

State-of-the-Art Computational Approaches for Characterizing Reaction Properties and Predicting Unknown Reactivity

Brett M. Savoie

Davidson Associate Professor of Chemical Engineering,
Purdue University

Students: Qiyuan Zhao, Tyler Pasut

P2SAC Fall Conference, Purdue University, 12/14/22

“Known Unknowns” and “Unknown Unknowns”

$A \rightarrow B$

- To safely plan a known reaction, we need access to solid thermodynamic data (e.g., ΔH_f , S° , C_v) to understand and classify risks.
- This is a “known unknown” in that we know the reaction, $A \rightarrow B$, but we need values for a few unknown variables.

$A \rightarrow ? \rightarrow B$; $A \rightarrow B + ?$; $A \rightarrow ?$

- $A \rightarrow ? \rightarrow B$, means that we know the net reaction, but there may be a consequential (e.g., potentially reactive) intermediate. Even if we have accurate thermodynamic data on A/B, neglecting the intermediate could be disastrous.
- The $A \rightarrow B + ?$ (unknown side-reaction) and $A \rightarrow ?$ (unknown main product), problems have similar “unknown unknown” characteristics.

“Known Unknowns” and “Unknown Unknowns”



TAFFI Component Increment Theory (TCIT)

- This is a “known unknown” in that we know the reaction, $A \rightarrow B$, but we need values for a few unknown variables.



Yet Another Reaction Program (YARP)

- The $A \rightarrow ? \rightarrow B$ means that we know the net reaction, but there may be a consequential (e.g., potentially reactive) intermediate. Even if we have accurate thermodynamic data on A/B, neglecting the intermediate could be disastrous.
- The $A \rightarrow B + ?$ (unknown side-reaction) and $A \rightarrow ?$ (unknown main product), problems have similar “unknown unknown” characteristics.

Challenges of Contemporary Group Theories

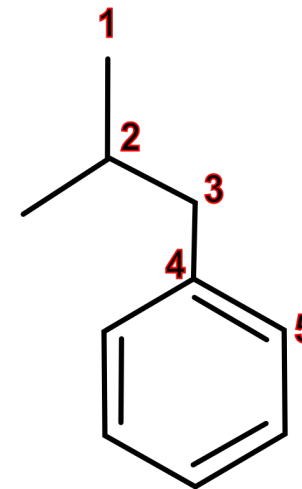
Benson Group Theory:

- The idea is to decompose molecular properties (ΔH_f , S° , C_v) as the sum of “group” contributions.
- Group contributions are calculated based on trusted experimental or computational data, and transferability is assumed.

Problems we want to address:

- **Specificity:** the definition of a “group” has never been formalized and inconsistent granularity is applied.
- **Provenance:** inconsistent thermodynamic data is available/used to determine group contributions.
- **Extensibility:** because of the provenance and specificity problems, it isn't possible to develop new groups in a consistent way.

From Anslyn and
Dougherty's
Textbook



1) C -(C)(H) ₃	2(-10.20)
2) C -(C) ₃ (H)	-1.90
3) C -(C _B)(C)(H) ₂	-4.86
4) C _B -(C)	5.51
5) C _B -(H)	5(3.30)

-5.15 kcal/mole
(-21.6 kJ/mole)

