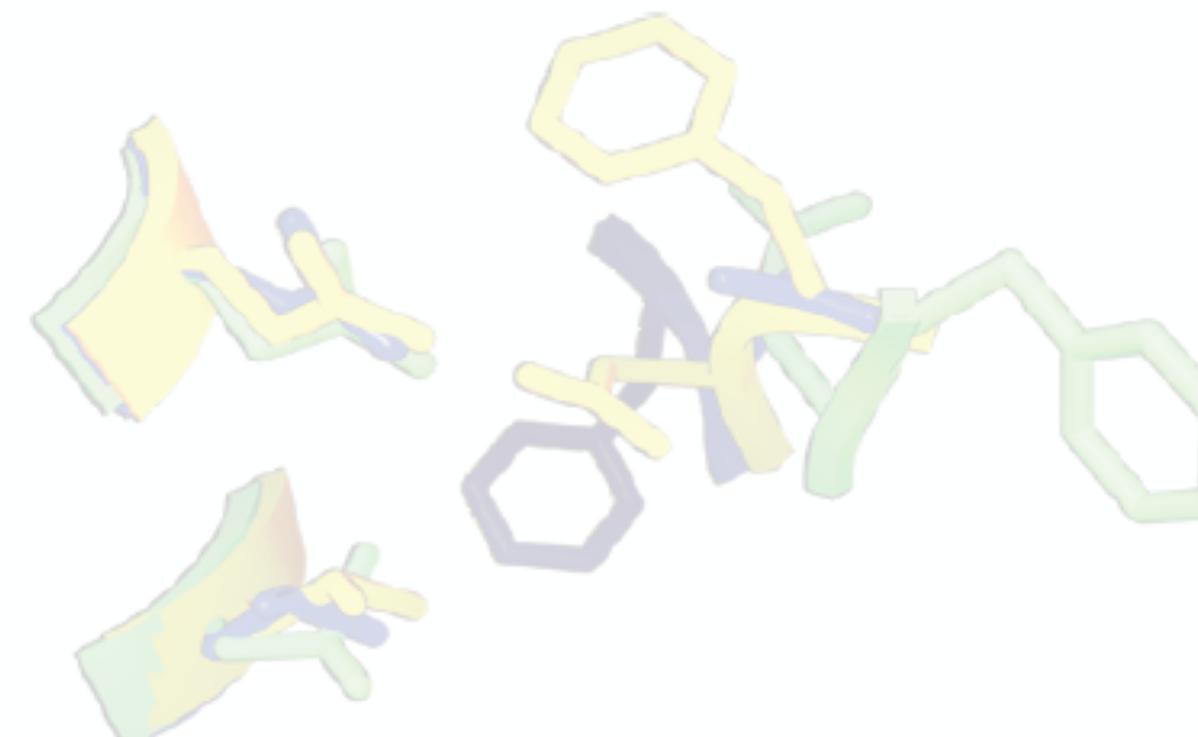
1. The **forcefield** does a poor job of modeling the physics of our system

$$V(\mathbf{q}) = \sum_{\text{bonds}} K_r(r - r_{eq})^2 + \sum_{\text{angles}} K_\theta(\theta - \theta_{eq})^2 + \sum_{\text{dihedrals}} \frac{V_n}{2}[1 + \cos(n\phi - \gamma)] + \sum_{i < j} \left[ \frac{A_{ij}}{R_{ij}^{12}} - \frac{B_{ij}}{R_{ij}^6} + \frac{q_i q_j}{\epsilon R_{ij}} \right]$$

2. We're missing some **essential chemical** in our simulations  
(e.g. protonation states, tautomers, covalent association)



3. We haven't **sampled** all of the relevant conformations



# INDUSTRY



**Jeff Wagner**



**Daniel Smith**



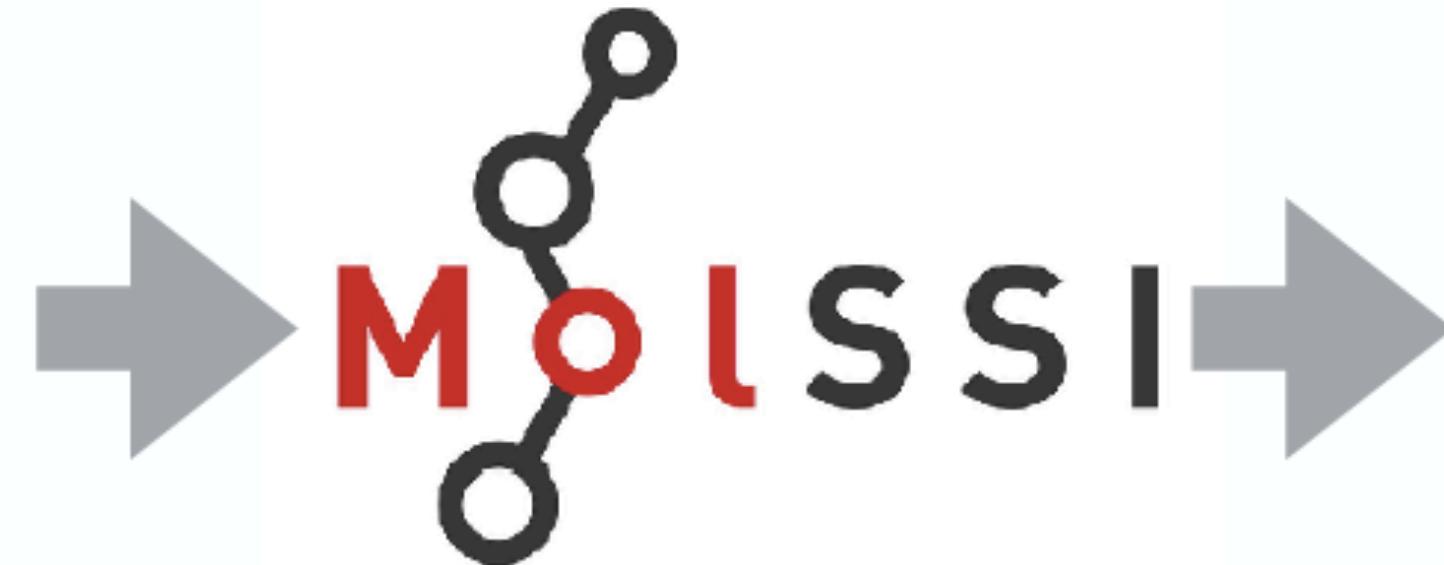
**CHRISTOPHER BAYLY**  
**OPENEYE SCIENTIFIC**



**KENNETH KROENLEIN**  
**NIST THERMODYNAMICS RESEARCH CENTER**  
(NIST is a US federal agency)

# The Open Forcefield Consortium

COORDINATING INTERMEDIARY



MOLECULAR SOFTWARE  
SCIENCES INSTITUTE  
coordination of funding  
while minimizing indirect costs



**JOHN CHODERA**

SLOAN KETTERING INSTITUTE



**MICHAEL GILSON**

UNIVERSITY OF CALIFORNIA, SAN DIEGO



**DAVID MOBLEY**

UNIVERSITY OF CALIFORNIA, IRVINE



**MICHAEL SHIRTS**

UNIVERSITY OF COLORADO, BOULDER



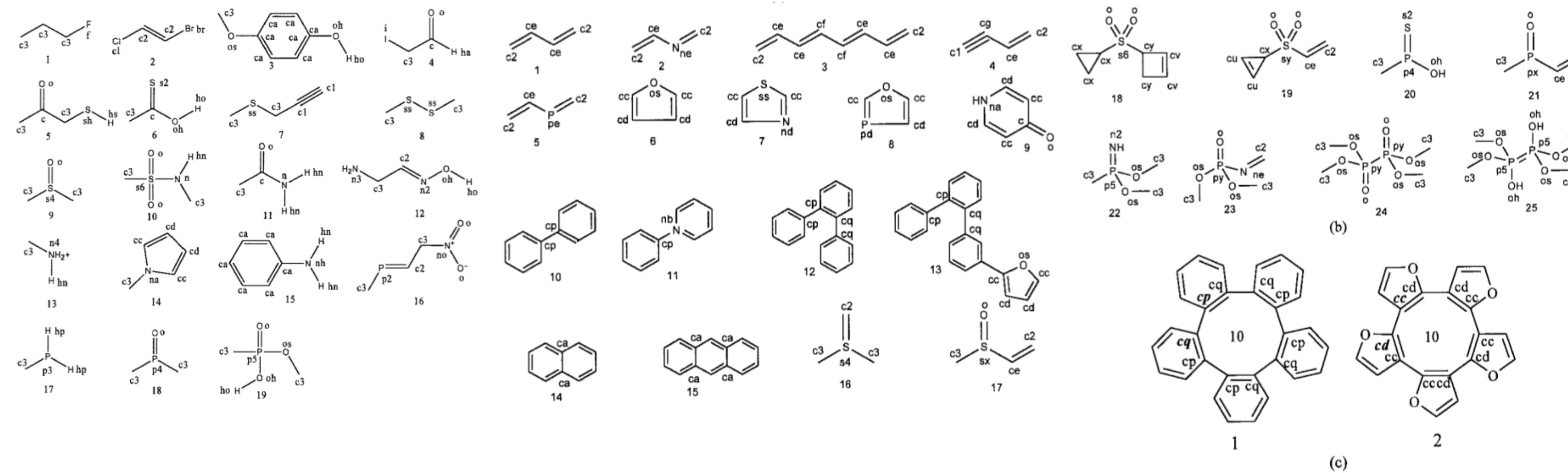
**LEE-PING WANG**

UNIVERSITY OF CALIFORNIA, DAVIS

# ACADEMIC

# As drug discovery explores new parts of chemical space, how can **forcefields** keep up?

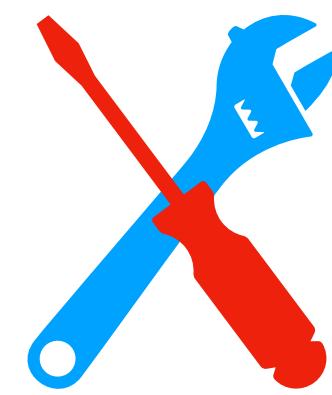
The Generalized Amber Forcefield (GAFF) was parameterized with this chemical universe:



The chemical space of potentially pharmacological active molecules is estimated to be  $\sim 0(10^{63})$   
Extension of this universe is nontrivial!

Wang J, Wolf RM, Caldwell JW, Kollman PA, and Case DA. J Comput Chem 25:1157, 2004.

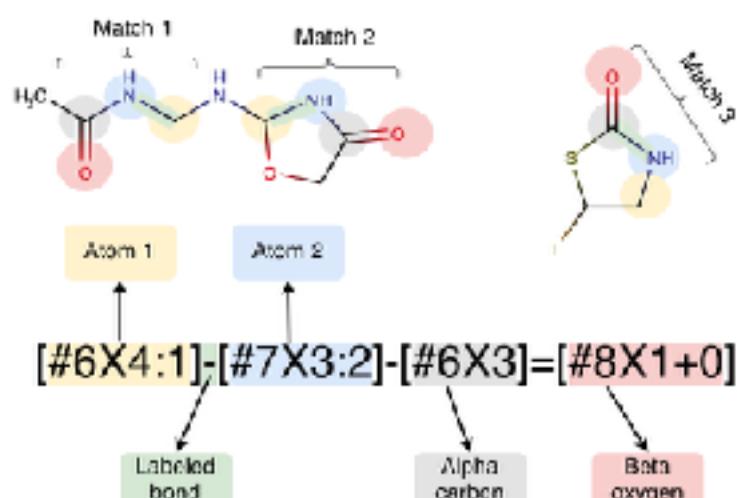
# Open Forcefield Initiative objectives



- Develop an open, scalable, and interoperable **toolkit** for automatically parameterizing forcefields



- Generate/curate **open datasets** necessary for producing high-accuracy small molecule forcefield

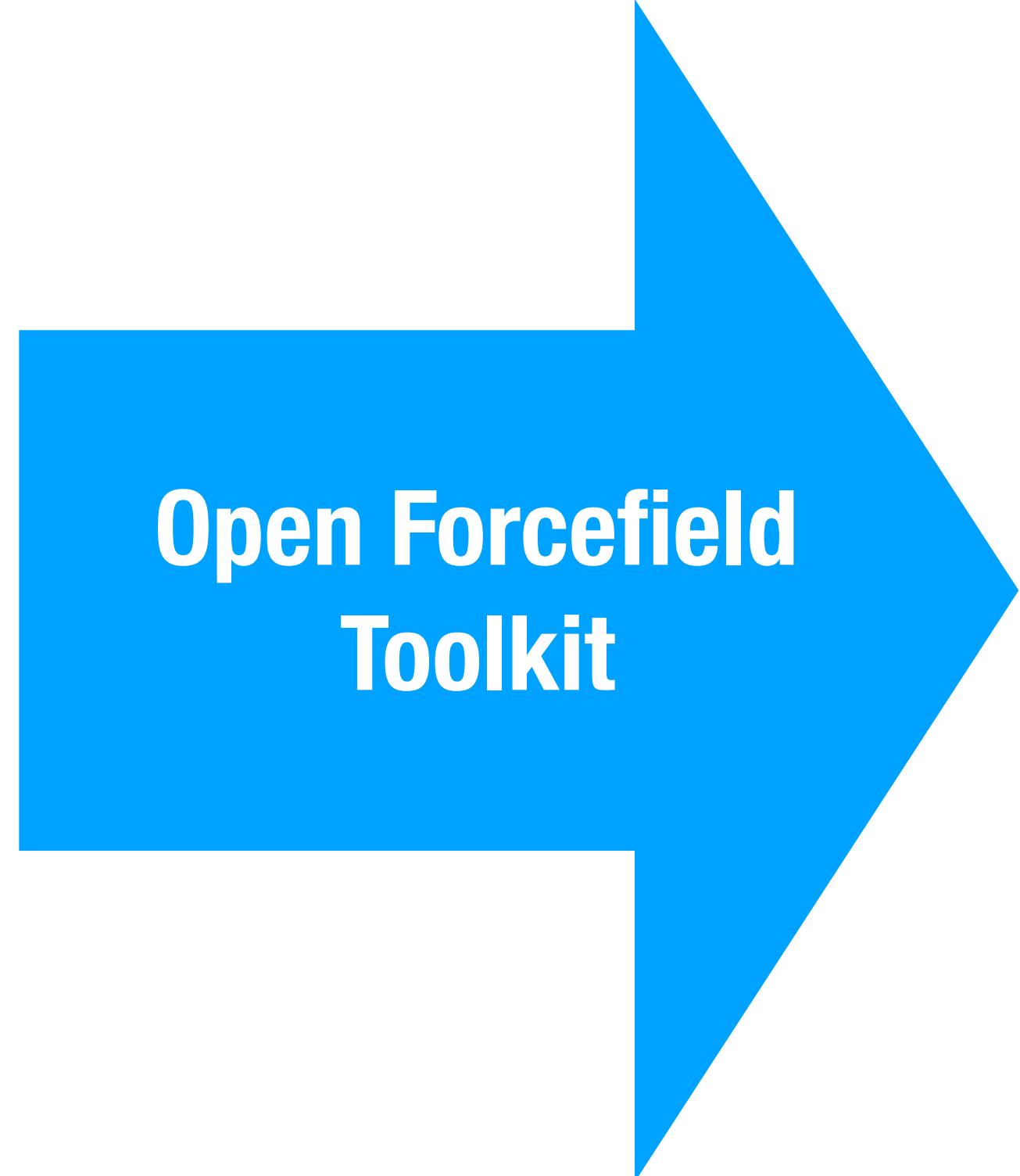


- Generate systematically-improved **forcefields**

# The Open Forcefield Toolkit provides a seamless way to parameterize systems

Molecule  
(various representations)

Forcefield  
(SMIRNOFF format)



Parameterized system  
(various representations)

Jeff Wagner



Andrea Rizzi



Binder enabled example

[https://github.com/openforcefield/openforcefield/tree/master/examples/  
using\\_smirnoff\\_with\\_amber\\_protein\\_forcefield](https://github.com/openforcefield/openforcefield/tree/master/examples/using_smirnoff_with_amber_protein_forcefield)