Experimental ΔH_f : -5.15 +/- 0.34 kcal/mol

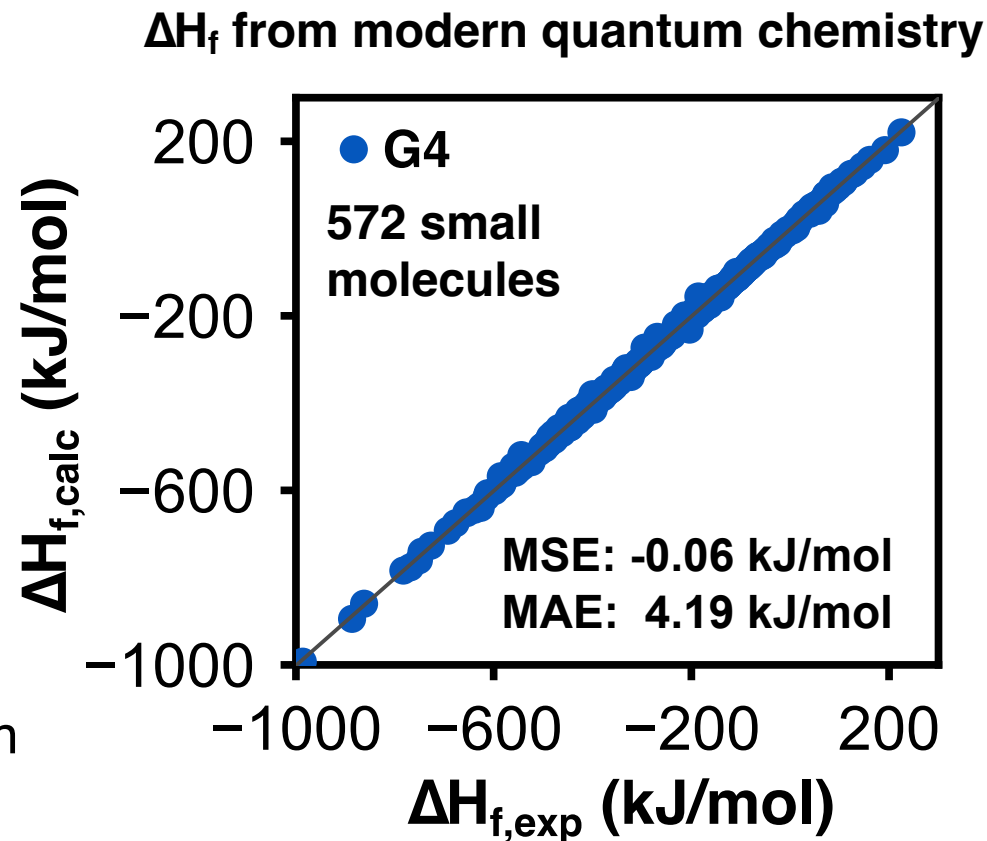
Challenges of Contemporary Group Theories

Benson Group Theory:

- The idea is to decompose molecular properties (ΔH_f , S° , C_v) as the sum of “group” contributions.
- Group contributions are calculated based on trusted experimental or computational data, and transferability is assumed.

Problems we want to address:

- **Specificity:** the definition of a “group” has never been formalized and inconsistent granularity is applied.
- **Provenance:** inconsistent thermodynamic data is available/used to determine group contributions.
- **Extensibility:** because of the provenance and specificity problems, it isn't possible to develop new groups in a consistent way.



Zhao, Q.; Savoie, B. M.; Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory. *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Challenges of Contemporary Group Theories

Benson Group Theory:

- The idea is to decompose molecular properties (ΔH_f , S° , C_v) as the sum of “group” contributions.

- Group
- on trust
- data, a

Prob

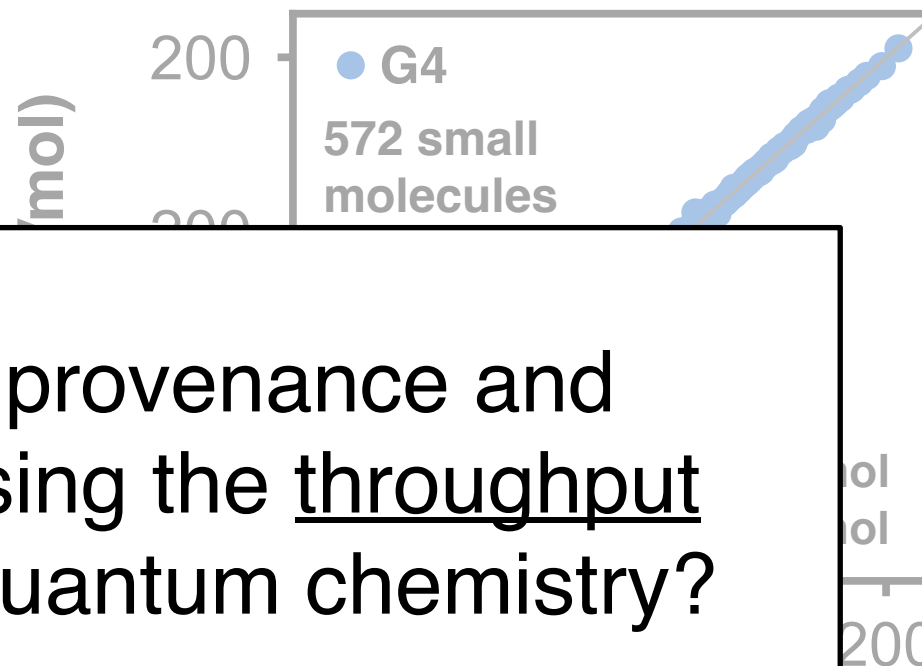
- Spec
- formal

Can we circumvent the provenance and extensibility challenges using the throughput and accuracy of modern quantum chemistry?

- **Provenance:** inconsistent thermodynamic data is available/used to determine group contributions.

- **Extensibility:** because of the provenance and specificity problems, it isn't possible to develop new groups in a consistent way.

ΔH_f from modern quantum chemistry



Zhao, Q.; Savoie, B. M.; Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory. *J. Chem. Info. Model.* **2020**, 60, 2199-2207

TAFFI Component Increment Theory (TCIT)

The fundamental idea

- Systematize component-definitions and model compound selection with rigorous graph-based typing.

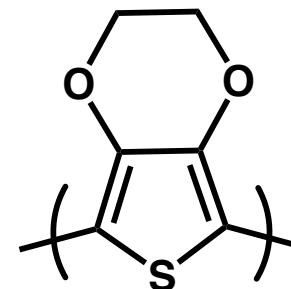
Zhao, Q.; Savoie, B. M.; "Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory". *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Zhao, Q.; Iovanac, N.; Savoie, B. M.; "Transferable Ring Corrections for Predicting Enthalpy of Formation of Cyclic Compounds" *J. Chem. Info. Model.* **2021**, 61, 5013-5027

Seo, B.; Lin, Z.-Y.; Zhao, Q.; Webb, M. A.; Savoie, B. M. Topology Automated Force-Field Interactions (TAFFI): A Framework for Developing Transferable Force Fields. *J. Chem. Inf. Model.* **2021**, 61 (10), 5013–5027. <https://doi.org/10.1021/acs.jcim.1c00491>.

**P2SAC
Publications**

TCIT is a component theory
(2-bond specific)



Topology Automated
Force Field Interactions



graph/structure
equivalence



S	0	1	0	0	1	0	0	0	0	0	0	0	0
C	1	0	1	0	0	0	0	0	0	0	0	0	0
C	0	1	0	1	0	1	0	0	0	0	0	0	0
C	0	0	1	0	1	0	0	0	0	0	1	0	0
C	1	0	0	1	0	0	0	0	0	0	0	0	0
O	0	0	1	0	0	0	1	0	0	0	0	0	0
C	0	0	0	0	0	1	0	1	1	1	0	0	0
C	0	0	0	0	0	0	1	0	0	0	1	1	1
H	0	0	0	0	0	0	1	0	0	0	0	0	0
H	0	0	0	0	0	0	1	0	0	0	0	0	0
O	0	0	0	1	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0

**Adjacency
matrix for
PEDOT
monomer**

TAFFI Component Increment Theory (TCIT)

The fundamental idea

- Systematize component-definitions and model compound selection with rigorous graph-based typing.