# The SMIRKS Native Open Force Field (SMIRNOFF) format avoids the complexity of atom typing

RDKit enables OFFToolkit to be fully Open Sourced

Caitlin Bannan



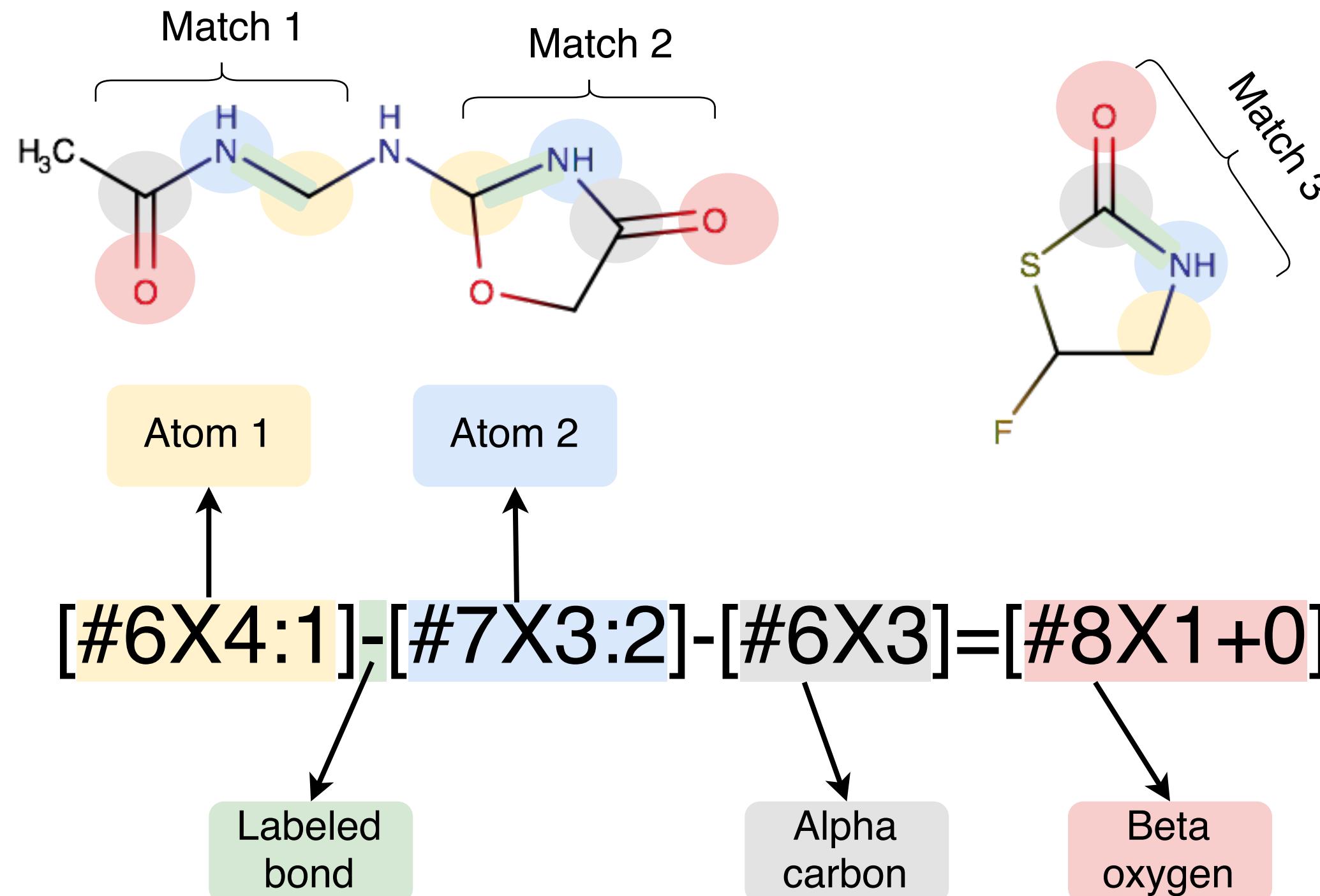
David Mobley



Christopher Bayly



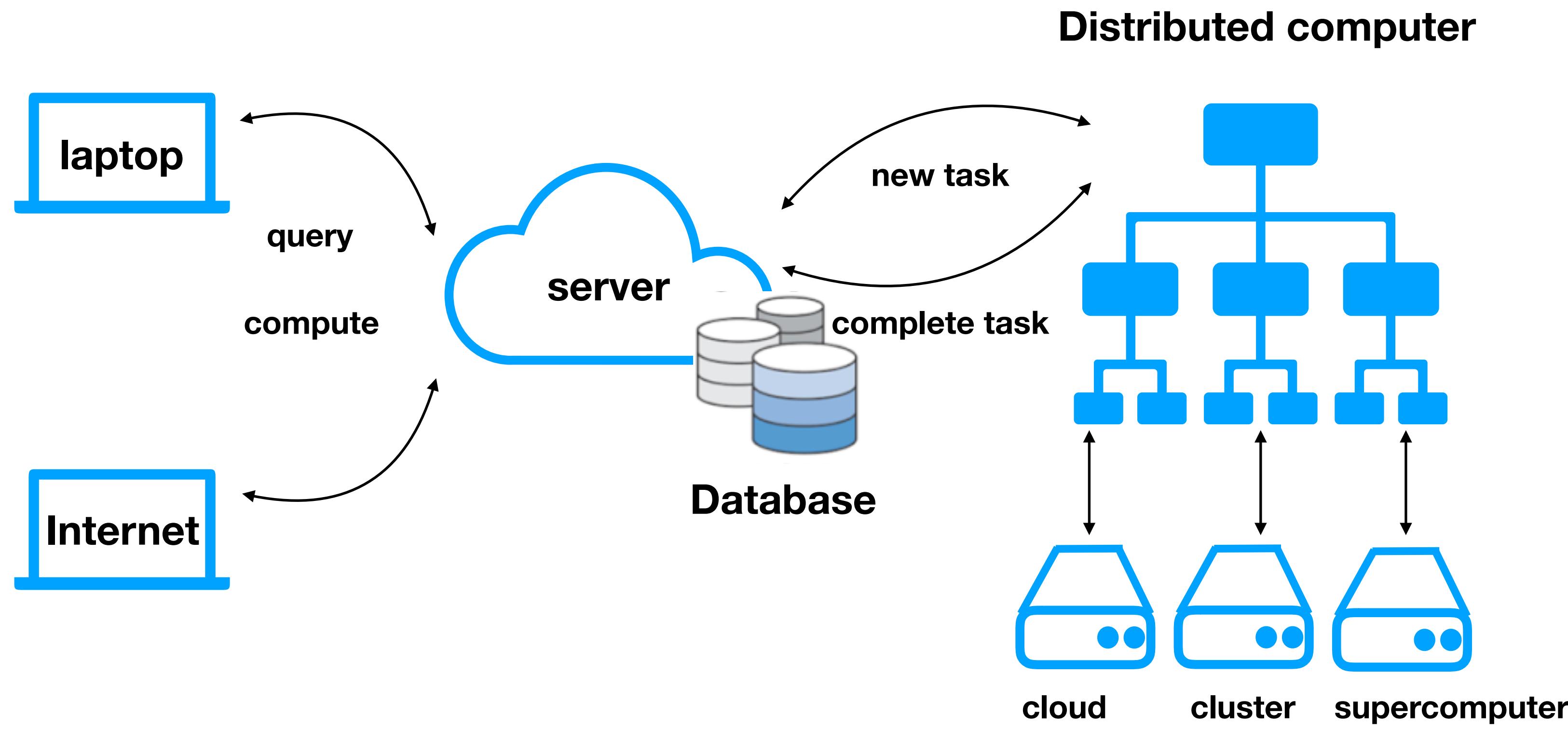
**MATCH BONDS DIRECTLY:**



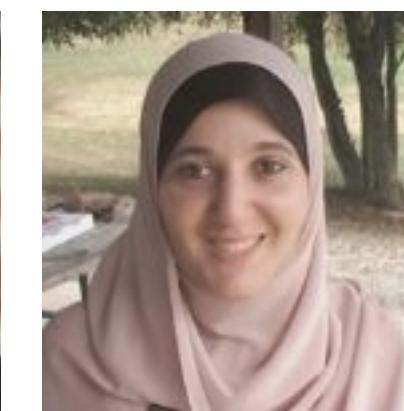
Mobley, D., et. al (2018). JCTC, 14(11), 6076–6092. <http://doi.org/10.1021/acs.jctc.8b00640>

Use of industry-standard SMARTS/SMIRKS chemical perception greatly simplifies tooling for parameter assignment while solving issues with extensibility and flexibility

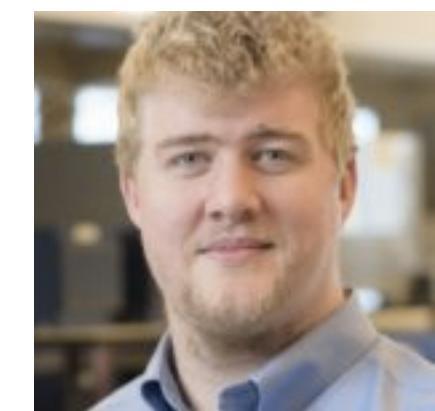
# QCArchive enables automatic generation and sharing of high-quality quantum chemistry data



Daniel Smith

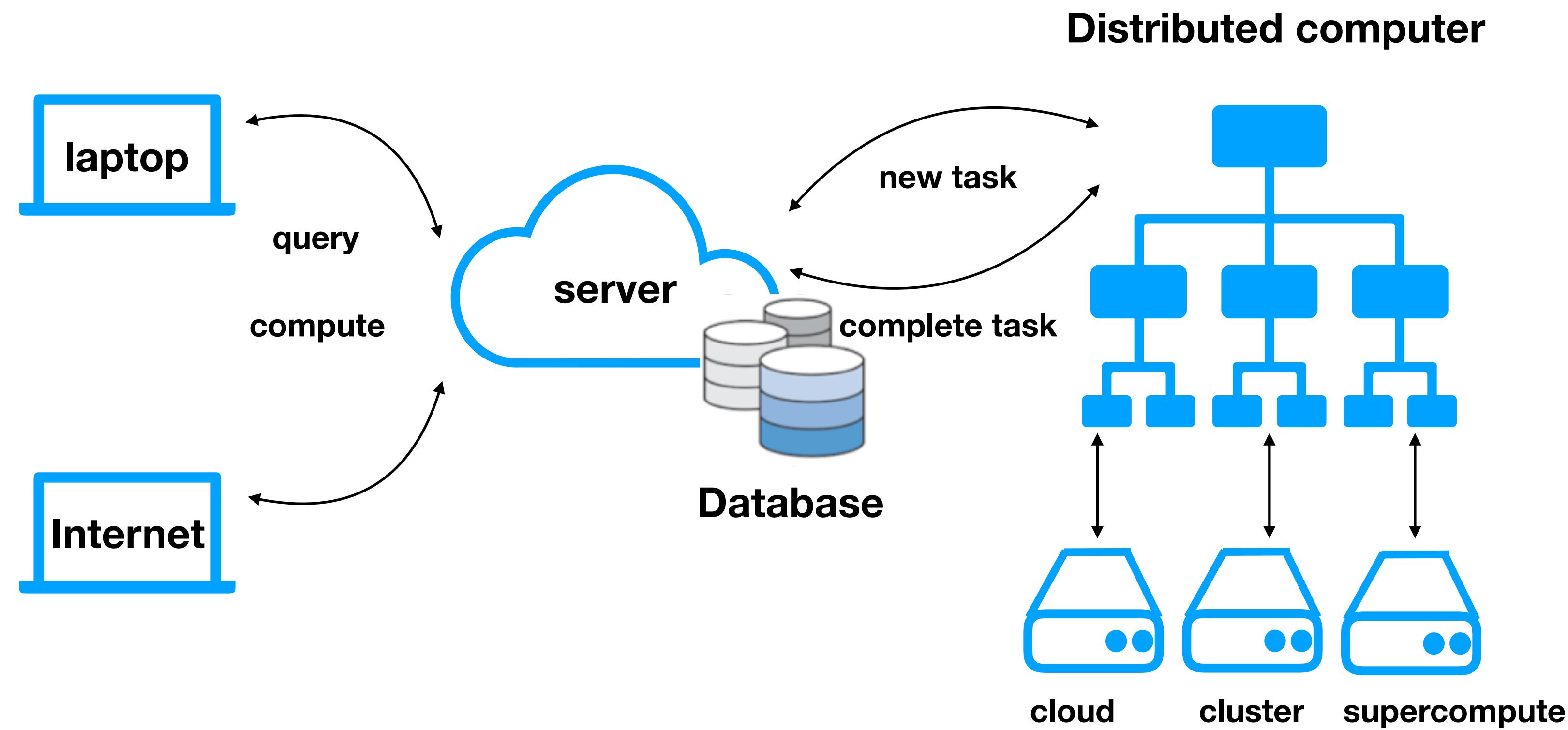


Doaa Altarawy



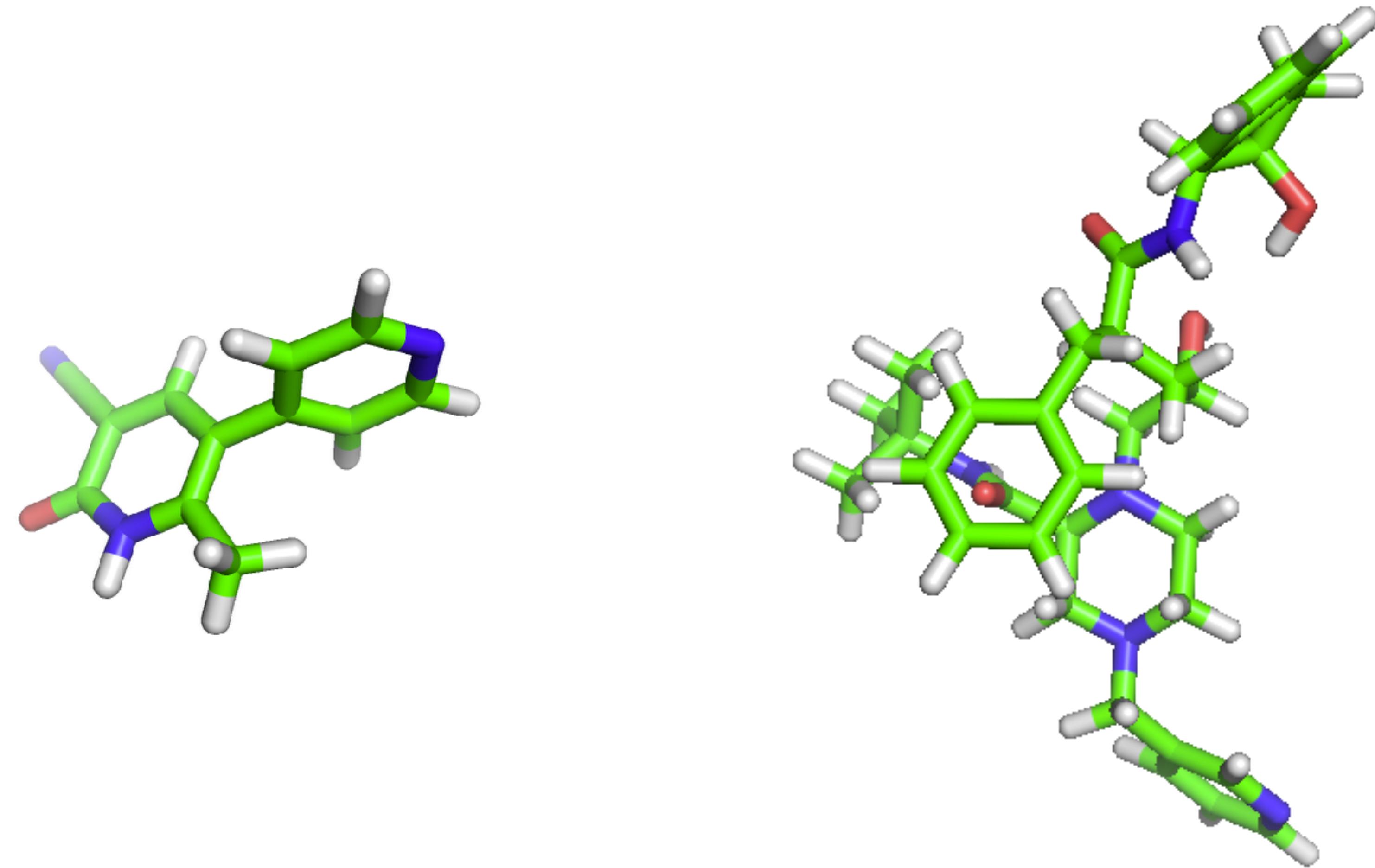
Levi Naden

# QCArchive enables automatic generation and sharing of high-quality quantum chemistry data

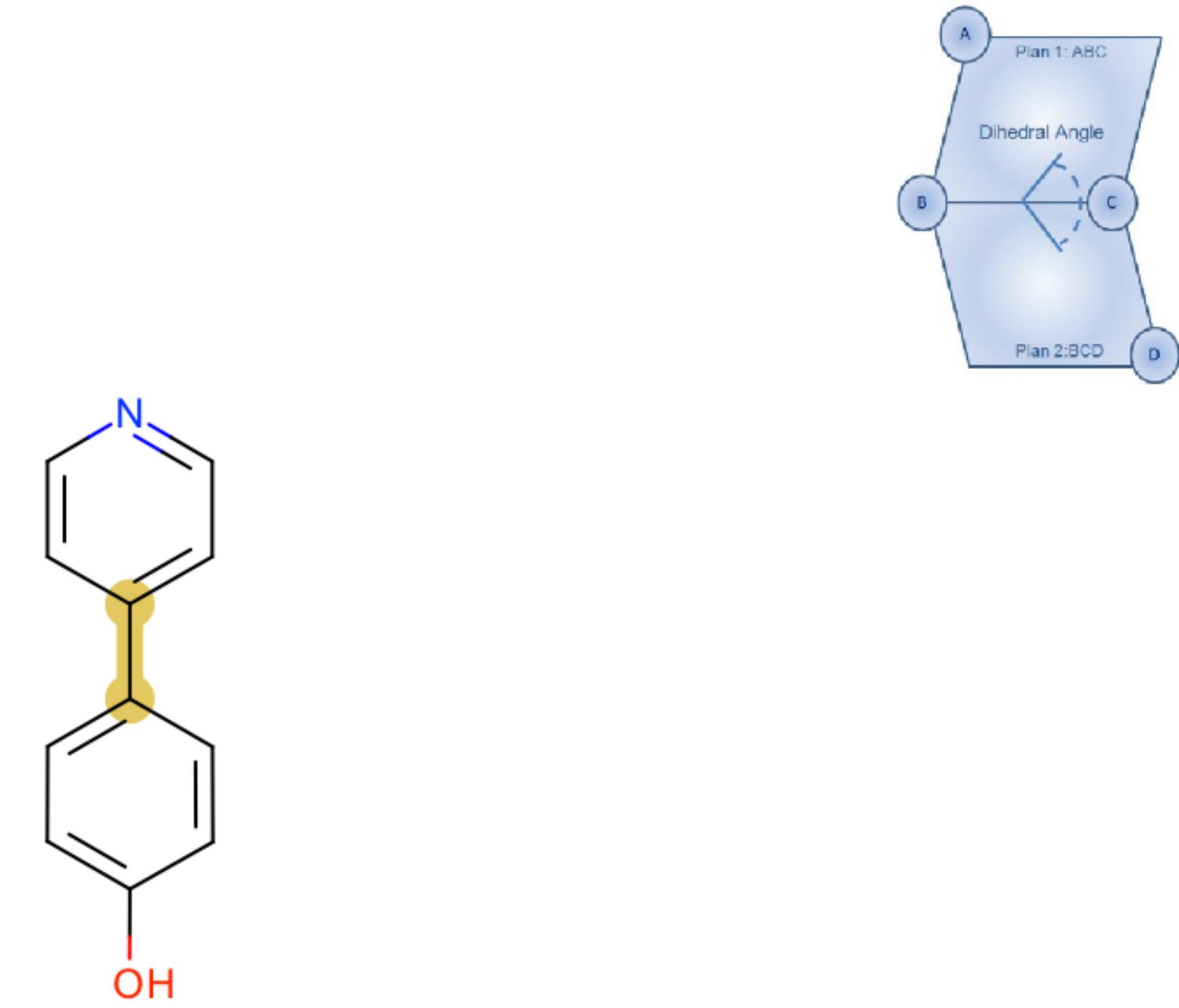
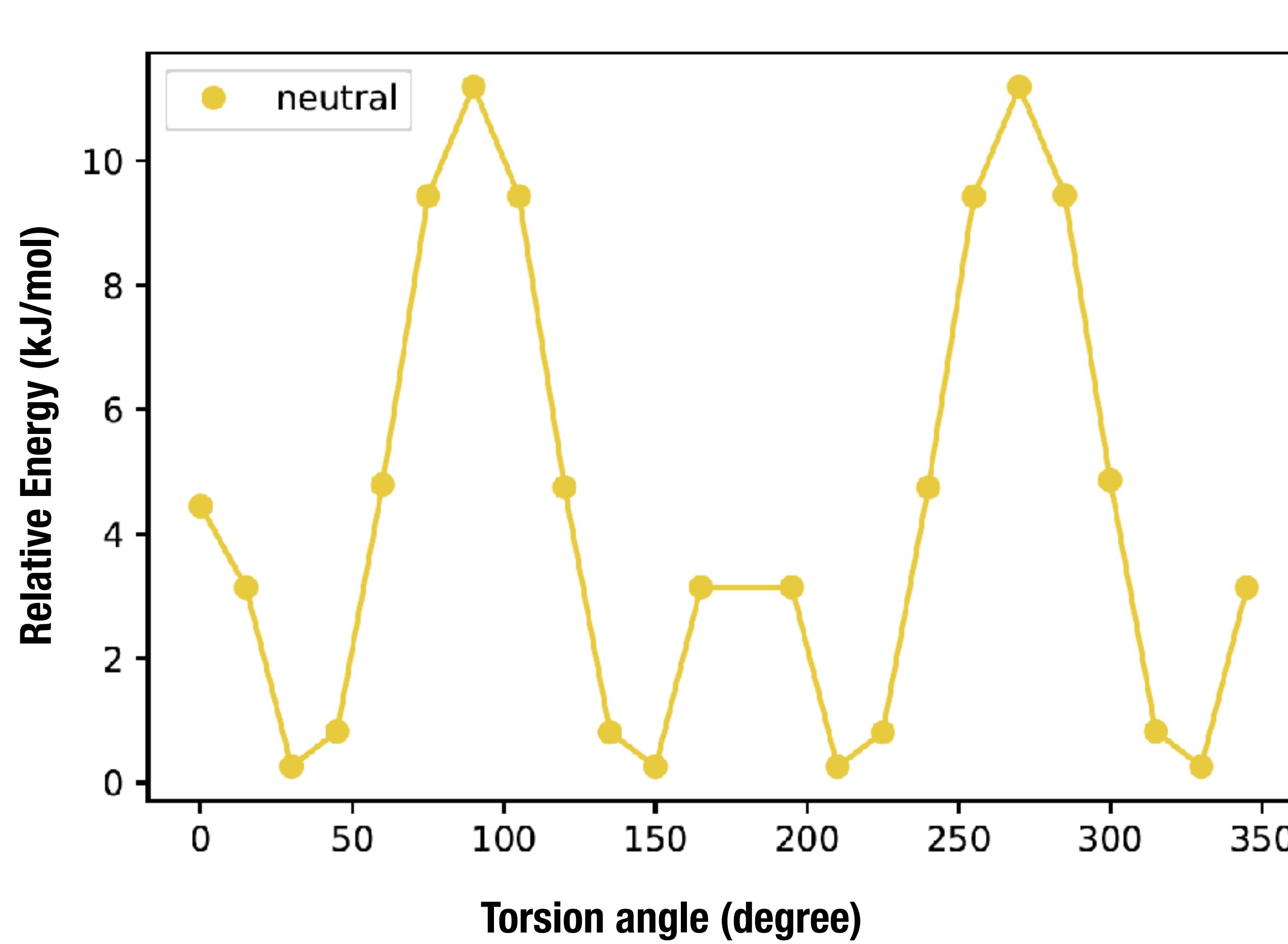