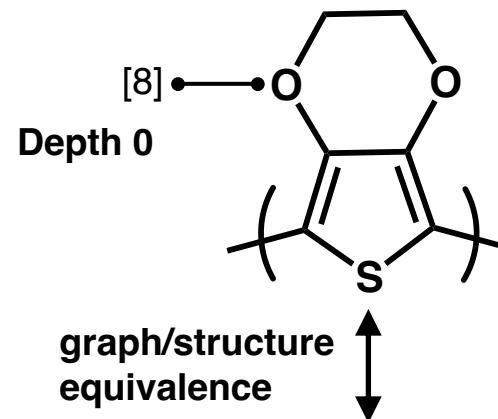
Zhao, Q.; Savoie, B. M.; "Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory". *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Zhao, Q.; Iovanac, N.; Savoie, B. M.; "Transferable Ring Corrections for Predicting Enthalpy of Formation of Cyclic Compounds" *J. Chem. Info. Model.* **2021**, 61, 5013-5027

Seo, B.; Lin, Z.-Y.; Zhao, Q.; Webb, M. A.; Savoie, B. M. Topology Automated Force-Field Interactions (TAFFI): A Framework for Developing Transferable Force Fields. *J. Chem. Inf. Model.* **2021**, 61 (10), 5013–5027. <https://doi.org/10.1021/acs.jcim.1c00491>.

**P2SAC
Publications**

TCIT is a component theory
(2-bond specific)



Topology Automated
Force Field Interactions



S	0	1	0	0	1	0	0	0	0	0	0	0	0
C	1	0	1	0	0	0	0	0	0	0	0	0	0
C	0	1	0	1	0	1	0	0	0	0	0	0	0
C	0	0	1	0	1	0	0	0	0	0	1	0	0
C	1	0	0	1	0	0	0	0	0	0	0	0	0
O	0	0	1	0	0	0	1	0	0	0	0	0	0
C	0	0	0	0	0	1	0	1	1	1	0	0	0
C	0	0	0	0	0	0	1	0	0	0	1	1	1
H	0	0	0	0	0	0	1	0	0	0	0	0	0
H	0	0	0	0	0	0	1	0	0	0	0	0	0
O	0	0	0	1	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0

**Adjacency
matrix for
PEDOT
monomer**

TAFFI Component Increment Theory (TCIT)

The fundamental idea

- Systematize component-definitions and model compound selection with rigorous graph-based typing.

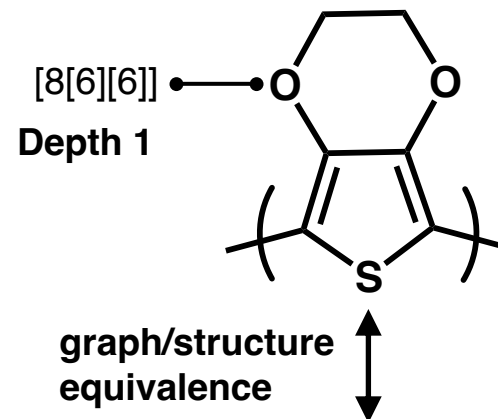
Zhao, Q.; Savoie, B. M.; "Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory". *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Zhao, Q.; Iovanac, N.; Savoie, B. M.; "Transferable Ring Corrections for Predicting Enthalpy of Formation of Cyclic Compounds" *J. Chem. Info. Model.* **2021**, 61, 5013-5027

Seo, B.; Lin, Z.-Y.; Zhao, Q.; Webb, M. A.; Savoie, B. M. Topology Automated Force-Field Interactions (TAFFI): A Framework for Developing Transferable Force Fields. *J. Chem. Inf. Model.* **2021**, 61 (10), 5013–5027. <https://doi.org/10.1021/acs.jcim.1c00491>.

**P2SAC
Publications**

TCIT is a component theory
(2-bond specific)



Topology Automated
Force Field Interactions



S	0	1	0	0	1	0	0	0	0	0	0	0	0
C	1	0	1	0	0	0	0	0	0	0	0	0	0
C	0	1	0	1	0	1	0	0	0	0	0	0	0
C	0	0	1	0	1	0	0	0	0	0	1	0	0
C	1	0	0	1	0	0	0	0	0	0	0	0	0
O	0	0	1	0	0	0	1	0	0	0	0	0	0
C	0	0	0	0	0	1	0	1	1	1	0	0	0
C	0	0	0	0	0	0	1	0	0	0	1	1	1
H	0	0	0	0	0	0	1	0	0	0	0	0	0
H	0	0	0	0	0	0	1	0	0	0	0	0	0
O	0	0	0	1	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0

**Adjacency
matrix for
PEDOT
monomer**

TAFFI Component Increment Theory (TCIT)

The fundamental idea

- Systematize component-definitions and model compound selection with rigorous graph-based typing.

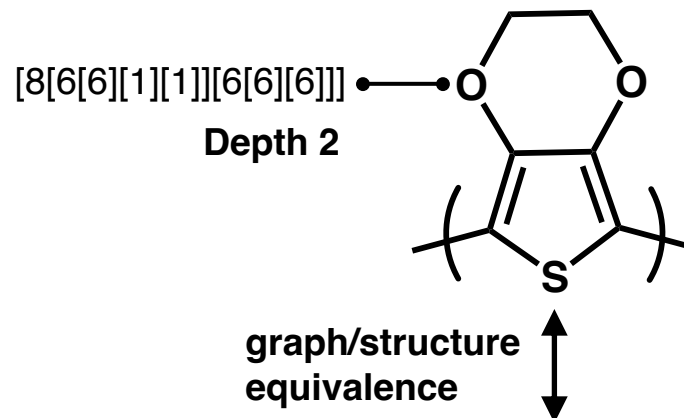
Zhao, Q.; Savoie, B. M.; “Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory”. *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Zhao, Q.; Iovanac, N.; Savoie, B. M.; “Transferable Ring Corrections for Predicting Enthalpy of Formation of Cyclic Compounds” *J. Chem. Info. Model.* **2021**, 61, 5013-5027

Seo, B.; Lin, Z.-Y.; Zhao, Q.; Webb, M. A.; Savoie, B. M. Topology Automated Force-Field Interactions (TAFFI): A Framework for Developing Transferable Force Fields. *J. Chem. Inf. Model.* **2021**, 61 (10), 5013–5027. <https://doi.org/10.1021/acs.jcim.1c00491>.

**P2SAC
Publications**

TCIT is a component theory
(2-bond specific)



Topology Automated
Force Field Interactions



S	0	1	0	0	1	0	0	0	0	0	0	0	0
C	1	0	1	0	0	0	0	0	0	0	0	0	0
C	0	1	0	1	0	1	0	0	0	0	0	0	0
C	0	0	1	0	1	0	0	0	0	0	1	0	0
C	1	0	0	1	0	0	0	0	0	0	0	0	0
O	0	0	1	0	0	0	1	0	0	0	0	0	0
C	0	0	0	0	0	1	0	1	1	1	0	0	0
C	0	0	0	0	0	0	1	0	0	0	1	1	1
H	0	0	0	0	0	0	1	0	0	0	0	0	0
H	0	0	0	0	0	0	1	0	0	0	0	0	0
O	0	0	0	1	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0

**Adjacency
matrix for
PEDOT
monomer**

TAFFI Component Increment Theory (TCIT)

The fundamental idea

- Systematize component-definitions and model compound selection with rigorous graph-based typing.