Daniel Smith   Doaa Altarawy   Levi Naden

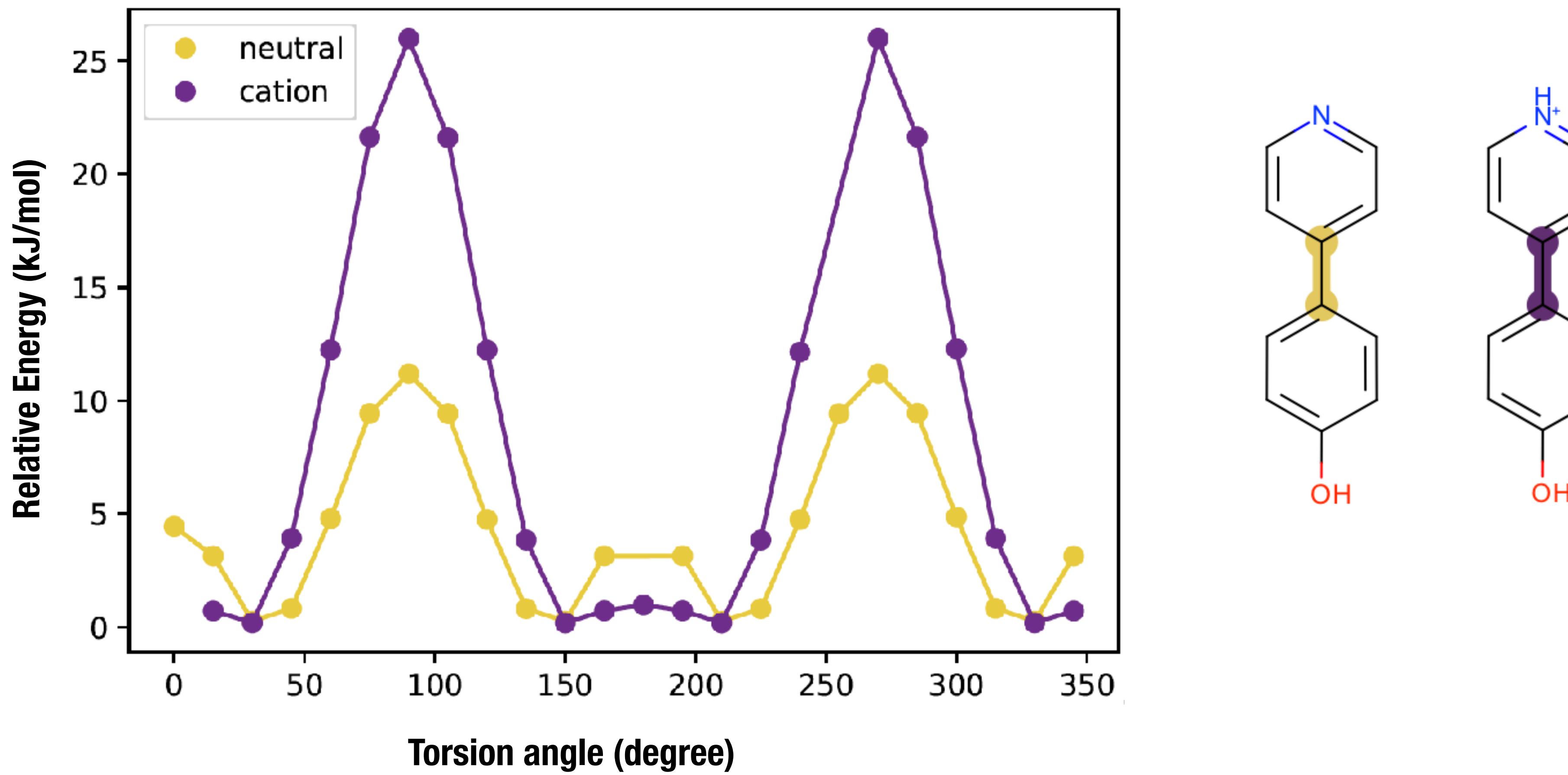
**Fragmenting molecules is necessary to avoid high computational cost of generating QC data**



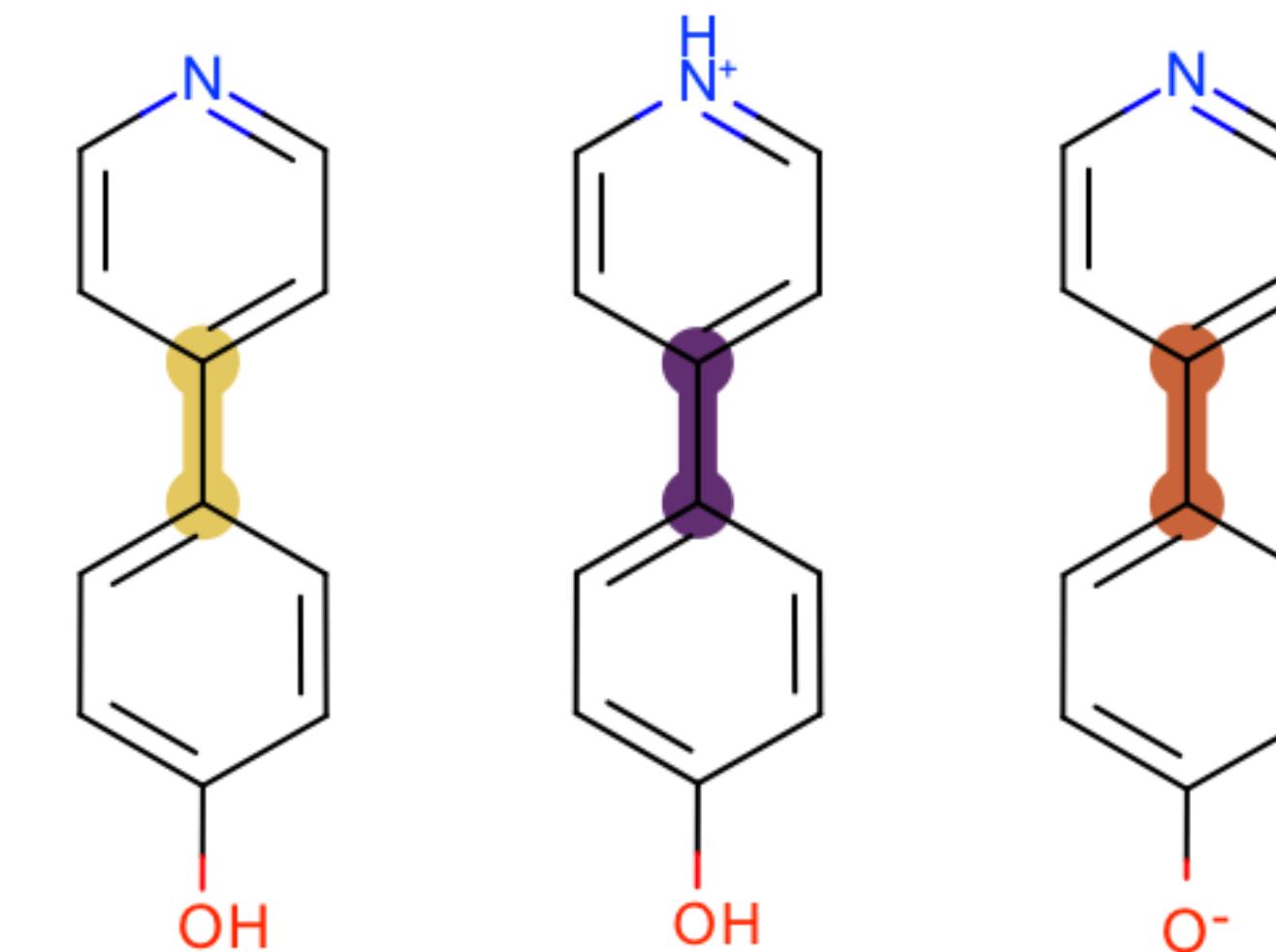
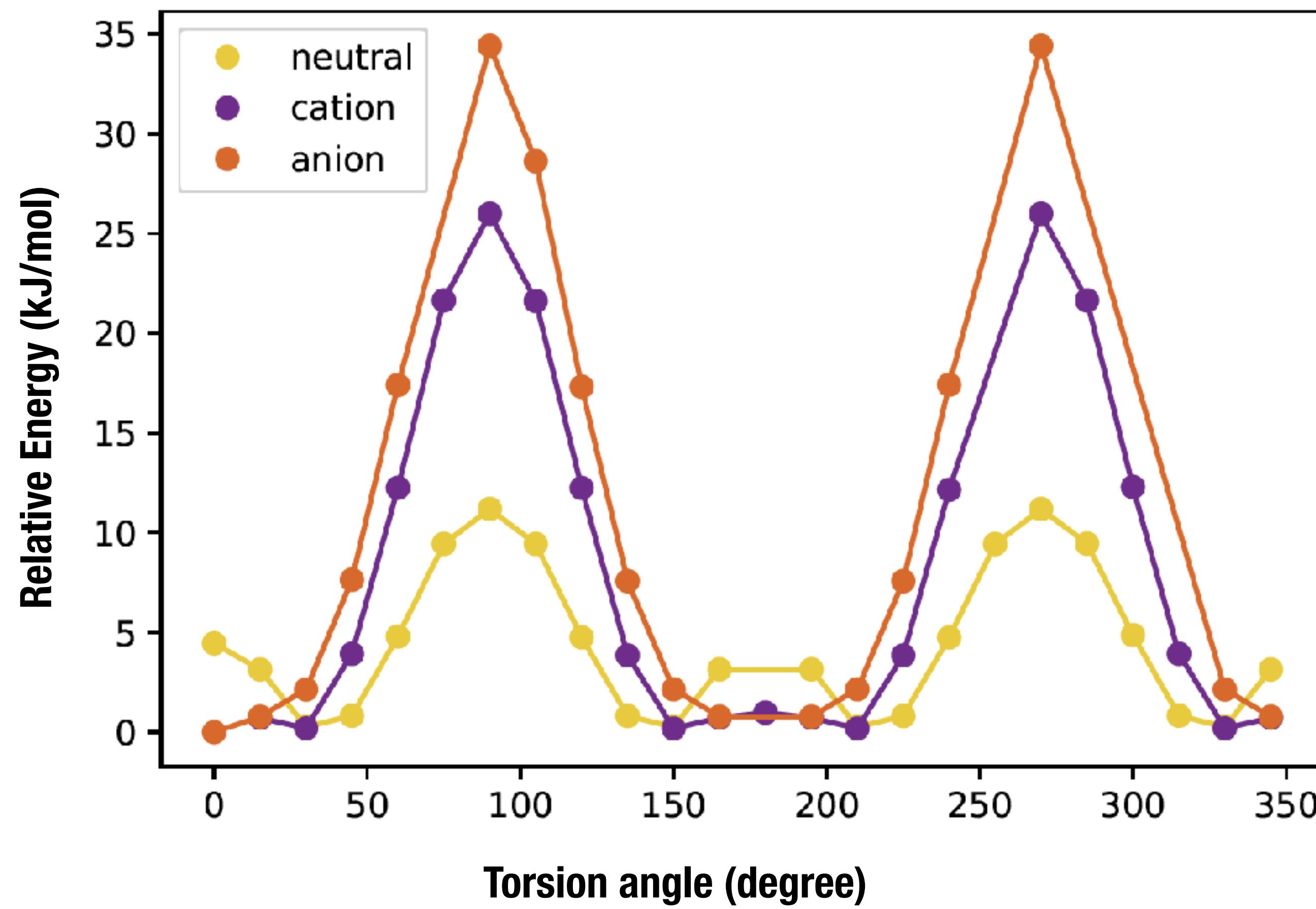
# The chemical environment of bonds is sensitive to small remote substituent changes



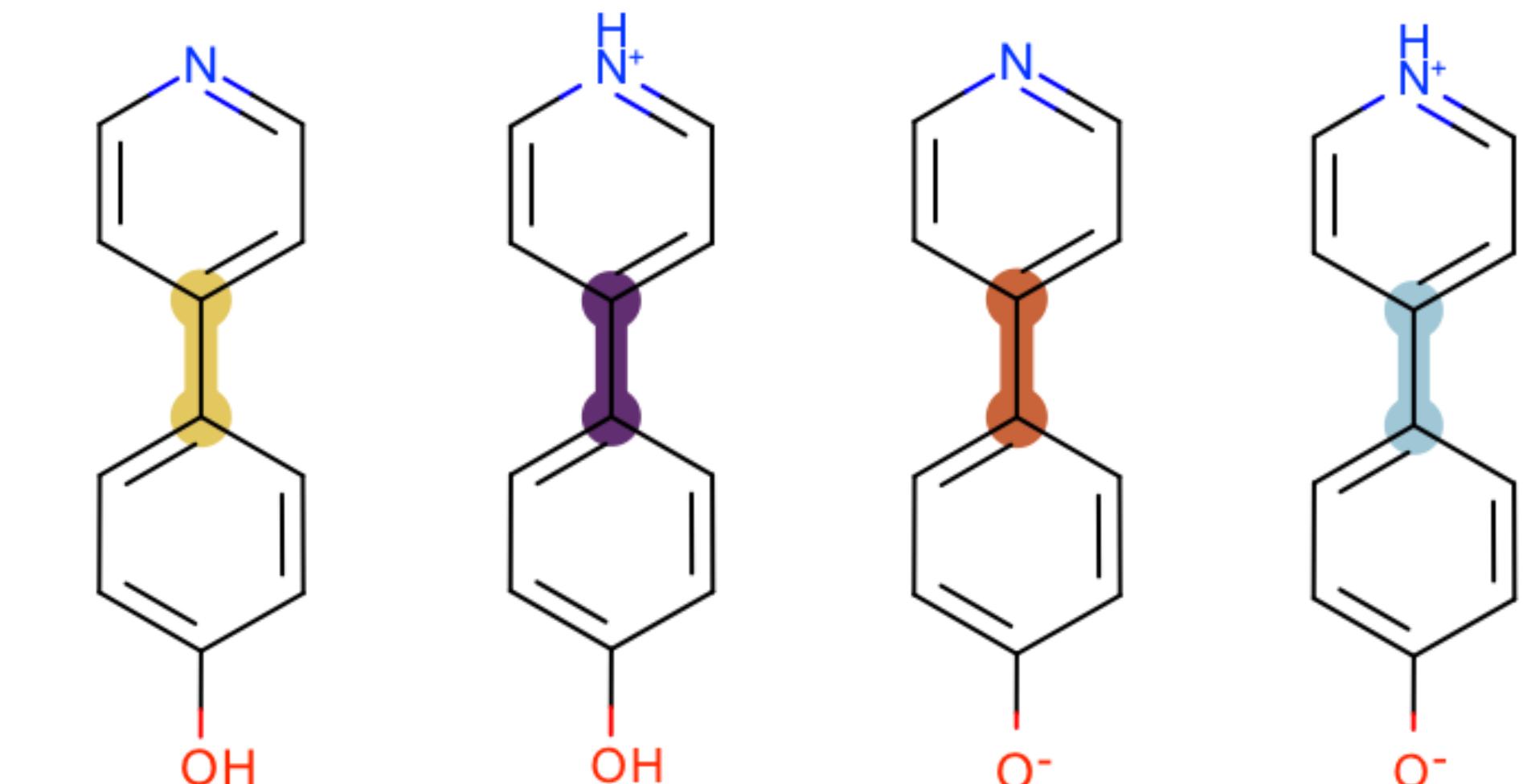
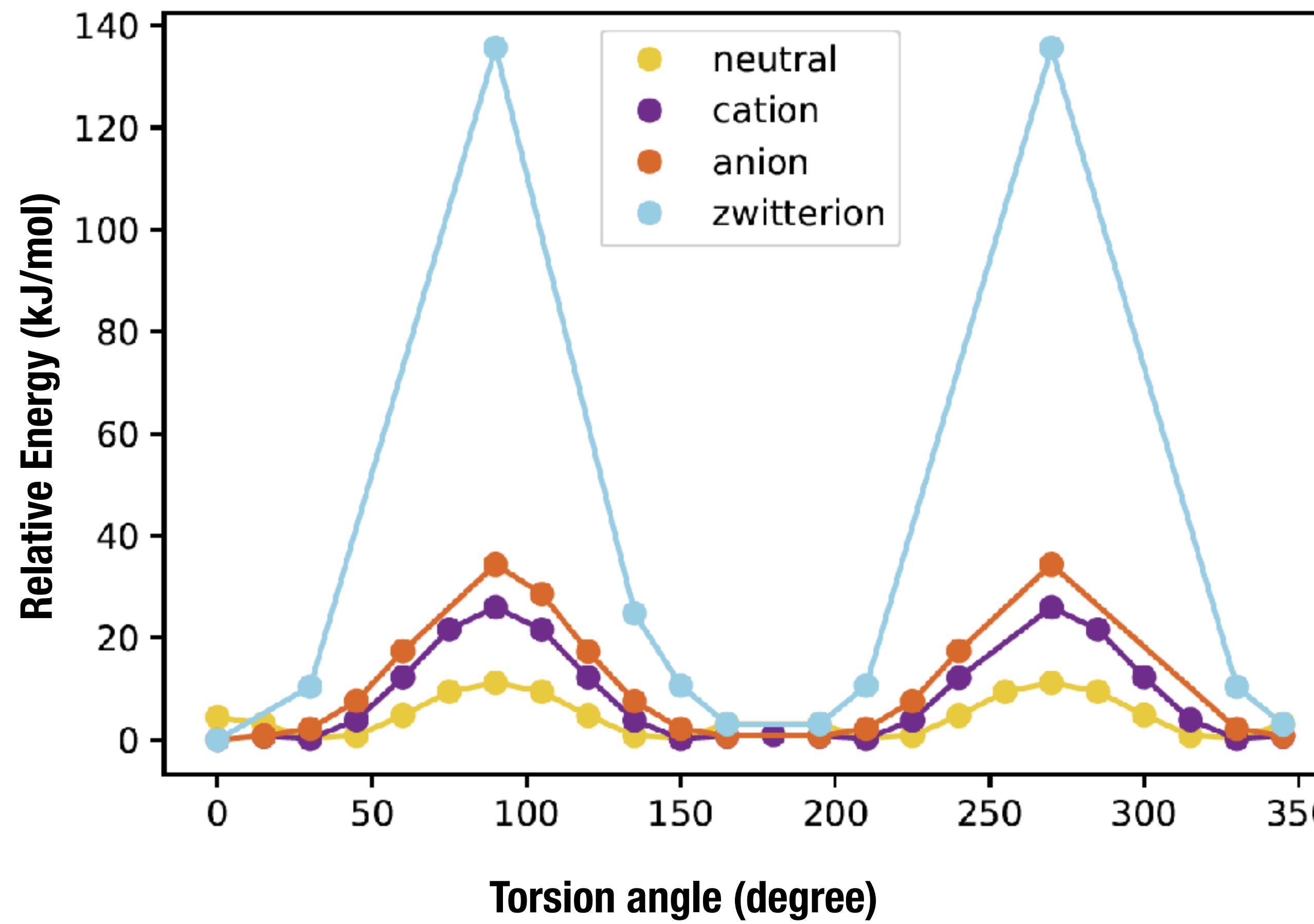
# The chemical environment of bonds is sensitive to small remote substituent changes



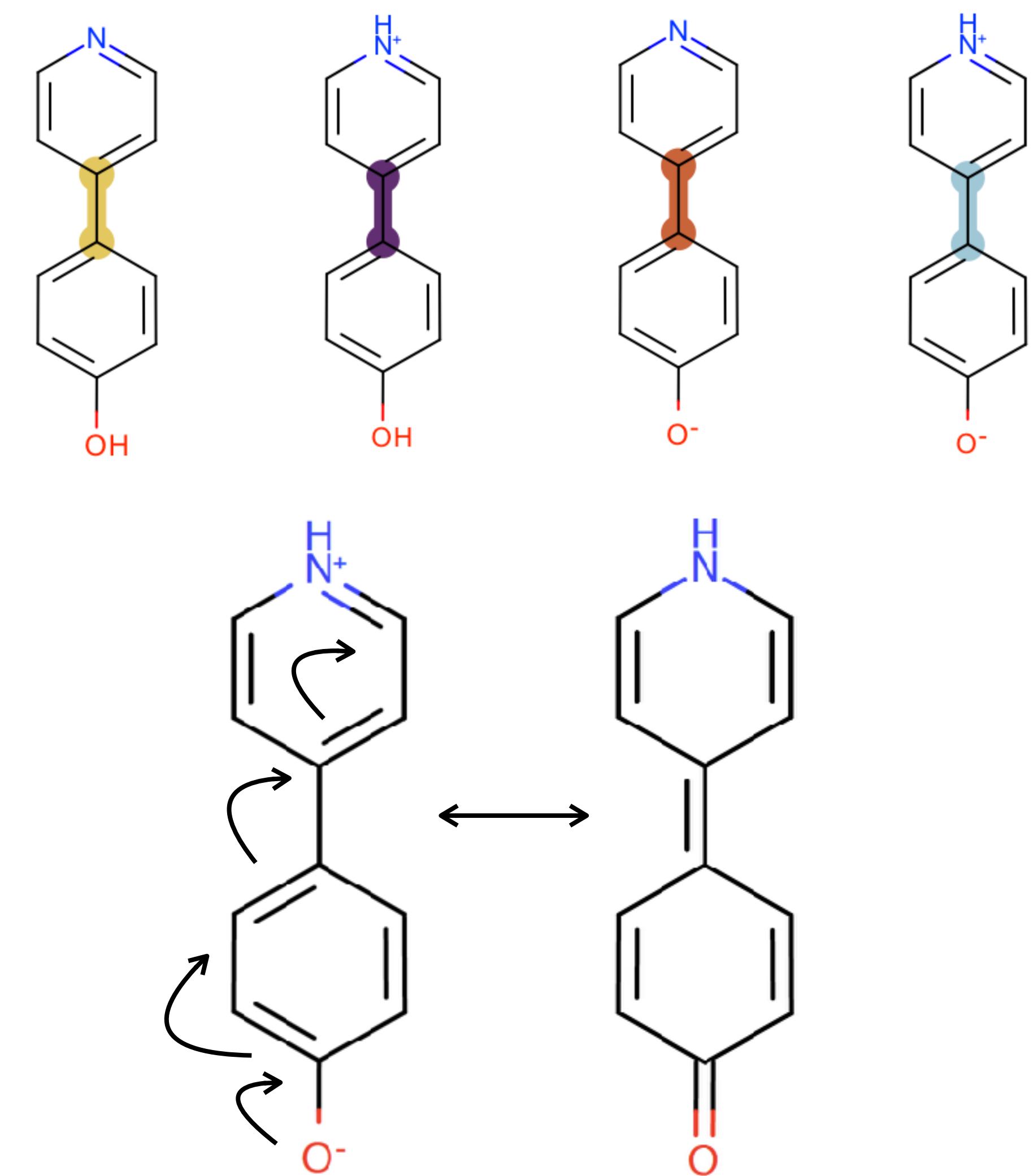
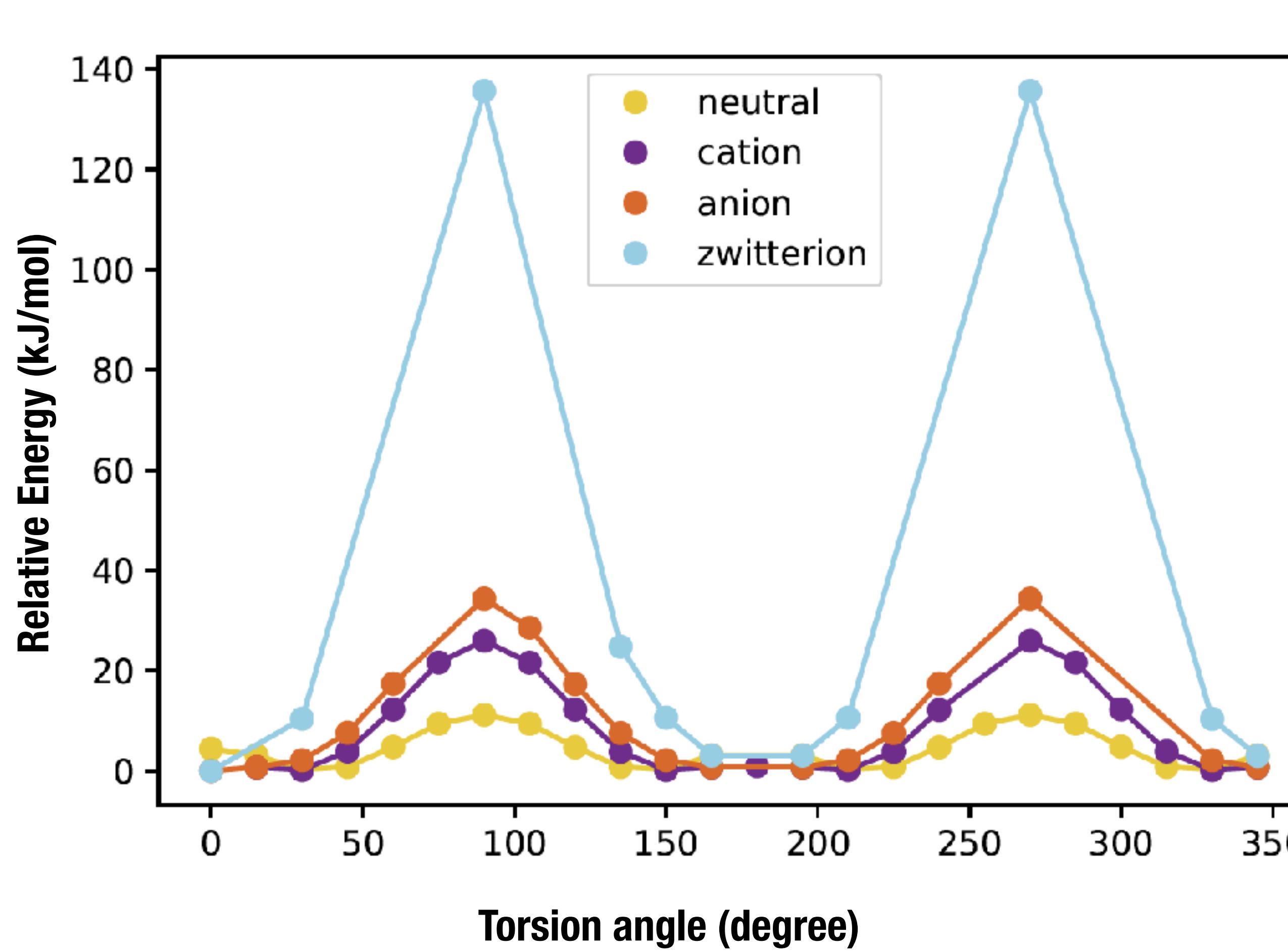
# The chemical environment of bonds is sensitive to small remote substituent changes



# Considering **chemical environment** is crucial when fragmenting bonds



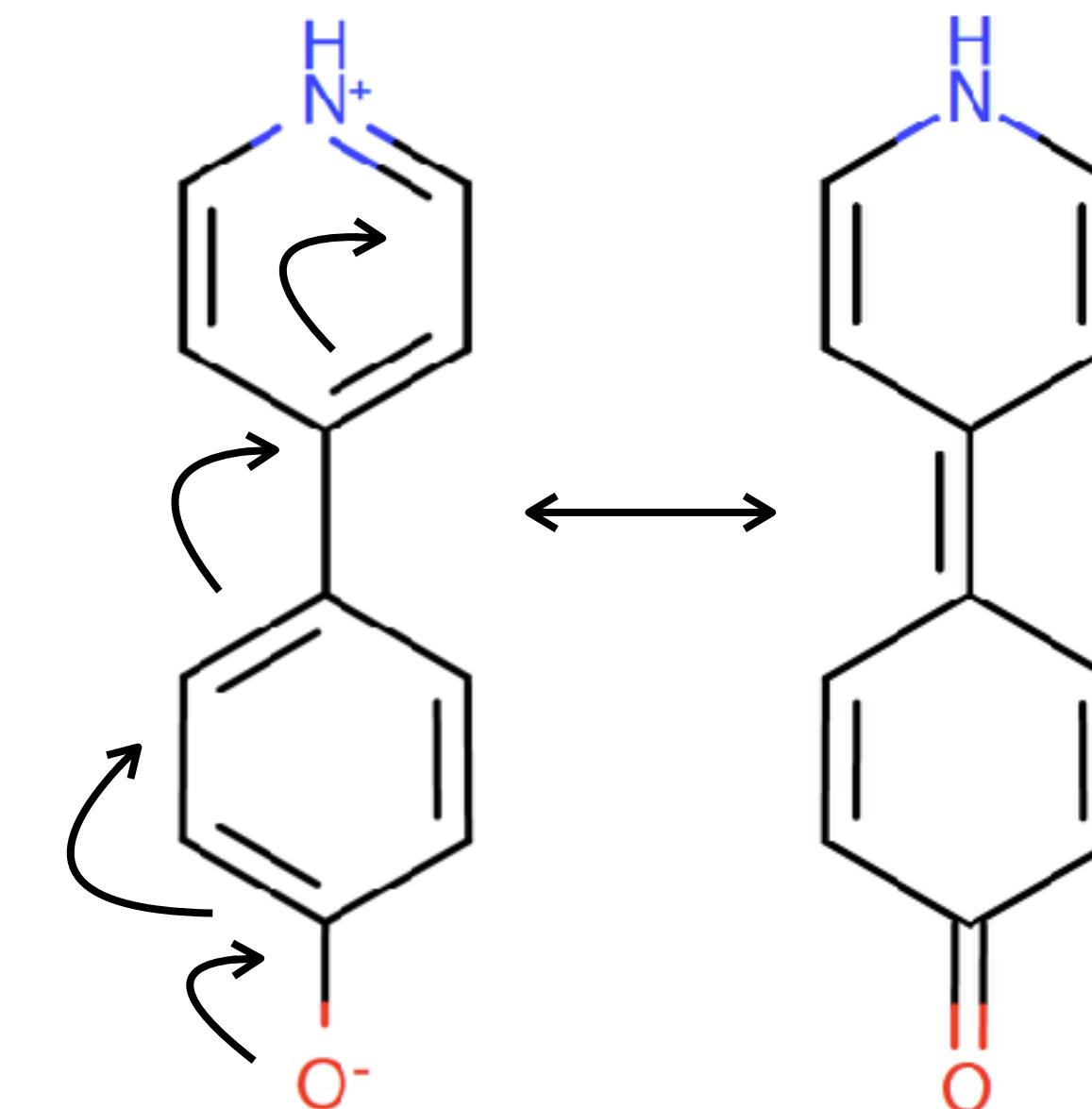
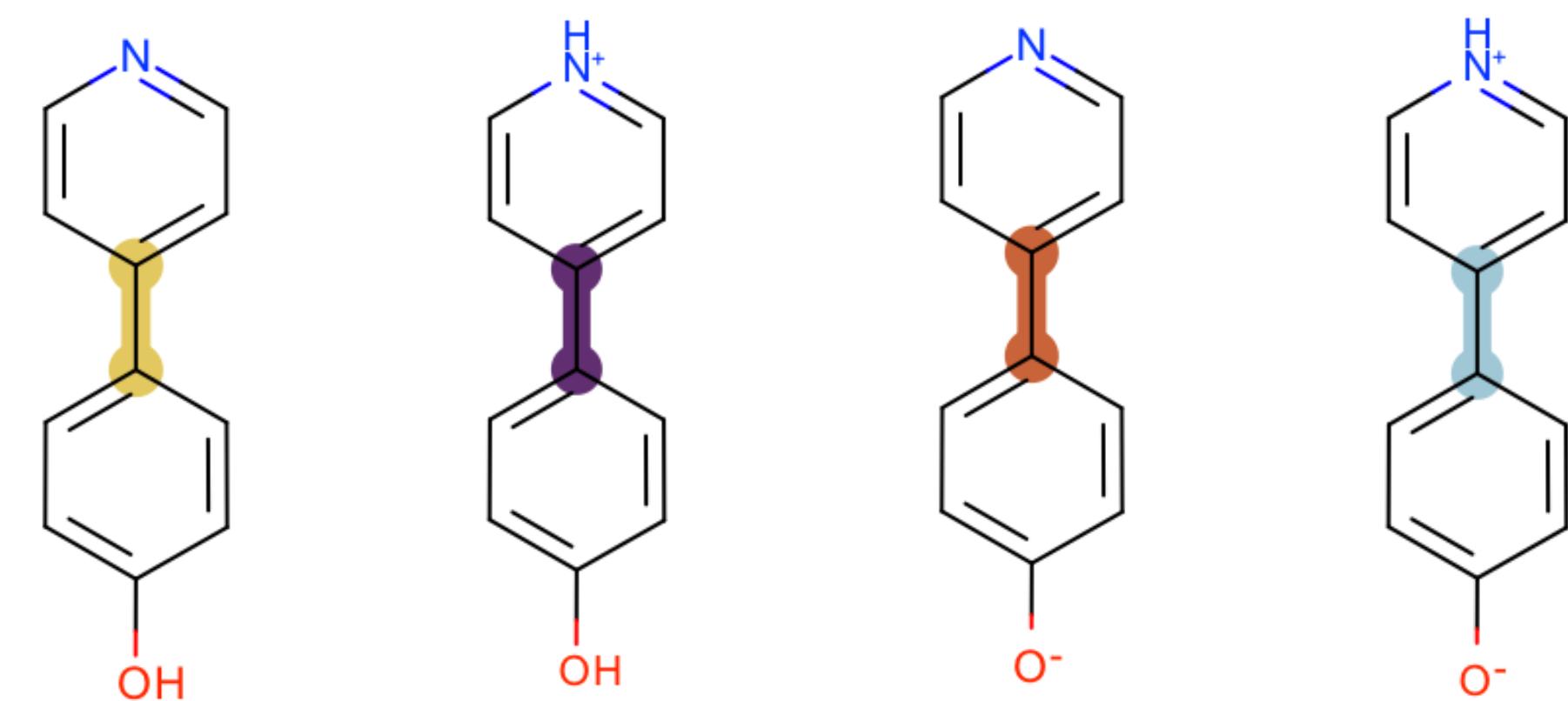
# Chemical environment impacts the torsion profile of small molecules



# Chemical environment impacts the torsion profile of small molecules

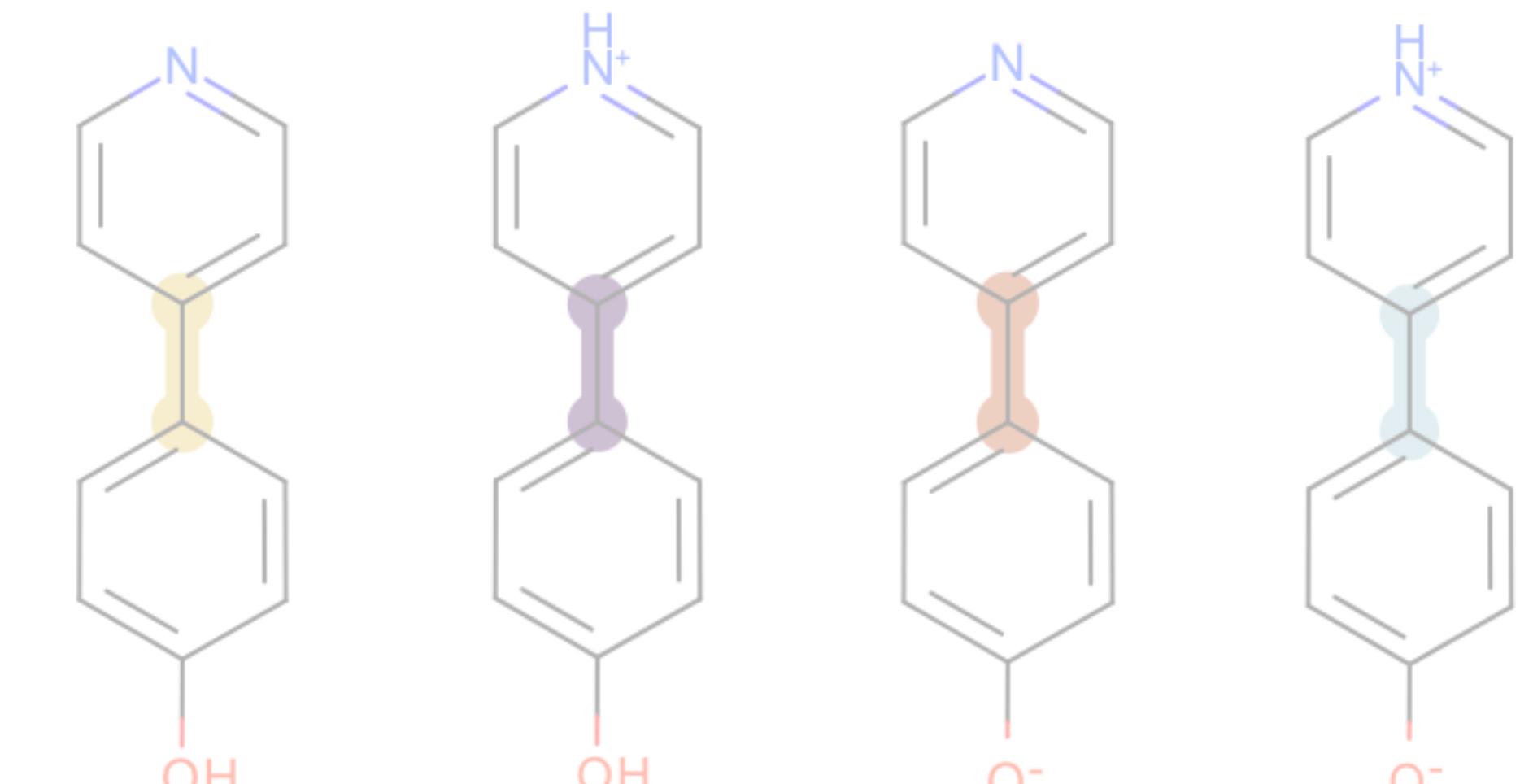
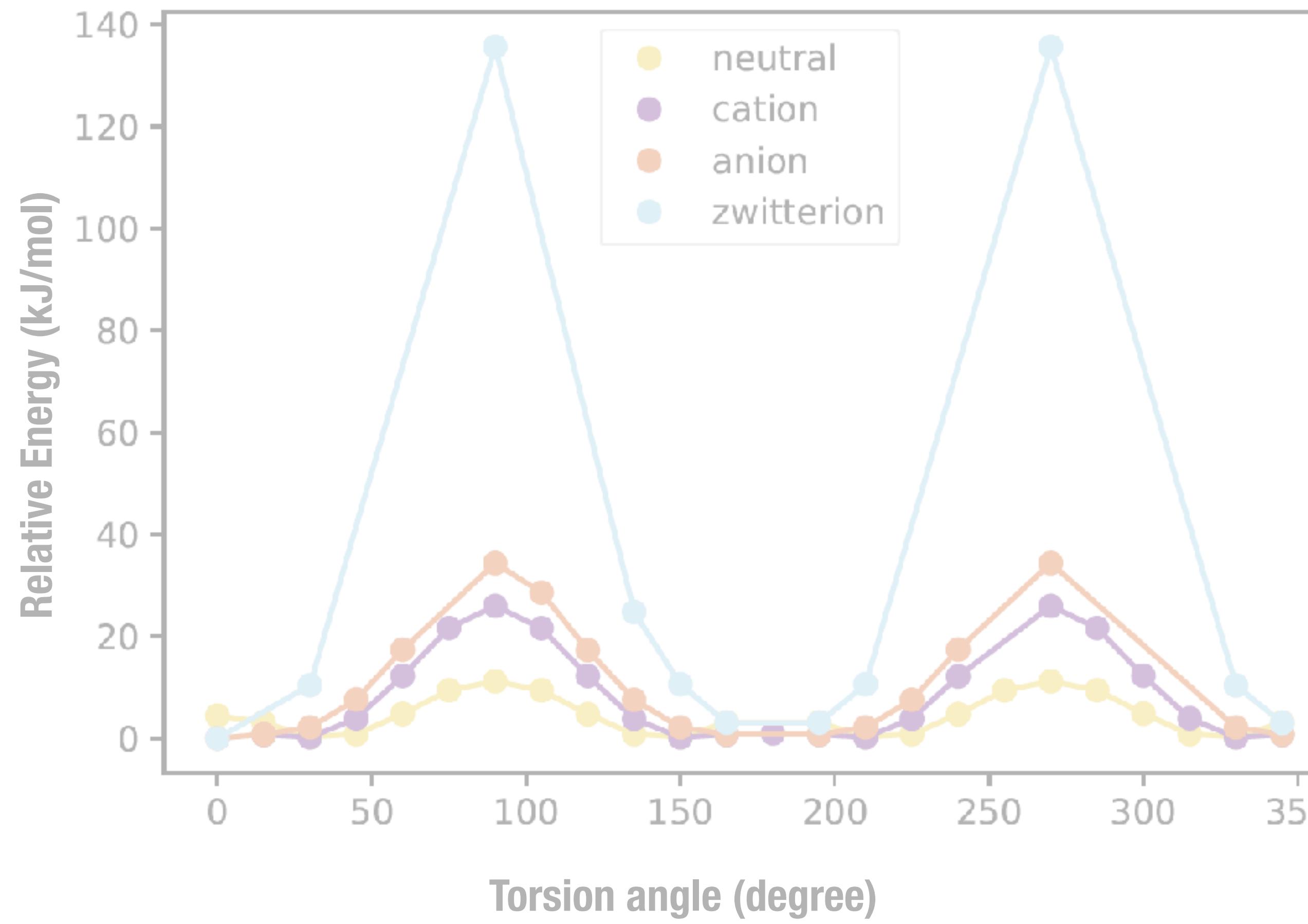
This issue poses several **problems** when  
**automating** molecule fragmentation:

1. Most cheminformatics tools label the central bond as **rotatable**
2. How do we ensure that we do not naively destroy the **chemical environment** by fragmenting a remote substituent?