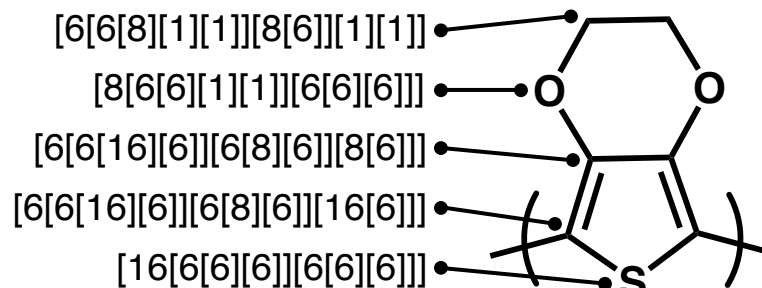
Zhao, Q.; Savoie, B. M.; "Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory". *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Zhao, Q.; Iovanac, N.; Savoie, B. M.; "Transferable Ring Corrections for Predicting Enthalpy of Formation of Cyclic Compounds" *J. Chem. Info. Model.* **2021**, 61, 5013-5027

Seo, B.; Lin, Z.-Y.; Zhao, Q.; Webb, M. A.; Savoie, B. M. Topology Automated Force-Field Interactions (TAFFI): A Framework for Developing Transferable Force Fields. *J. Chem. Inf. Model.* **2021**, 61 (10), 5013–5027. <https://doi.org/10.1021/acs.jcim.1c00491>.

**P2SAC
Publications**

TCIT is a component theory
(2-bond specific)



Topology Automated
Force Field Interactions



S	0	1	0	0	1	0	0	0	0	0	0	0	0
C	1	0	1	0	0	0	0	0	0	0	0	0	0
C	0	1	0	1	0	1	0	0	0	0	0	0	0
C	0	0	1	0	1	0	0	0	0	0	1	0	0
C	1	0	0	1	0	0	0	0	0	0	0	0	0
O	0	0	1	0	0	0	1	0	0	0	0	0	0
C	0	0	0	0	0	1	0	1	1	1	0	0	0
C	0	0	0	0	0	0	1	0	0	0	1	1	1
H	0	0	0	0	0	0	1	0	0	0	0	0	0
H	0	0	0	0	0	0	1	0	0	0	0	0	0
O	0	0	0	1	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0

**Adjacency
matrix for
PEDOT
monomer**

TAFFI Component Increment Theory (TCIT)

The fundamental idea

- Systematize component-definitions and model compound selection with rigorous graph-based typing.
- Two-bond specificity should improve both the accuracy and transferability of the resulting components.
- Parameterizing a component model **would not be feasible with only experimental data.**

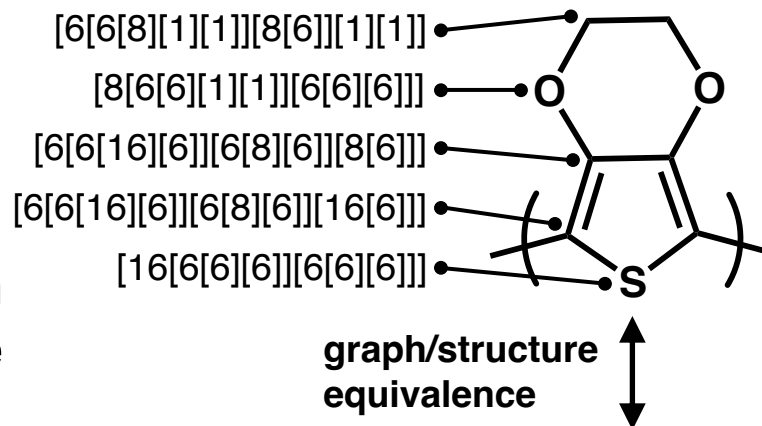
Zhao, Q.; Savoie, B. M.; "Enthalpy of Formation Prediction via a fully Self-Consistent Component Increment Theory". *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Zhao, Q.; Iovanac, N.; Savoie, B. M.; "Transferable Ring Corrections for Predicting Enthalpy of Formation of Cyclic Compounds" *J. Chem. Info. Model.* **2021**, 61, 5013-5027

Seo, B.; Lin, Z.-Y.; Zhao, Q.; Webb, M. A.; Savoie, B. M. Topology Automated Force-Field Interactions (TAFFI): A Framework for Developing Transferable Force Fields. *J. Chem. Inf. Model.* **2021**, 61 (10), 5013–5027. <https://doi.org/10.1021/acs.jcim.1c00491>.