# The AM1 Wiberg Bond Order is a cheap measure of electron population overlap between two atoms

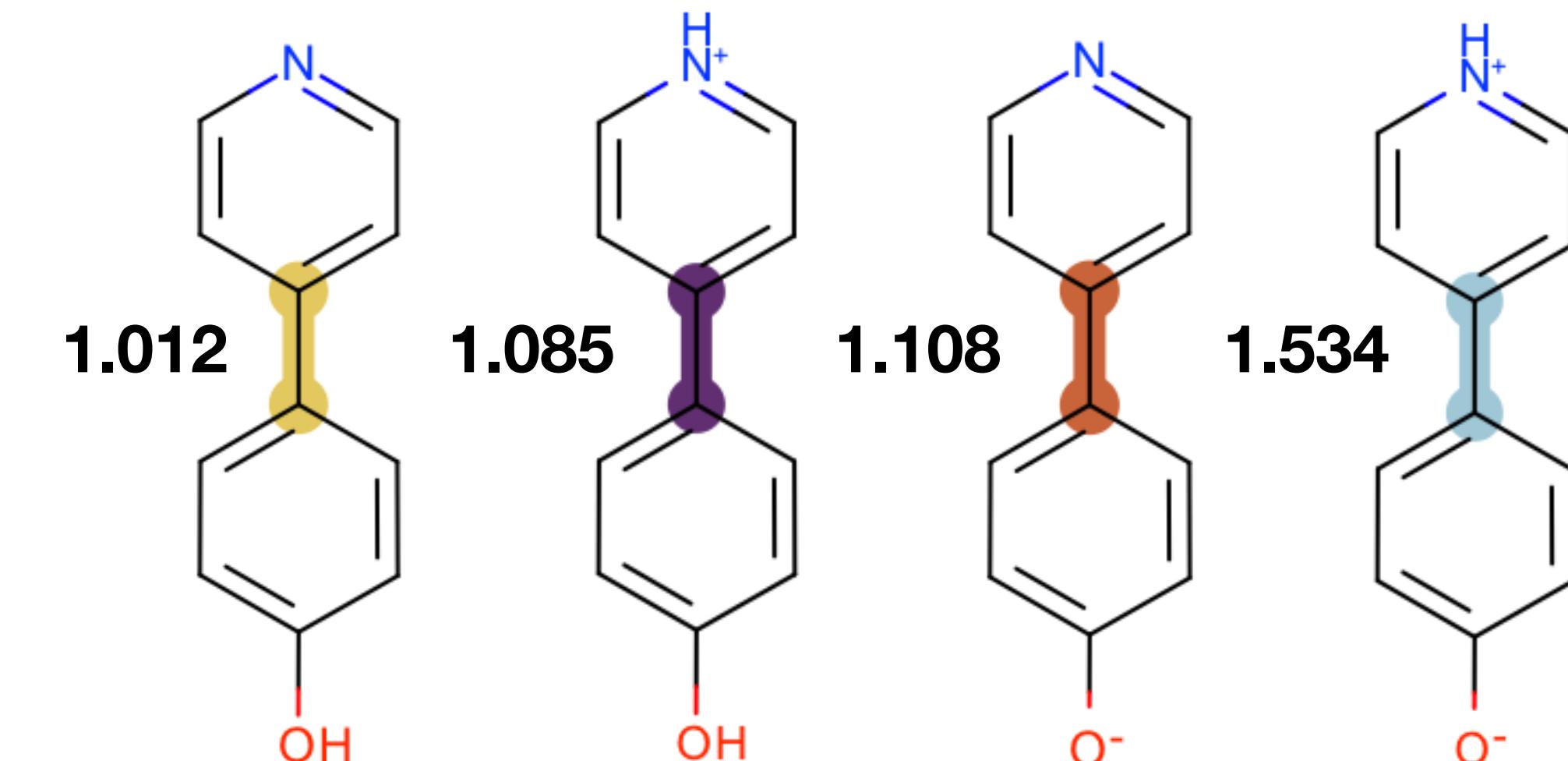
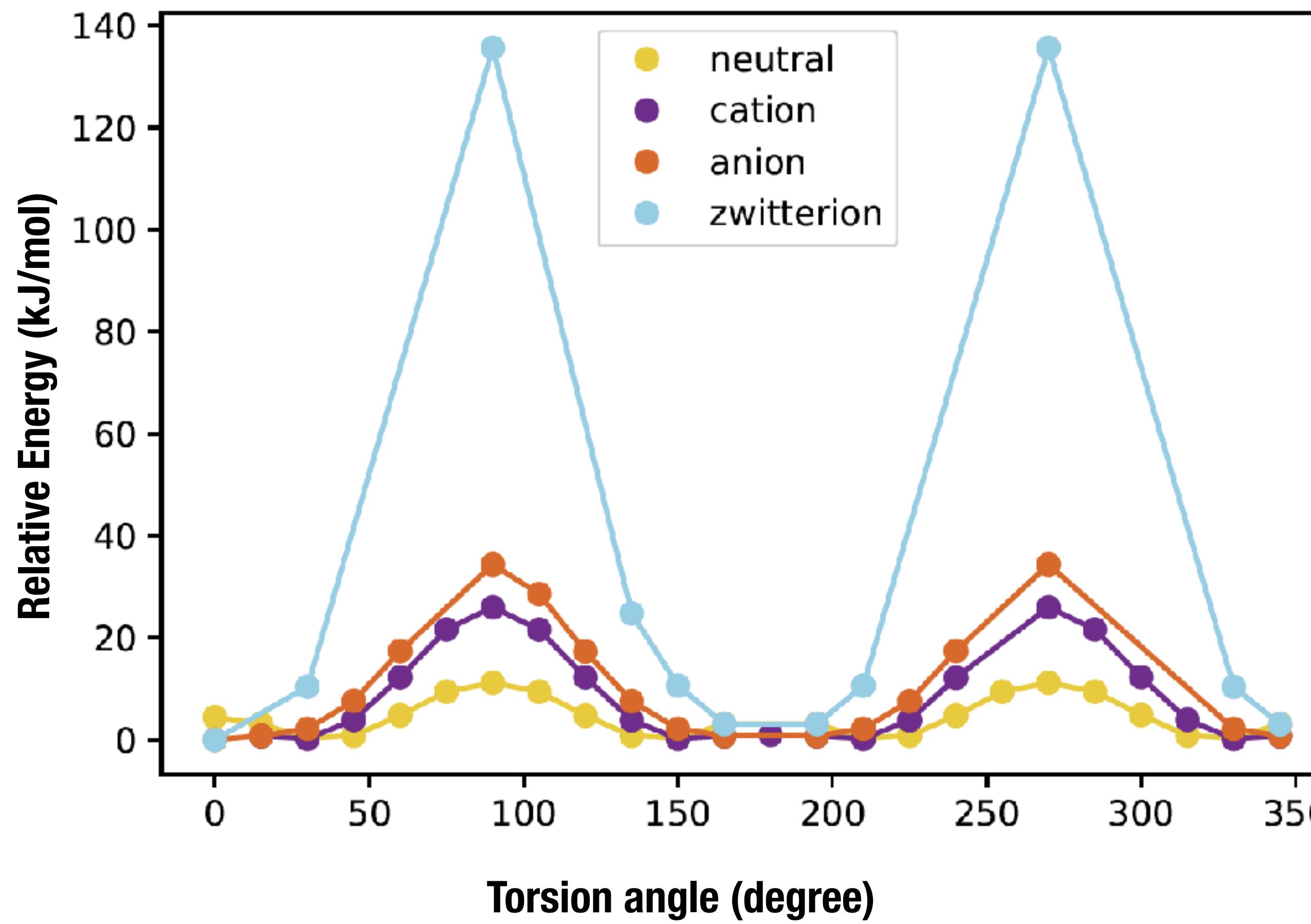
$$W_{AB} = \sum_{\mu \in A} \sum_{\nu \in B} |D_{\mu\nu}|^2.$$



$$D_{\mu\nu} = 2 \sum_i^{\text{occ.}} C_{\mu i} C_{\nu i}^*,$$

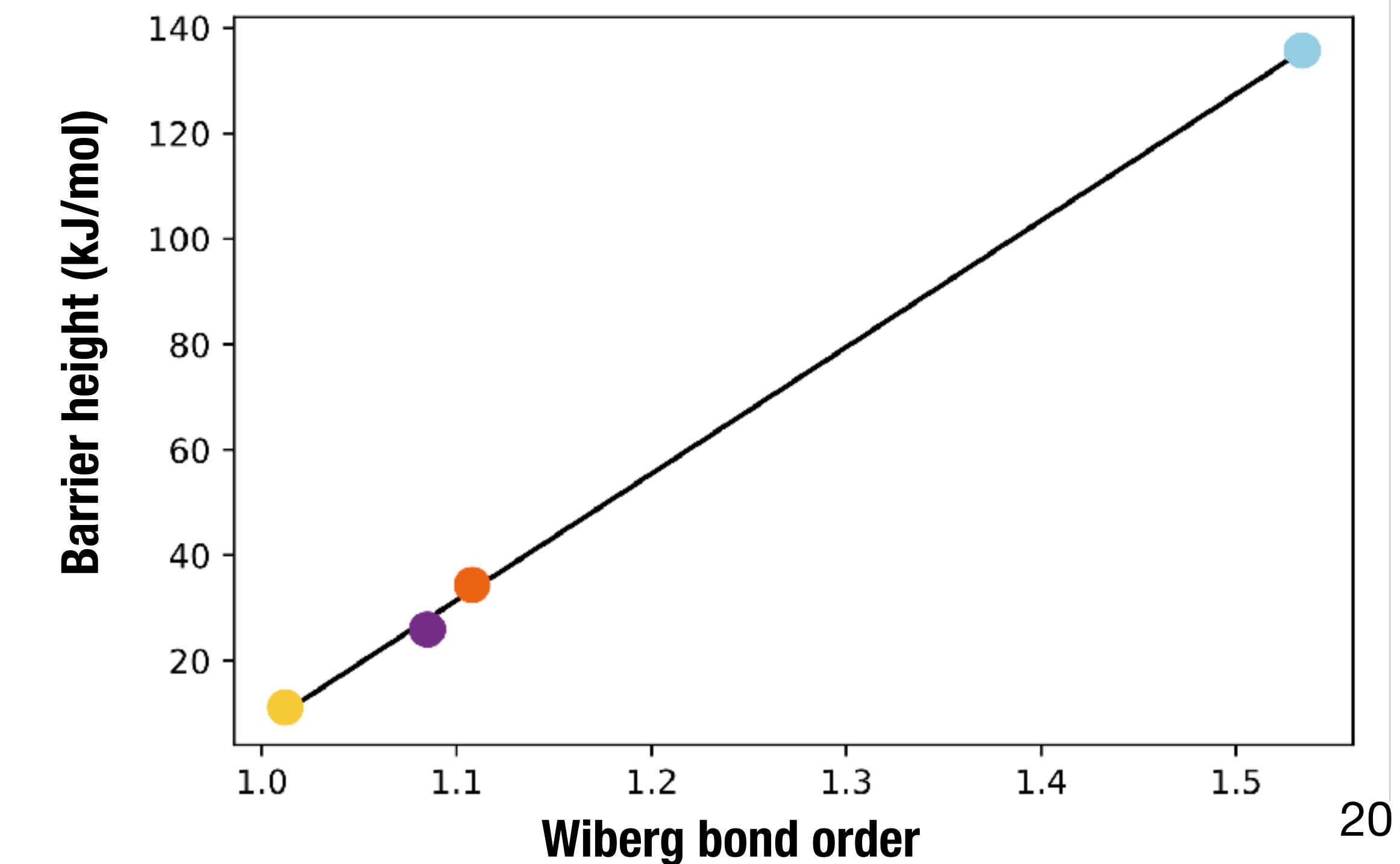
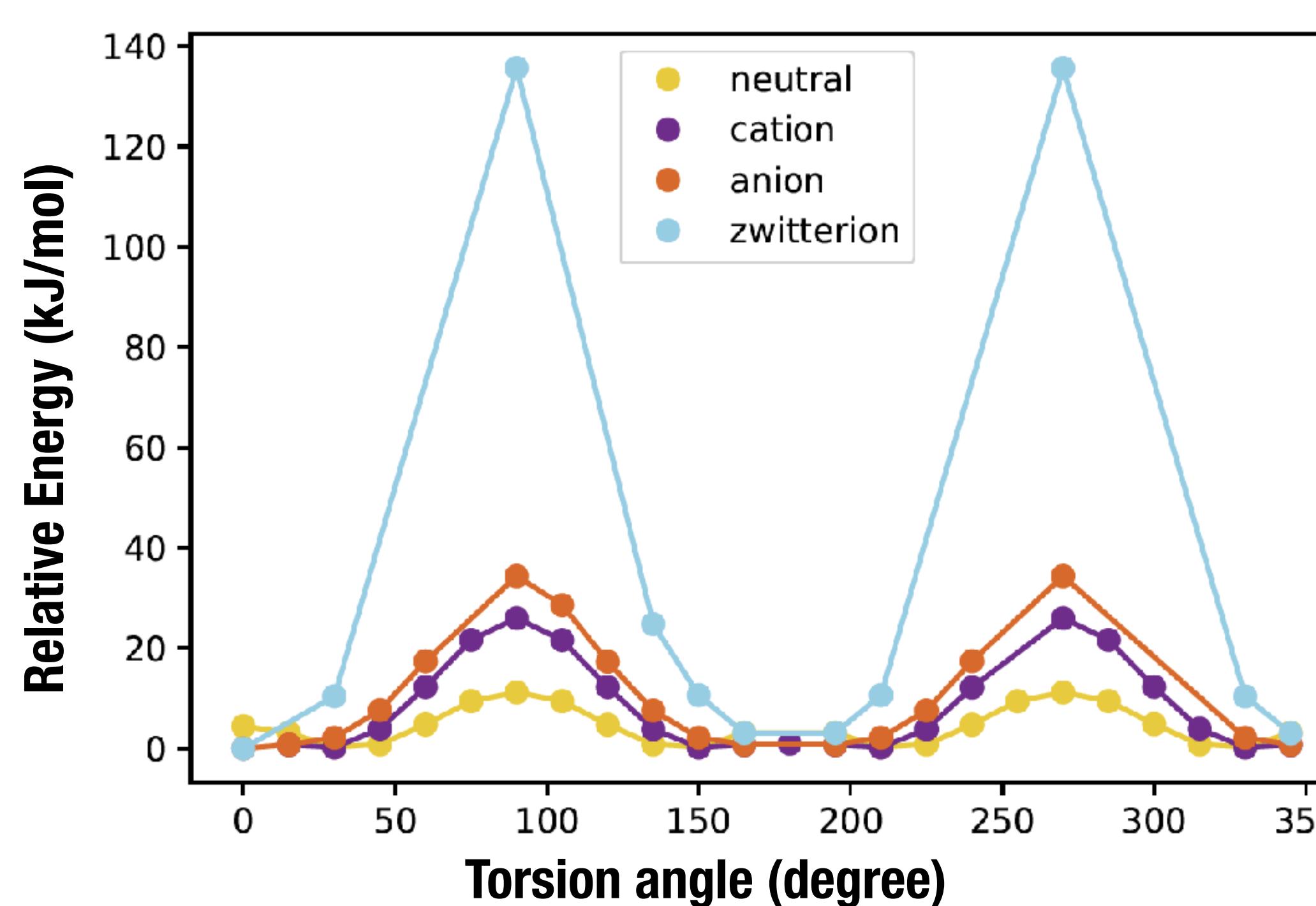
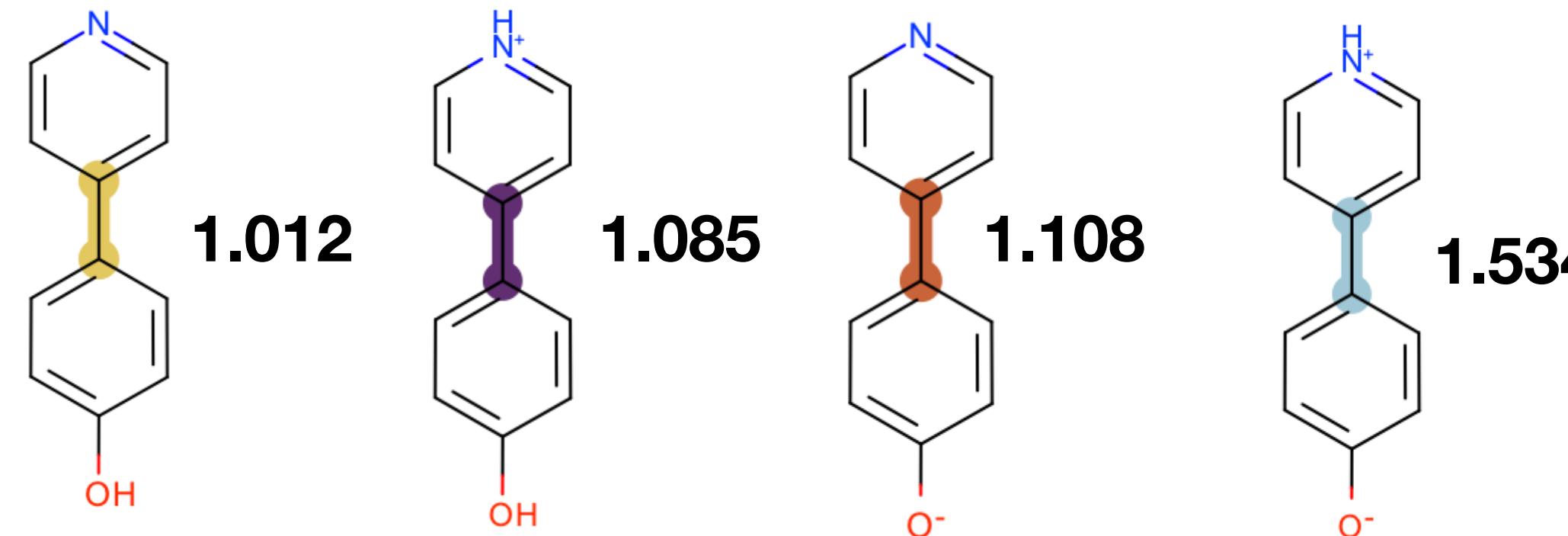
# The Wiberg Bond Order is a cheap measure of electron population overlap between two atoms

$$W_{AB} = \sum_{\mu \in A} \sum_{\nu \in B} |D_{\mu\nu}|^2.$$

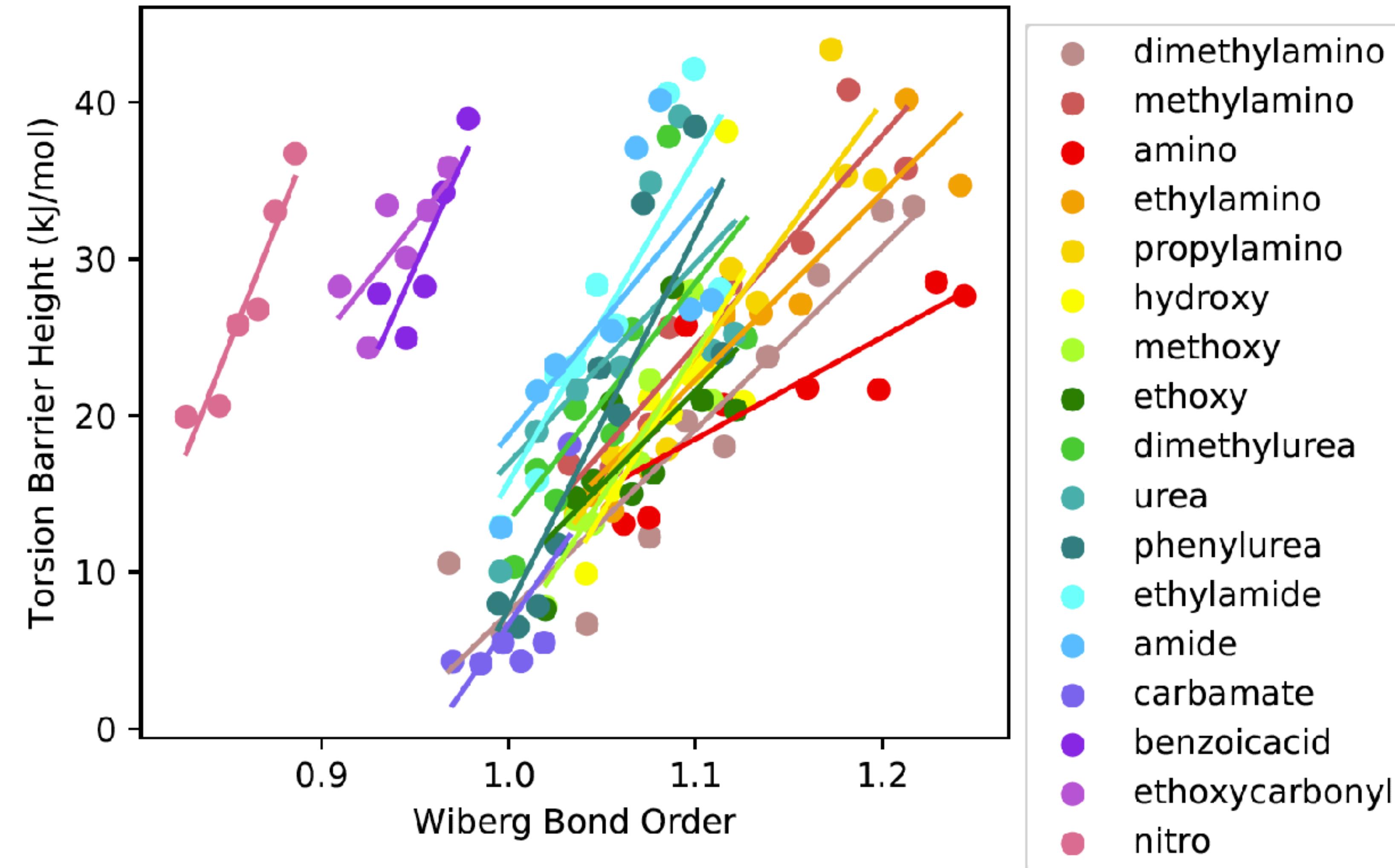


$$D_{\mu\nu} = 2 \sum_i^{\text{occ.}} C_{\mu i} C_{\nu i}^*,$$

# Torsion potential barrier heights are linear with Wiberg Bond Orders

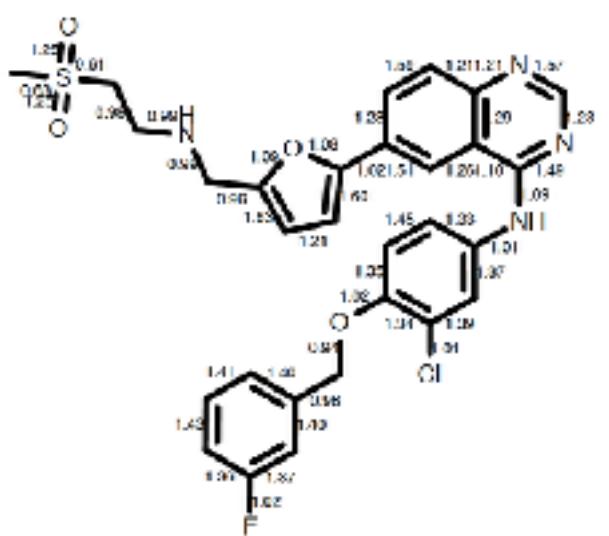


# Torsion barrier height increases with increasing Wiberg Bond Order for other functional groups



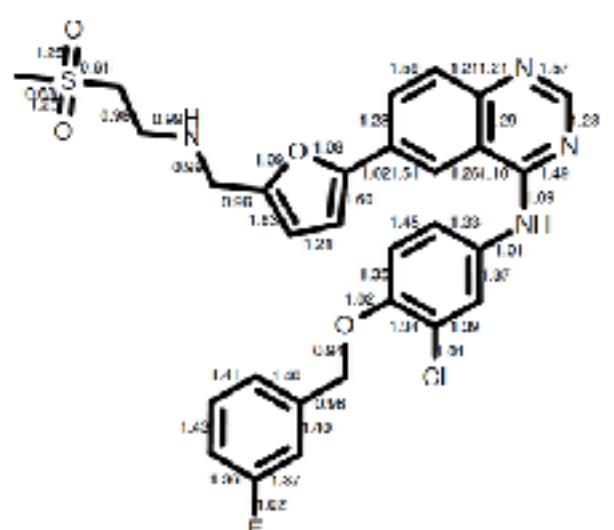
# Intelligent **fragmentation** can reduce the misrepresentation of torsions in QM database

Calculate Wiberg Bond Order  
from AM1 calculation

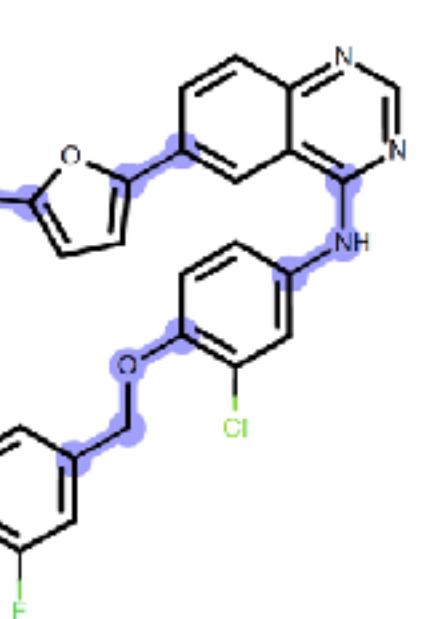


# Intelligent fragmentation can reduce the misrepresentation of torsions in QM database

Calculate Wiberg Bond Order  
from AM1 calculation



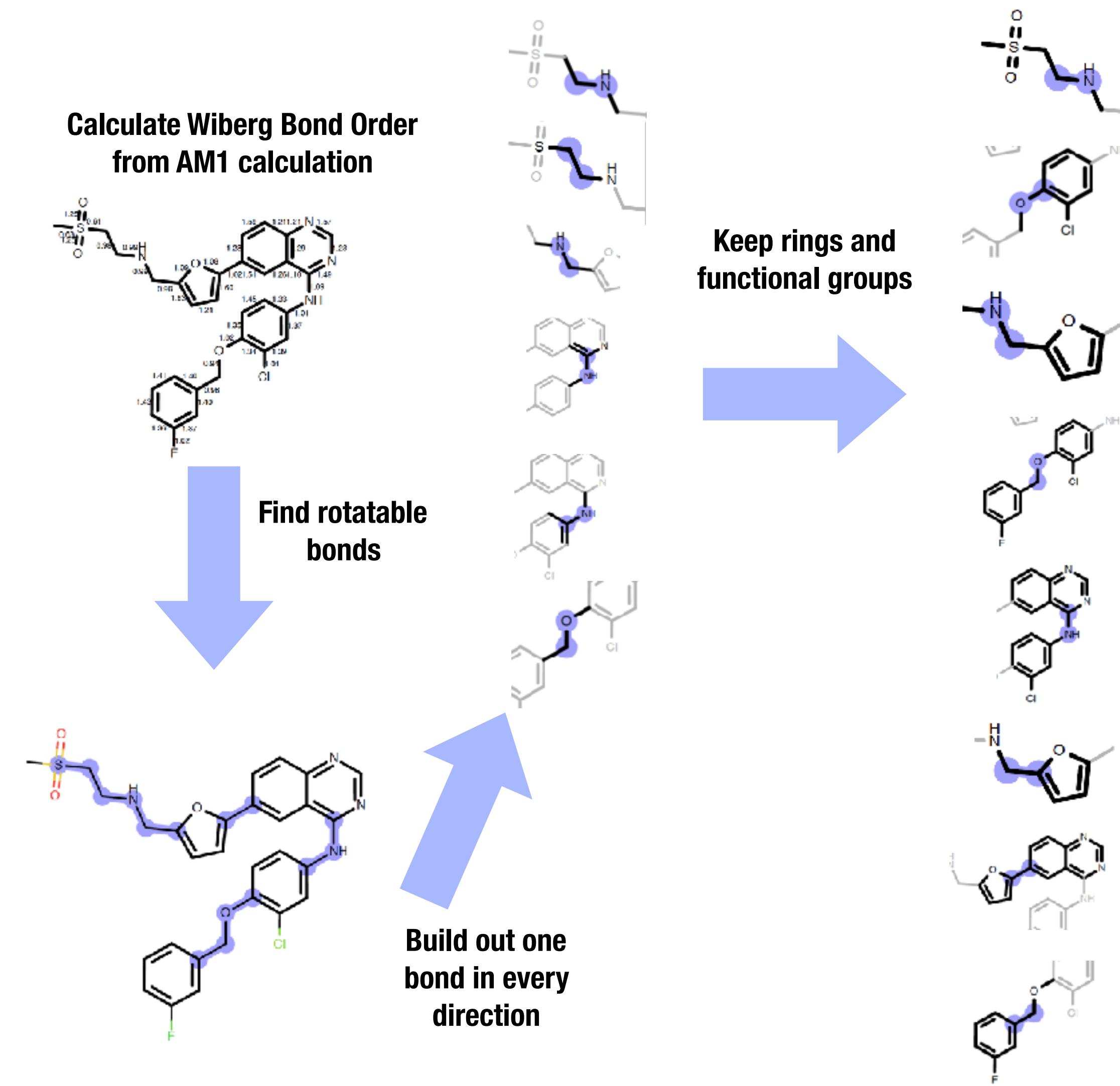
Find rotatable  
bonds



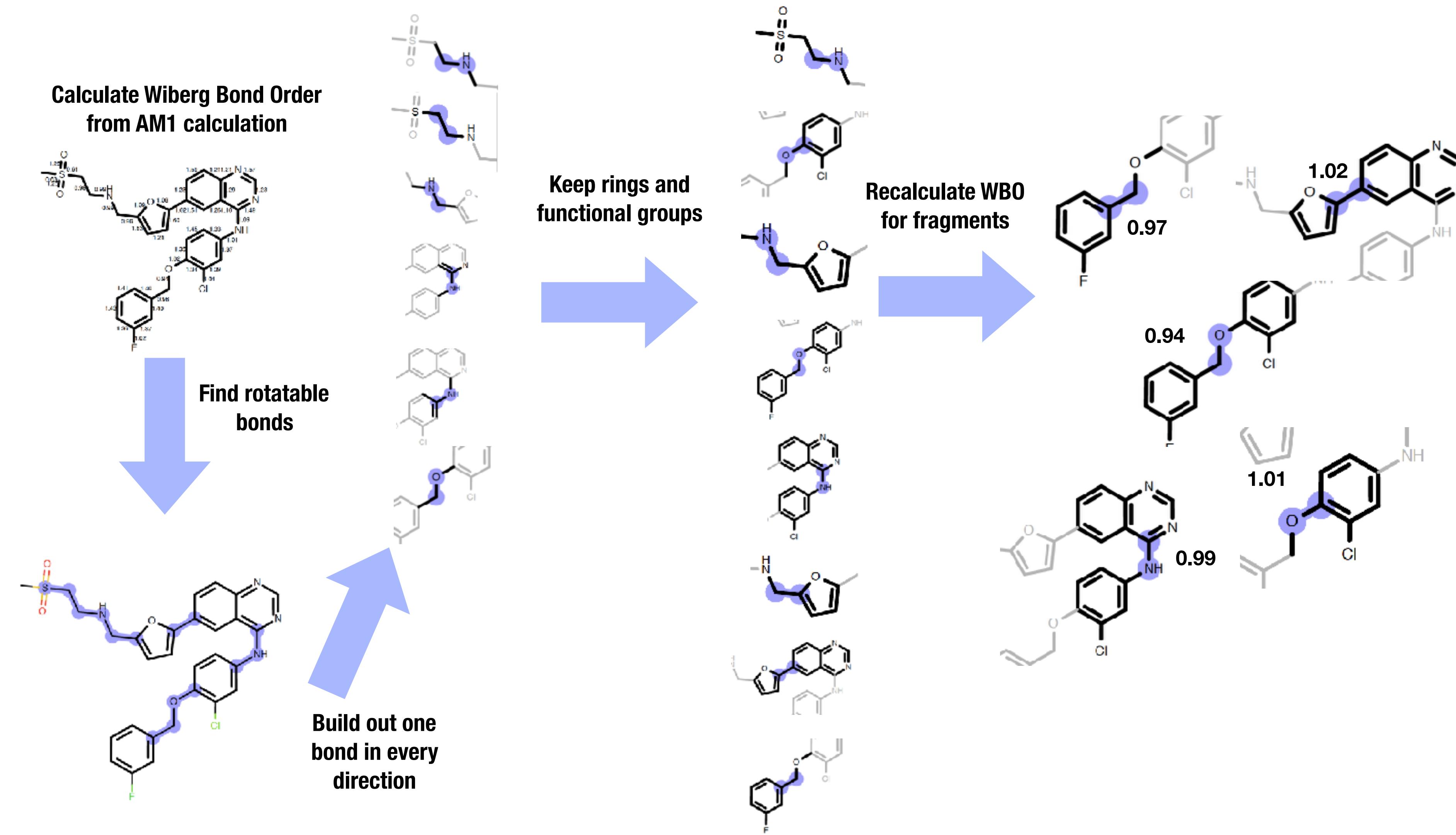
# Intelligent fragmentation can reduce the misrepresentation of torsions in QM database



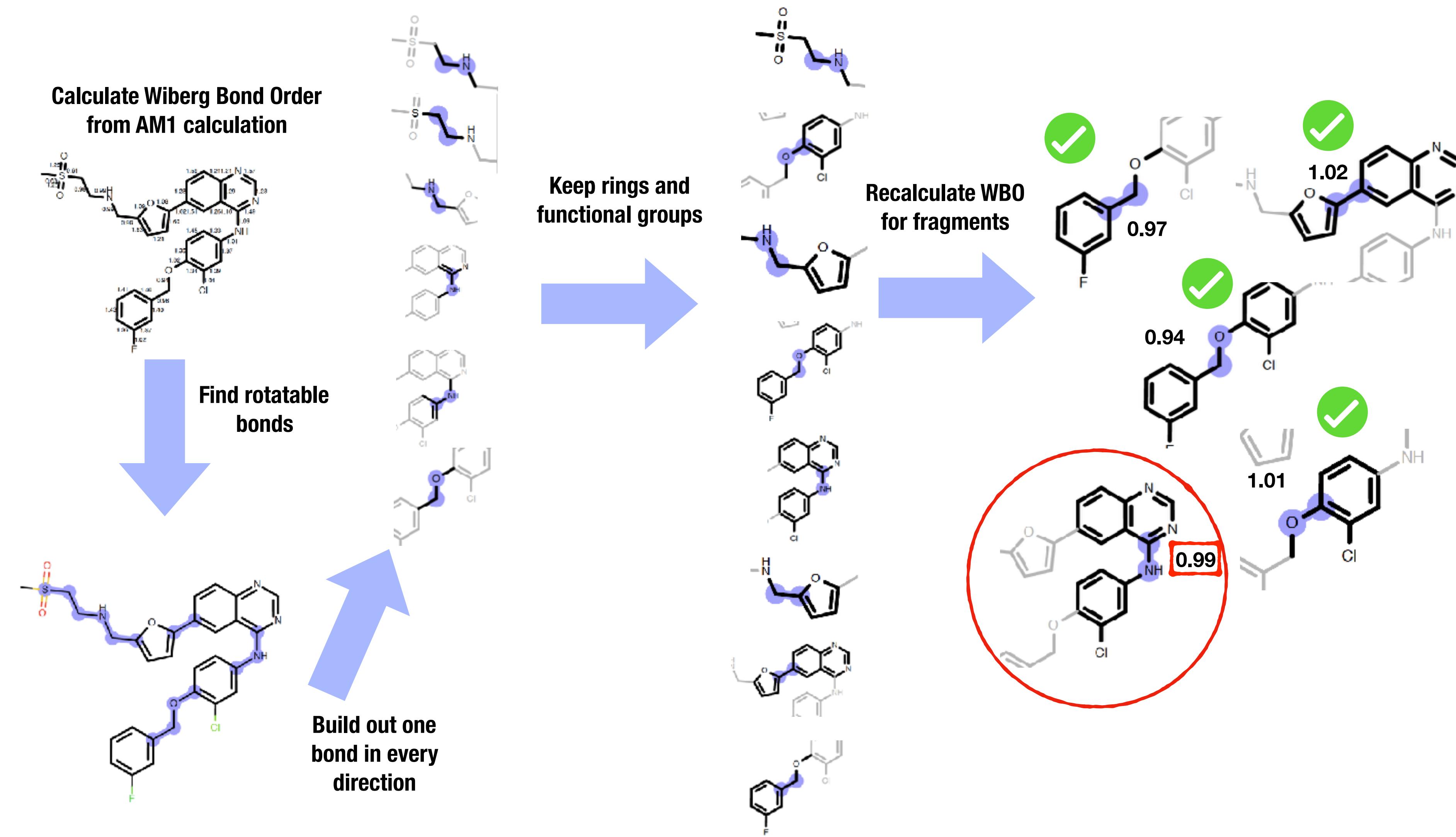
# Intelligent fragmentation can reduce the misrepresentation of torsions in QM database



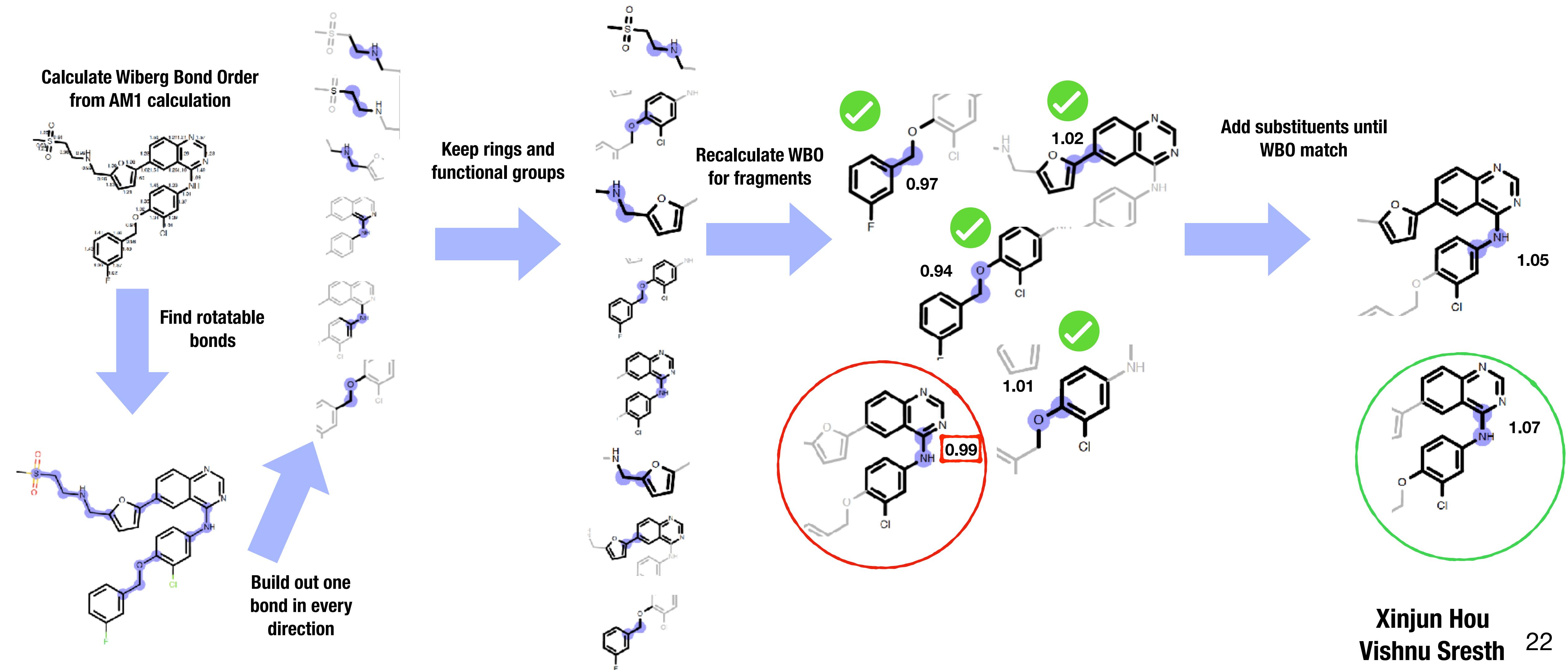
# Intelligent fragmentation can reduce the misrepresentation of torsions in QM database



# Intelligent fragmentation can reduce the misrepresentation of torsions in QM database



# Intelligent fragmentation can reduce the misrepresentation of torsions in QM database



# OFF is pursuing several routes to integrate **open source** AM1 WBO calculations.

## 1. AmberTools SQM

Dustin Tracy from the Roitberg lab added WBO to SQM.