**P2SAC
Publications**

TCIT is a component theory
(2-bond specific)



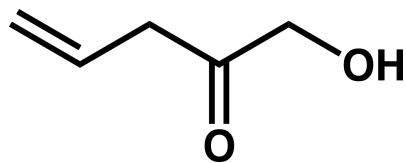
Topology Automated
Force Field Interactions



S	0	1	0	0	1	0	0	0	0	0	0	0	0
C	1	0	1	0	0	0	0	0	0	0	0	0	0
C	0	1	0	1	0	1	0	0	0	0	0	0	0
C	0	0	1	0	1	0	0	0	0	0	1	0	0
C	1	0	0	1	0	0	0	0	0	0	0	0	0
O	0	0	1	0	0	0	1	0	0	0	0	0	0
C	0	0	0	0	0	1	0	1	1	1	0	0	0
C	0	0	0	0	0	0	1	0	0	0	1	1	1
H	0	0	0	0	0	0	1	0	0	0	0	0	0
H	0	0	0	0	0	0	1	0	0	0	0	0	0
O	0	0	0	1	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0
H	0	0	0	0	0	0	0	1	0	0	0	0	0

**Adjacency
matrix for
PEDOT
monomer**

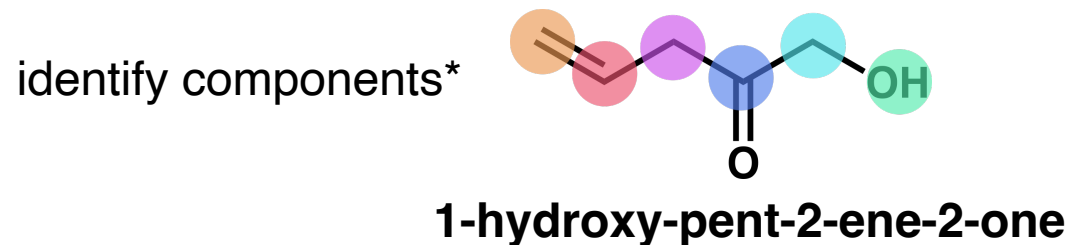
Graphical Decomposition of Model Compounds



1-hydroxy-pent-2-ene-2-one

How will we select molecules for parameterizing TCIT components?

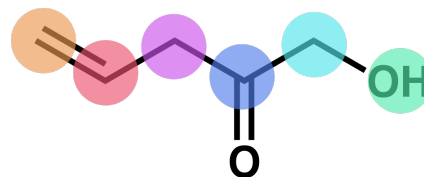
Graphical Decomposition of Model Compounds



**How will we select
molecules for
parameterizing TCIT
components?**

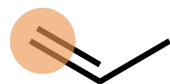
Graphical Decomposition of Model Compounds