Needs testing and integration into an AmberTool conda package so it can be integrated with OFFToolKit

## 2. Graph Nets trained on AM1 WBO

Can potentially be a lot cheaper than AM1



Yuanqing Wang

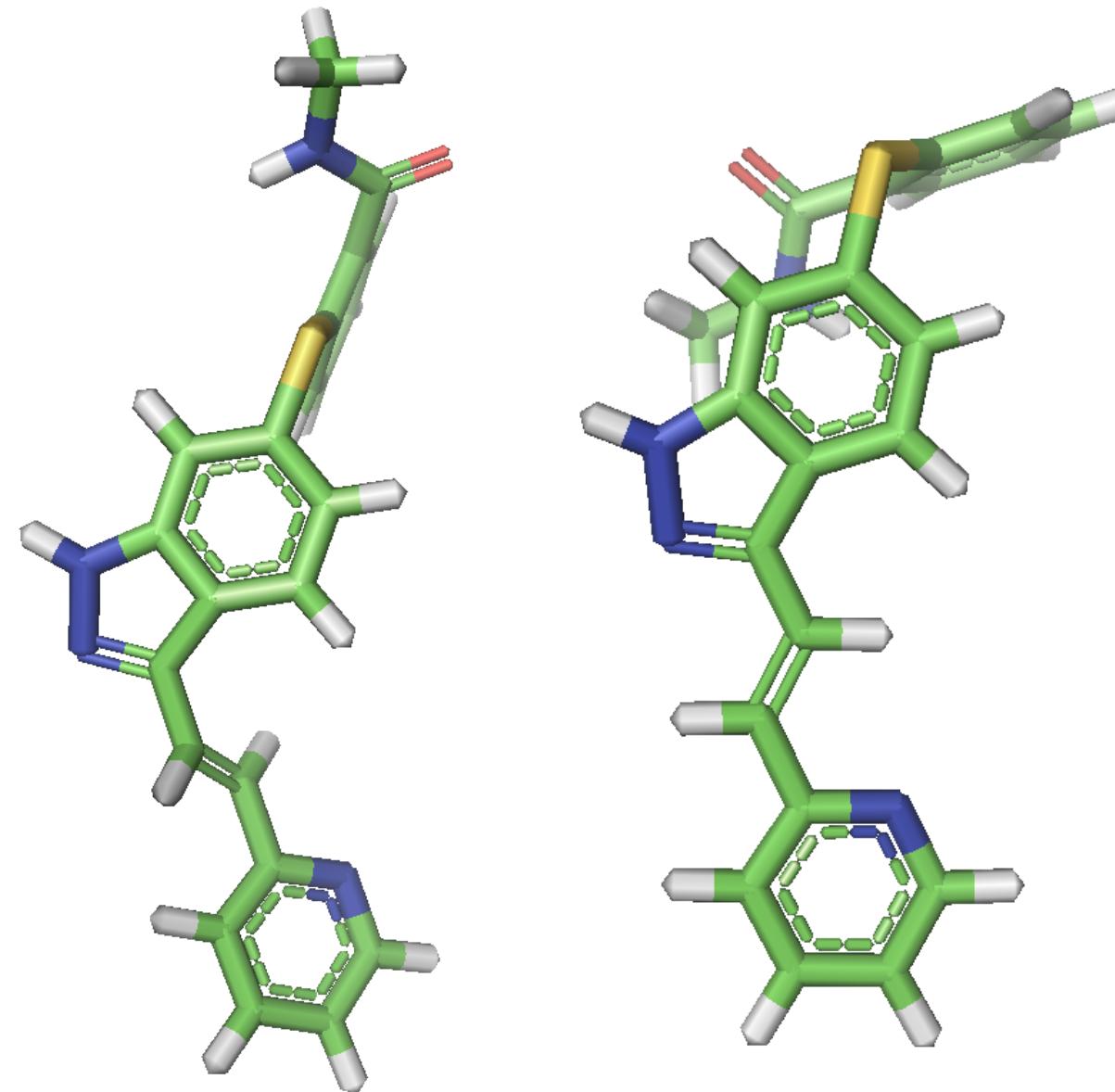


Adrian Roitberg

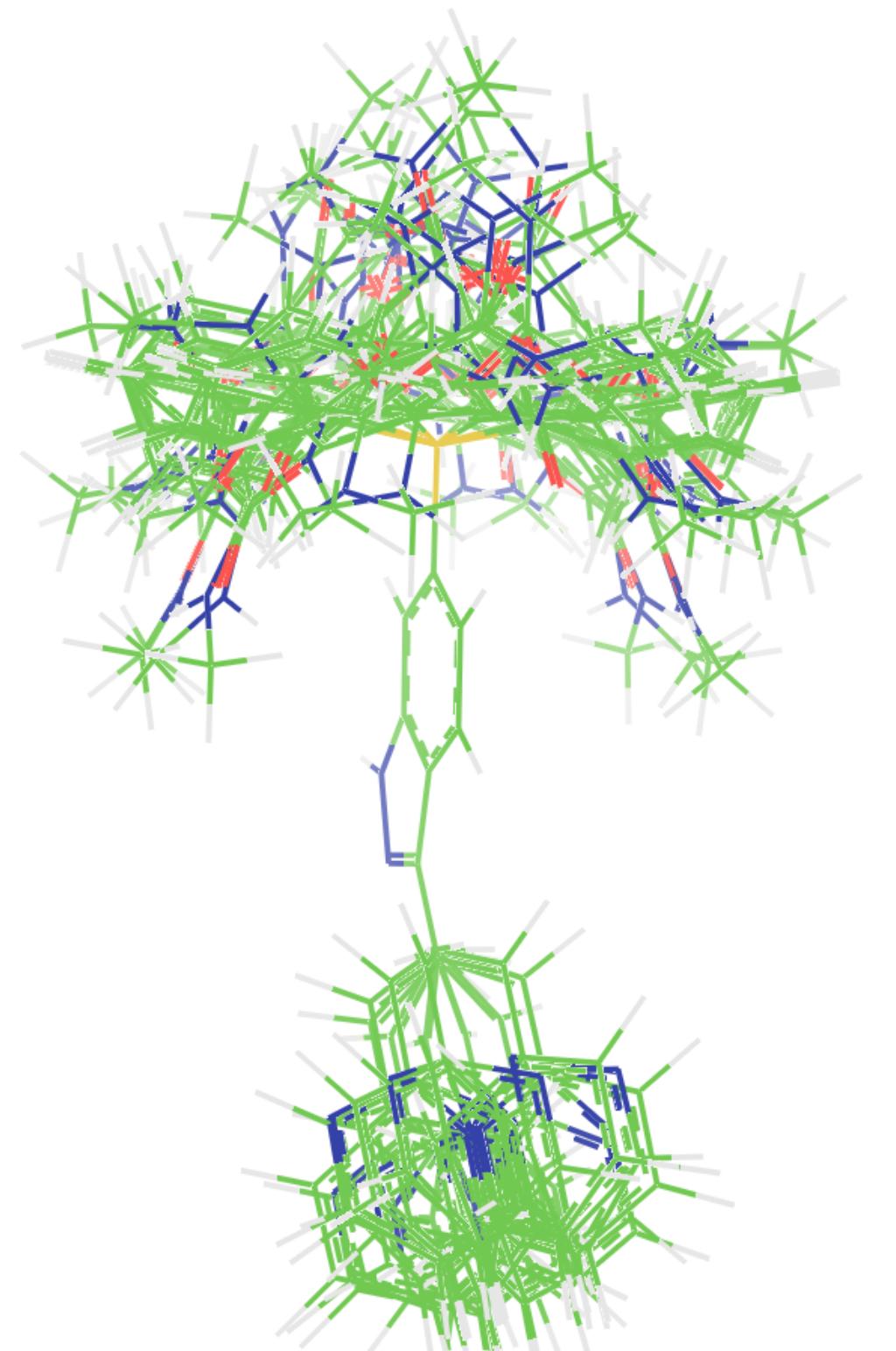


Dustin Tracy

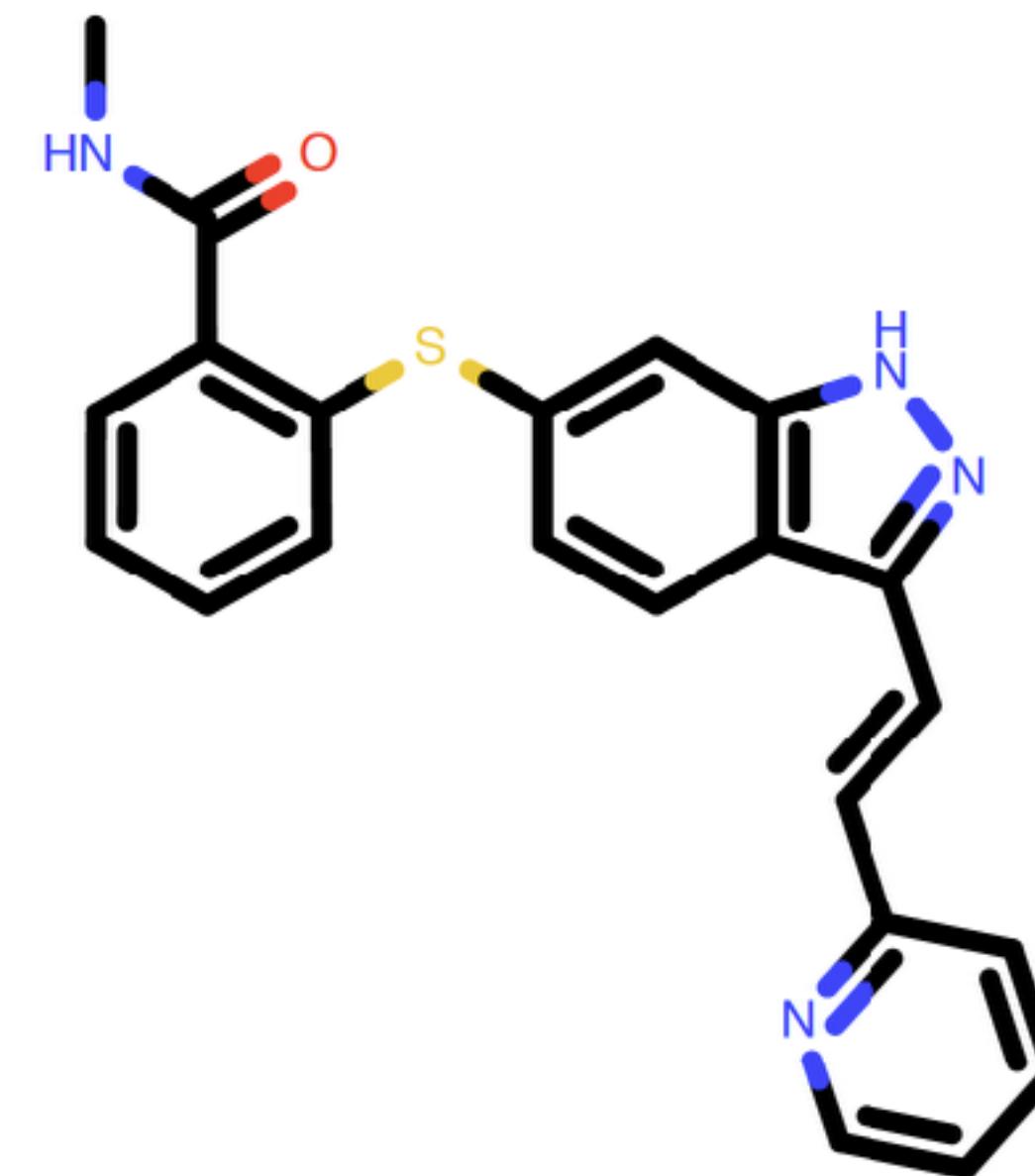
# Indexing molecules for QC database



Quantum chemistry represents molecules by their coordinates.

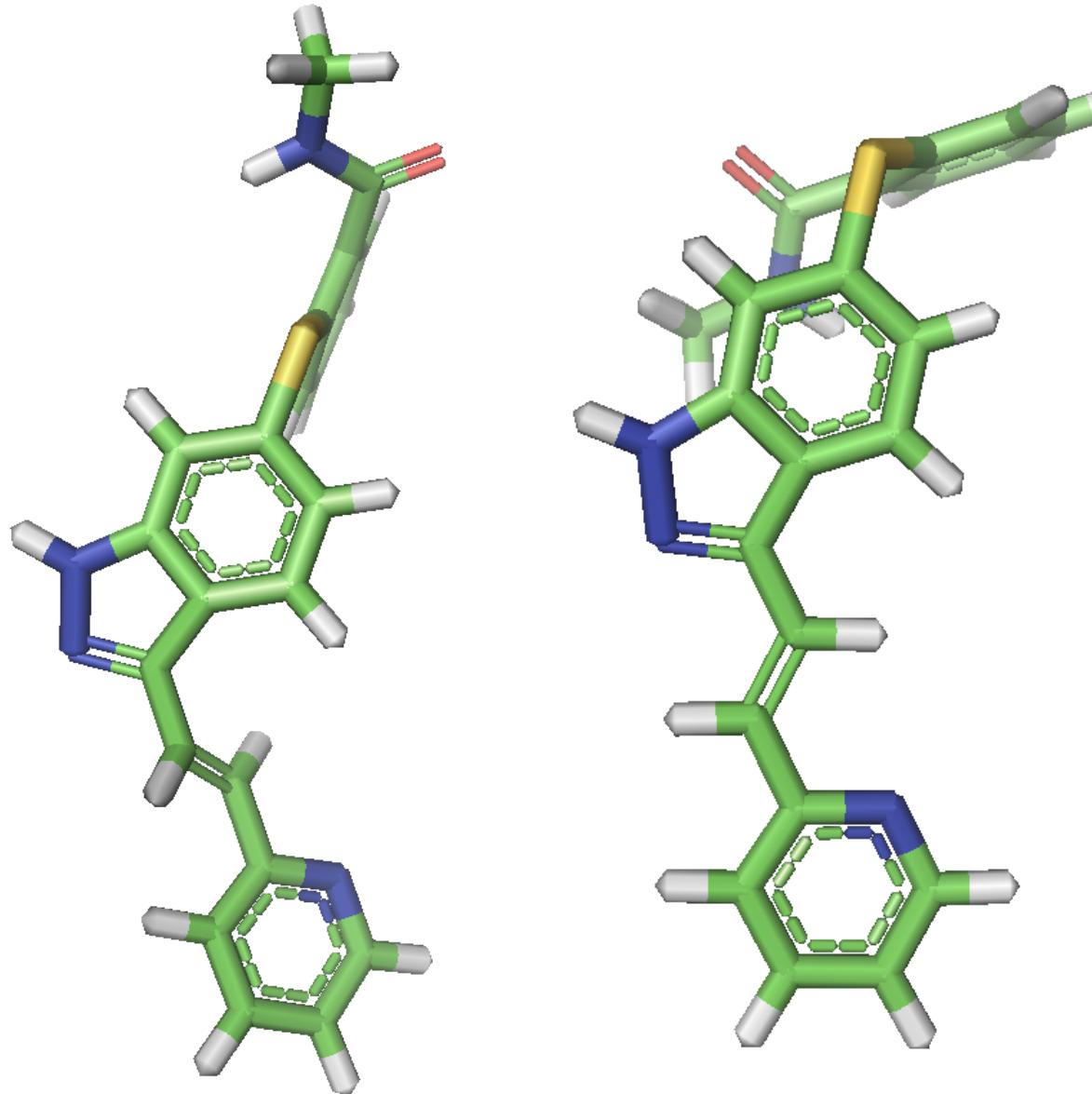


Molecular mechanics represents molecules as conformational distributions

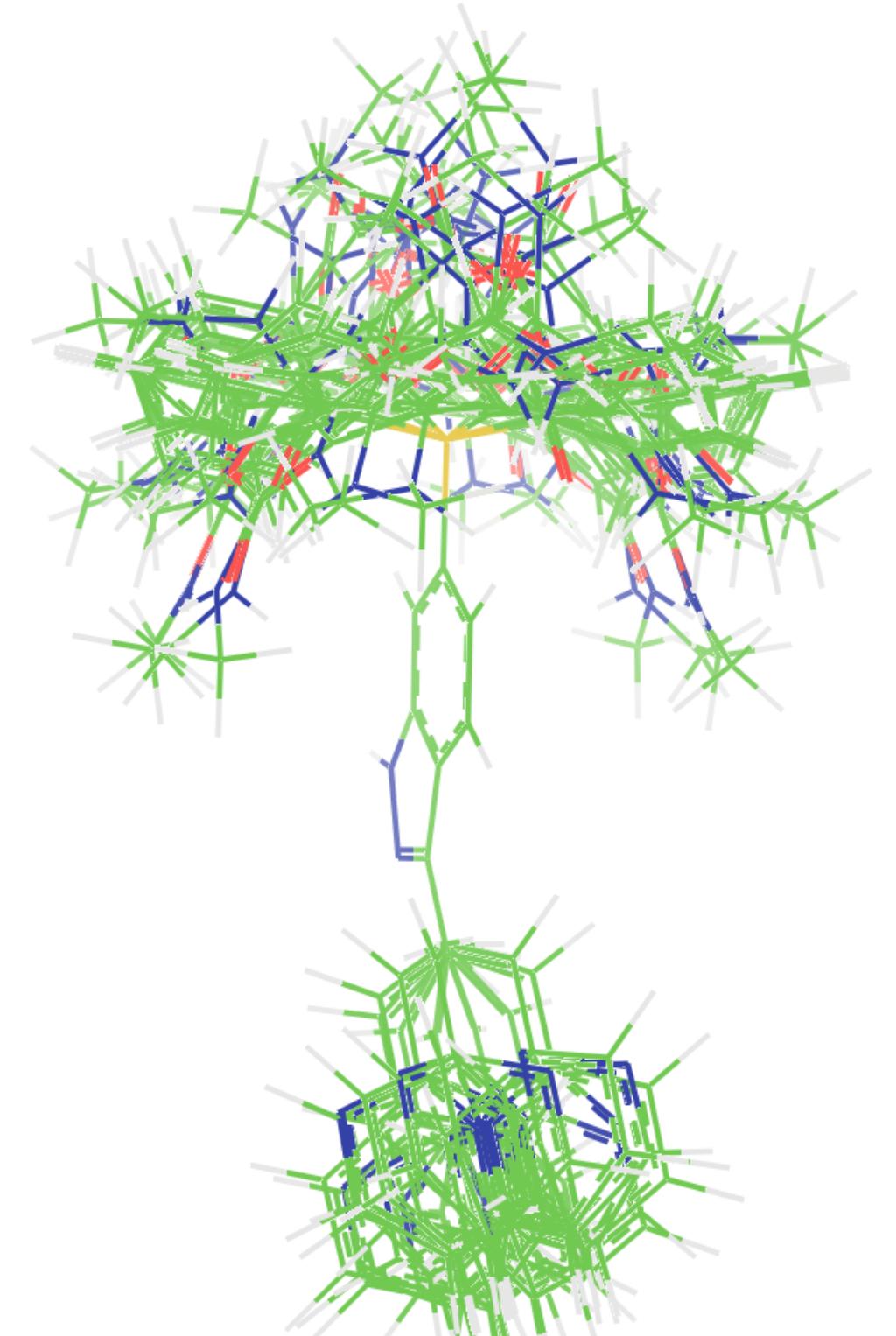


Cheminformatics represent molecules as graphs

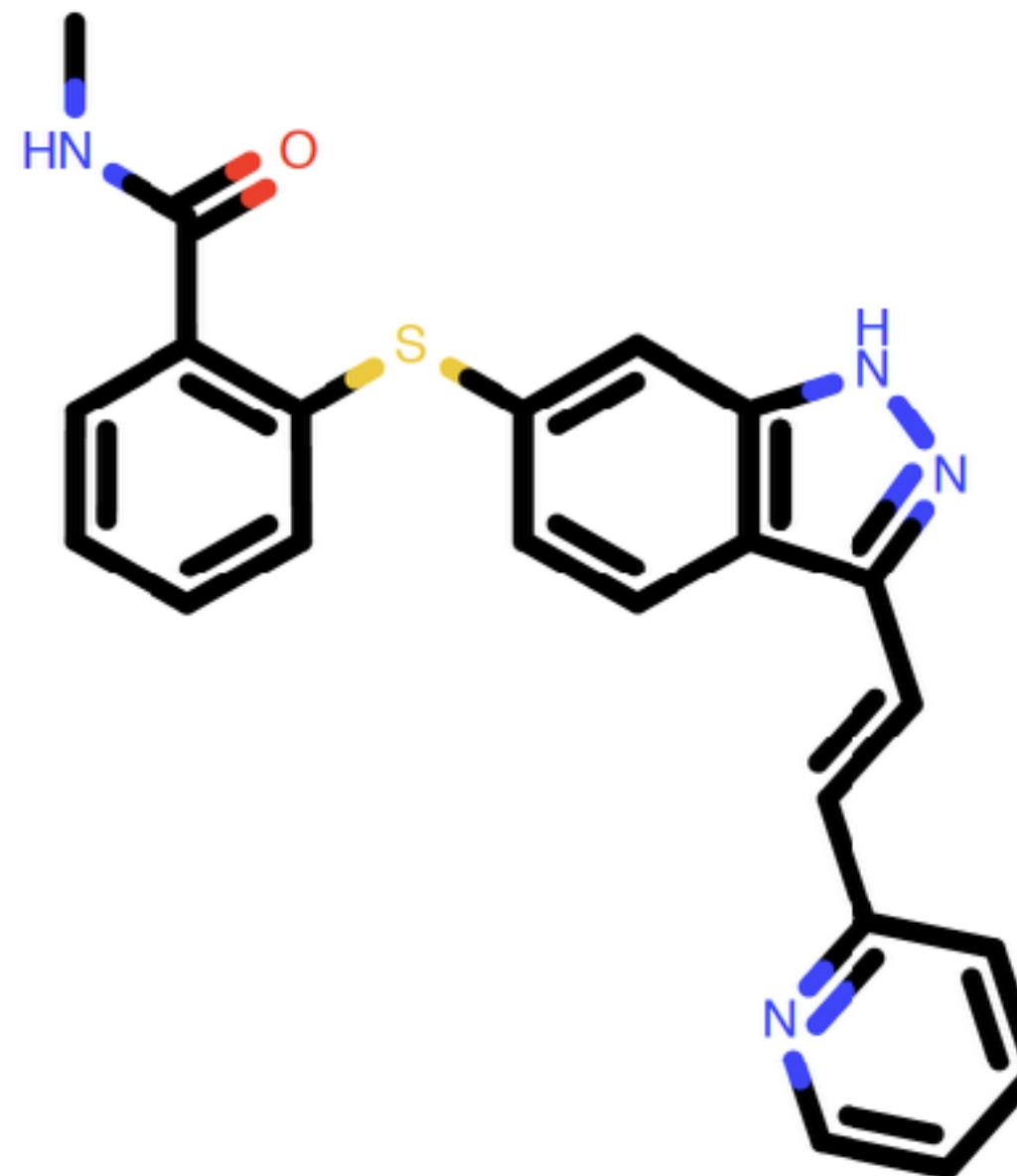
# CMILES indices link different representation of molecules



Quantum chemistry represents molecules by their coordinates.



Molecular mechanics represents molecules as conformational distributions



Cheminformatics represent molecules as graphs

**SMILES** and **InChI** are not attached to coordinates so calculations with different geometries can be grouped together

**SMILES:** Simplified Molecular Input Line Entry Specification

**InChI:** The IUPAC International Chemical Identifier

**cmiles provides indices that ensure broad **usability** and **sustainability** of the database**

**cmiles**

---



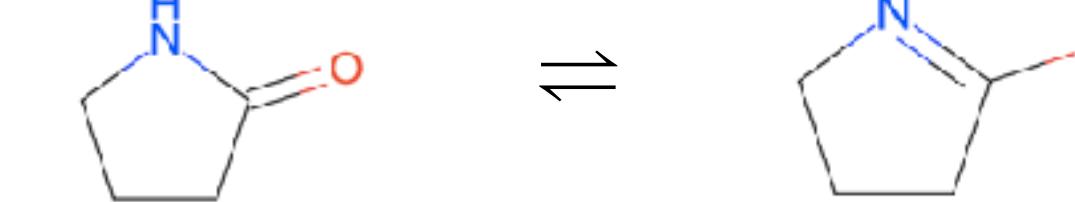
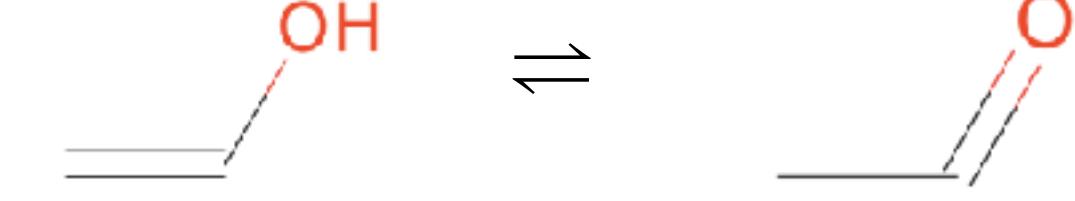
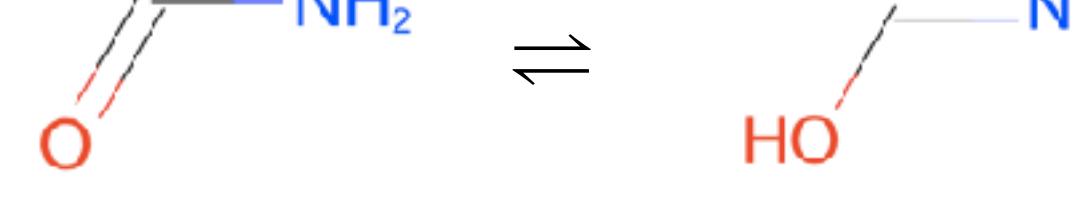
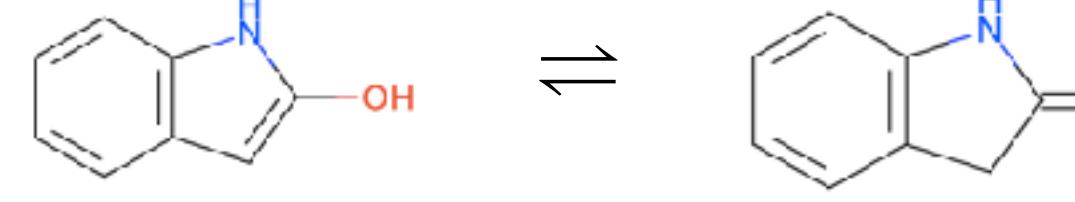
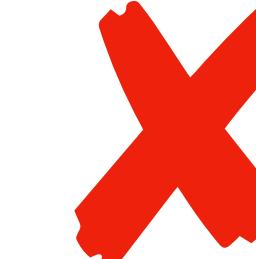
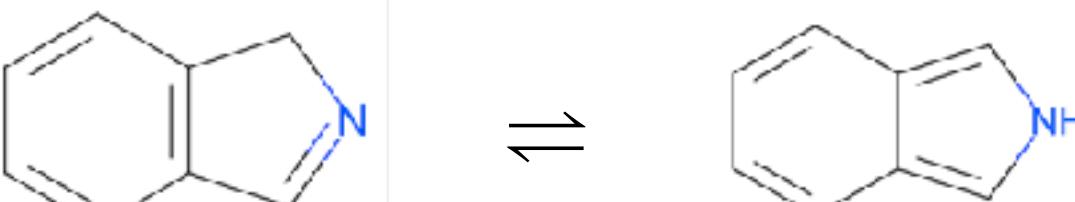
**SMILES must be canonical to avoid redundancy and search failures**

**Canonical SMILES are only canonical with respect to toolkit **and** toolkit version**

**cmiles is set up to test if updates in dependent cheminformatics toolkits changed the canonicalization algorithm.**

**~ 0.5 % of SMILES tested changed with RDKit updates - mostly stereochemistry**

# Standardizing molecules for QCArchive

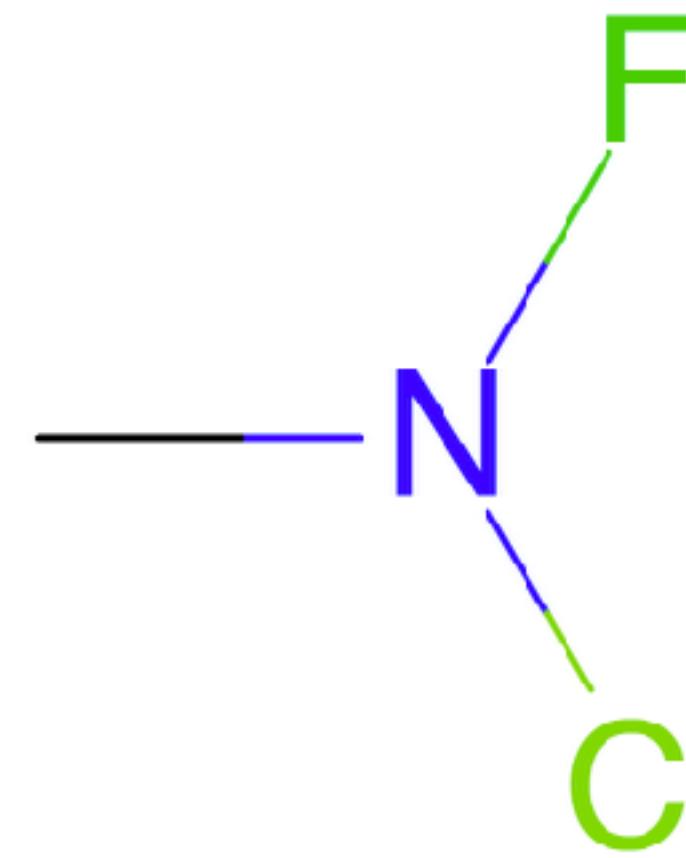
	Tautomer	InChI	RDKit	OpenEye
<b>lactam - lactim</b>				
<b>keto - enol</b>				
<b>amide - imidic acid</b>				
<b>indole-2-ol - indole-2-one</b>				
<b>2H-isoindole - 1H-isoindole</b>				

# OFF Toolkit requires fully defined **stereochemistry** to avoid incorrect partial charges and/or partial bond orders

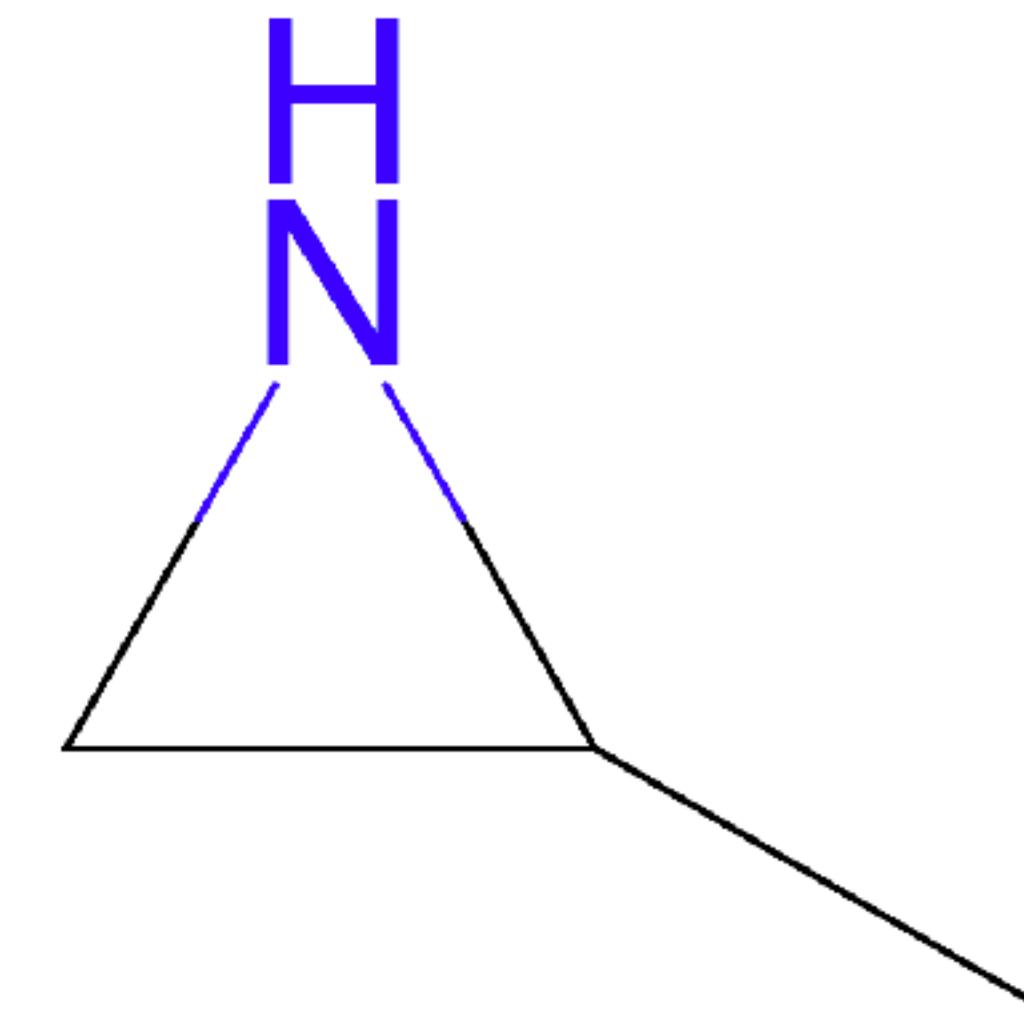
## Why is this an issue?

Stereochemical differences can lead to different partial charges or partial bond orders and potentially different parameter assignments

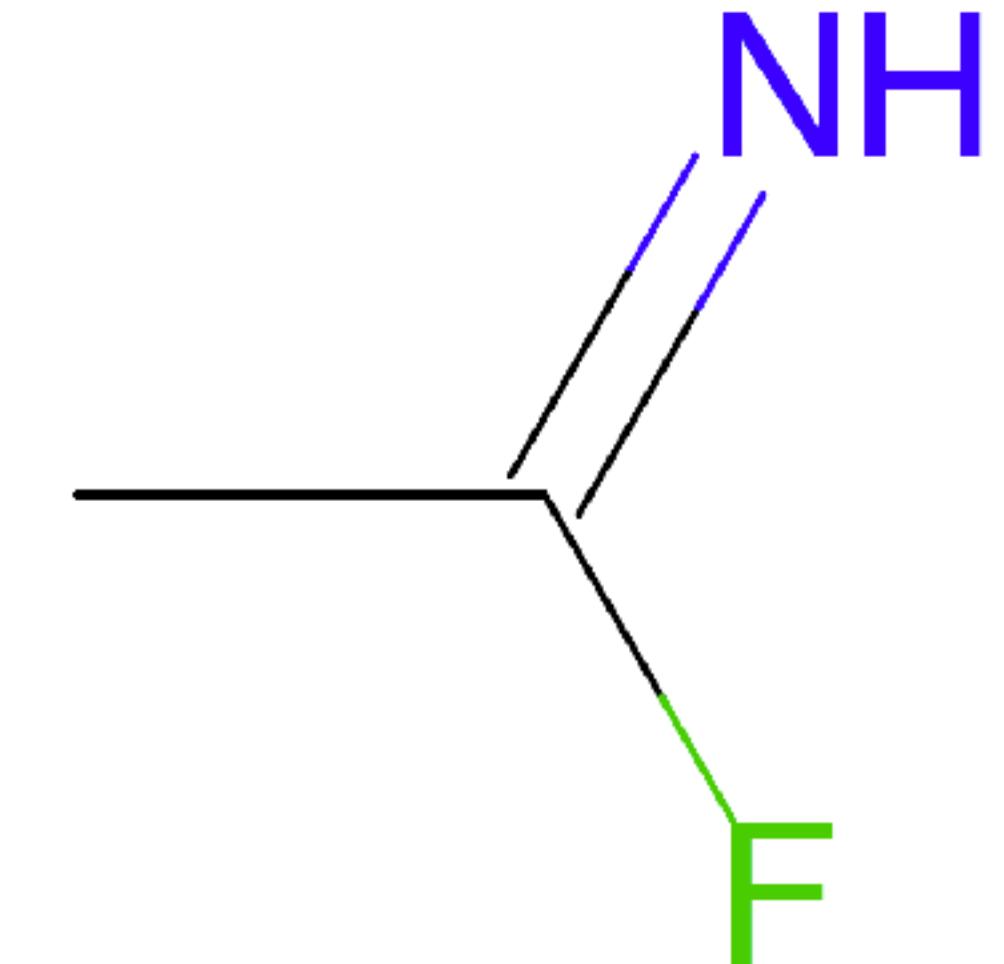
pyramidal nitrogen



N in three membered ring



primary imine

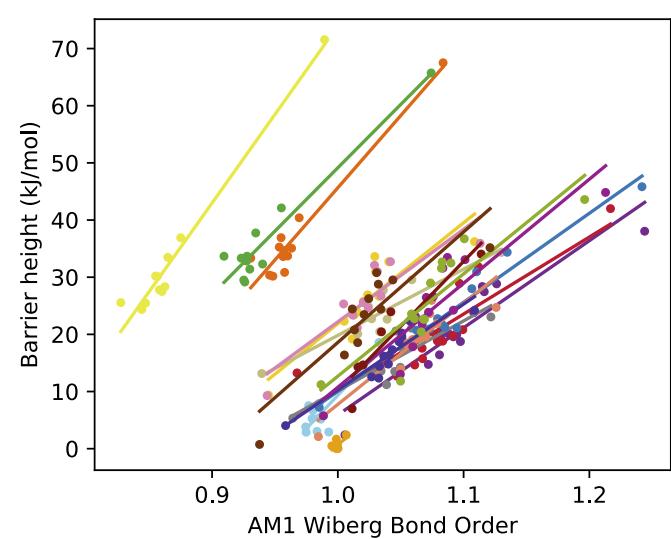


Conversion timescale:  $10^3 - 10^5$  per sec

OE does not consider this stereogenic

RDKit flags as potential stereogenic if hydrogen is explicit

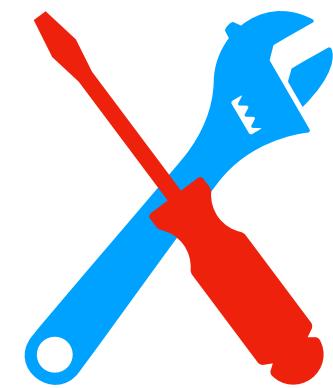
# Summary



1. WBO is a powerful, cheap measure that provides robust signal of a bond's chemical environment and can be used to automate fragmentation and potentially finding rotatable bonds.
2. Automatic submission to QCArchive and massively parallel computation allows multidimensional torsion scan
3. The QCArchive OFF datasets are indexed with multiple canonical identifiers with their provenance



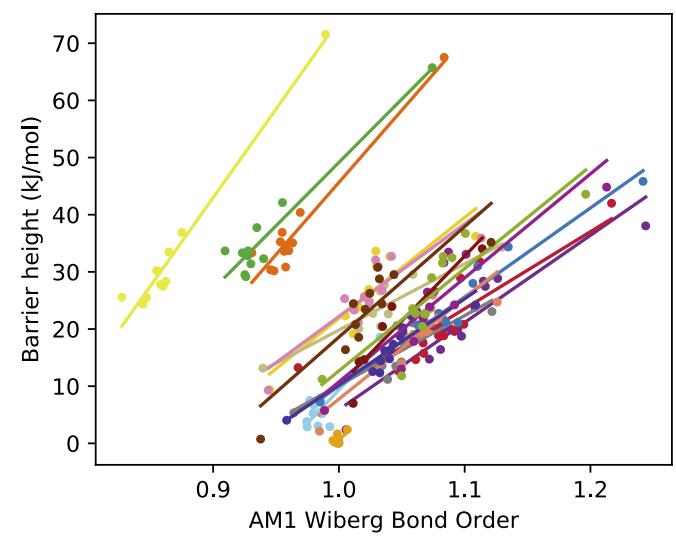
# Future Directions



- Integrate AM1 WBO and `fragmenter` with the rest of the OFF software stack

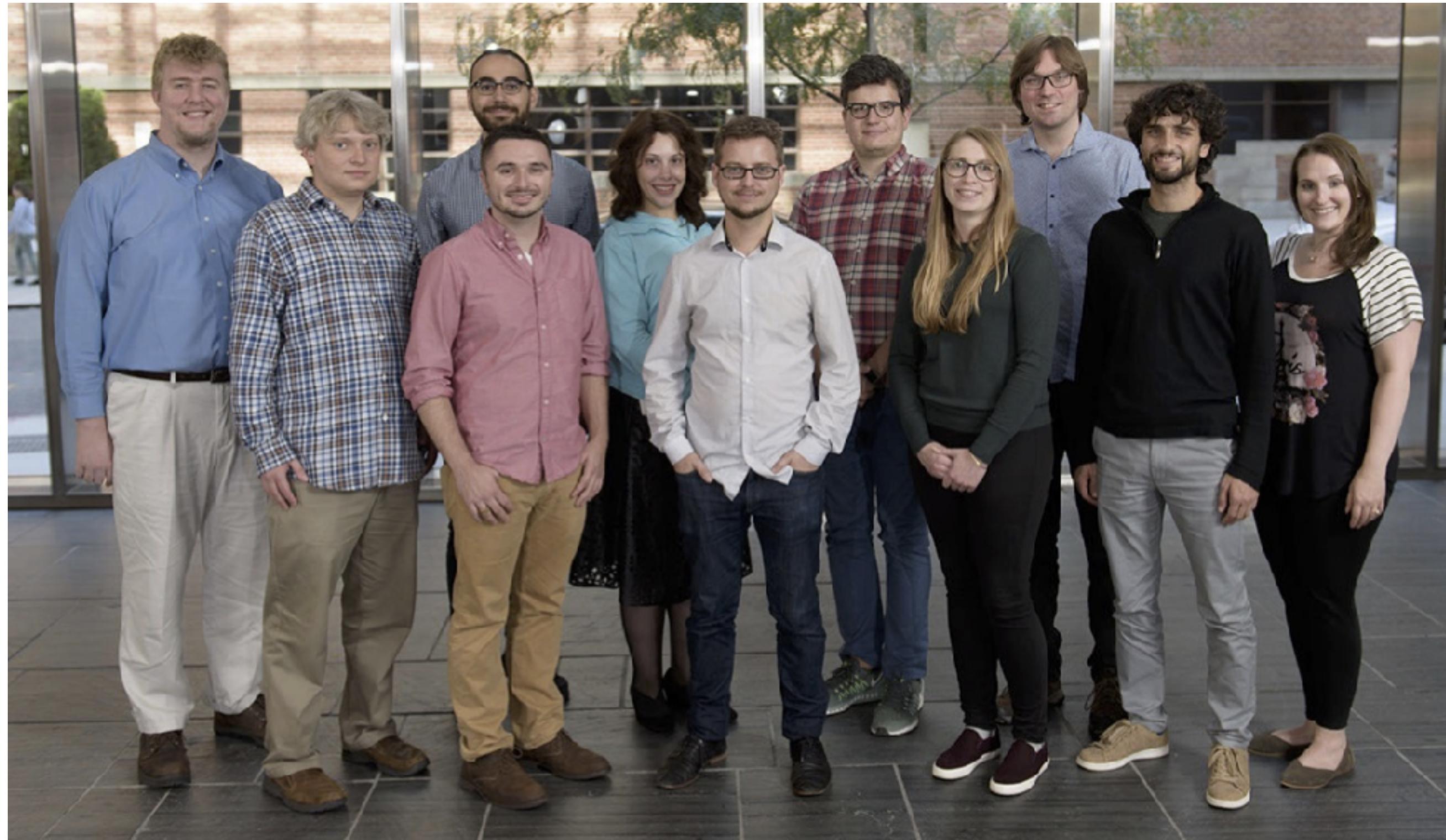


- Optimize fragmentation scheme via learning patterns not to fragment from the validation set



- Generalize the torsion barrier heights vs WBO relationship and interpolate torsion force constants

# Acknowledgment



**John Chodera**

**Josh Fass**

**Mehtap Işık**

**Andrea Rizzi**

**Bas Rustenberg**

**Marcus Wieder**

**Simon Boothroyd**

**Rafal Wiewiora**

**Melissa Boby**

**OpenEye**

**Christopher Bayly**

**#Torsions channel**

**Alberto Gobbi**

**Adrian Roitberg**

**Lee-Ping Wang**

**Yudong Qiu**

**David Mobley**

**Caitlin Bannan**

**MolSSI**

**Daniel Smith**

**Doaa Altarawy**

**Levi Naden**



<https://qcarchive.molssi.org/>

## Funding



# Questions