identify components*



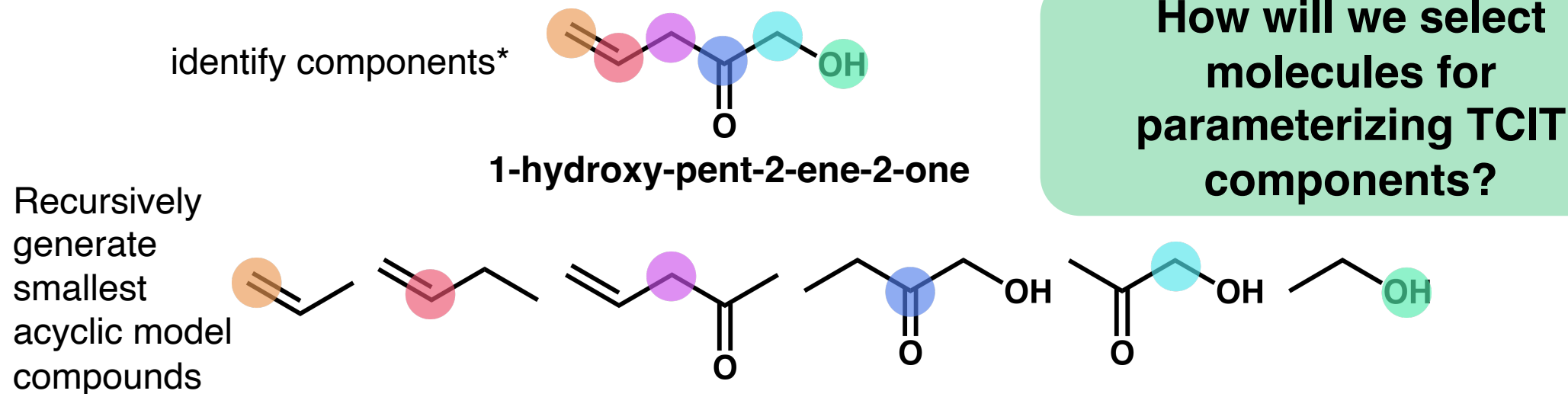
1-hydroxy-pent-2-ene-2-one

Recursively
generate
smallest
acyclic model
compounds

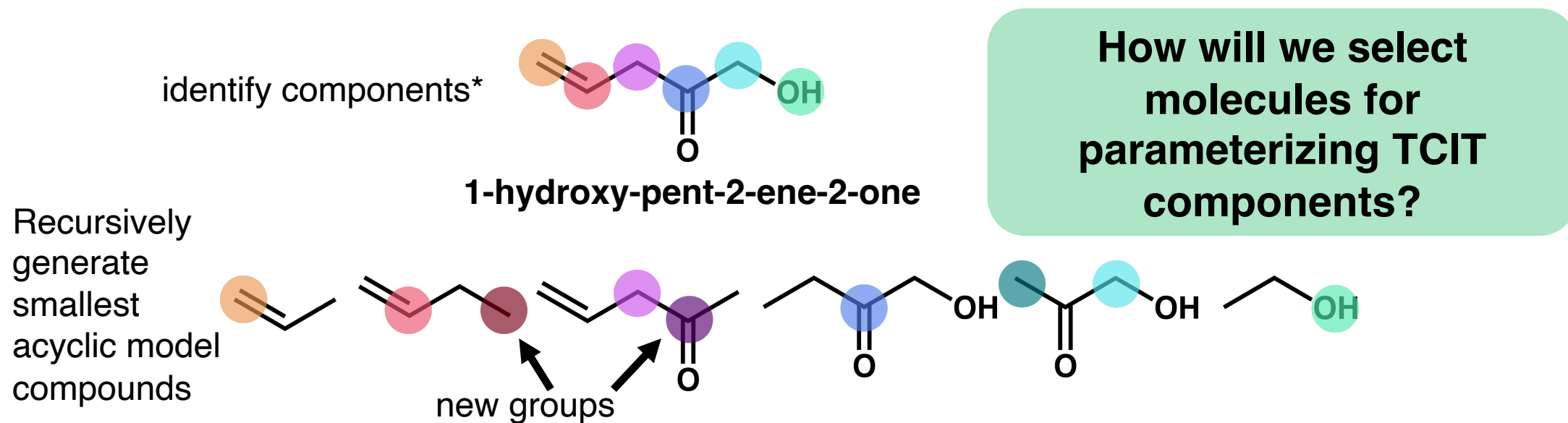


**How will we select
molecules for
parameterizing TCIT
components?**

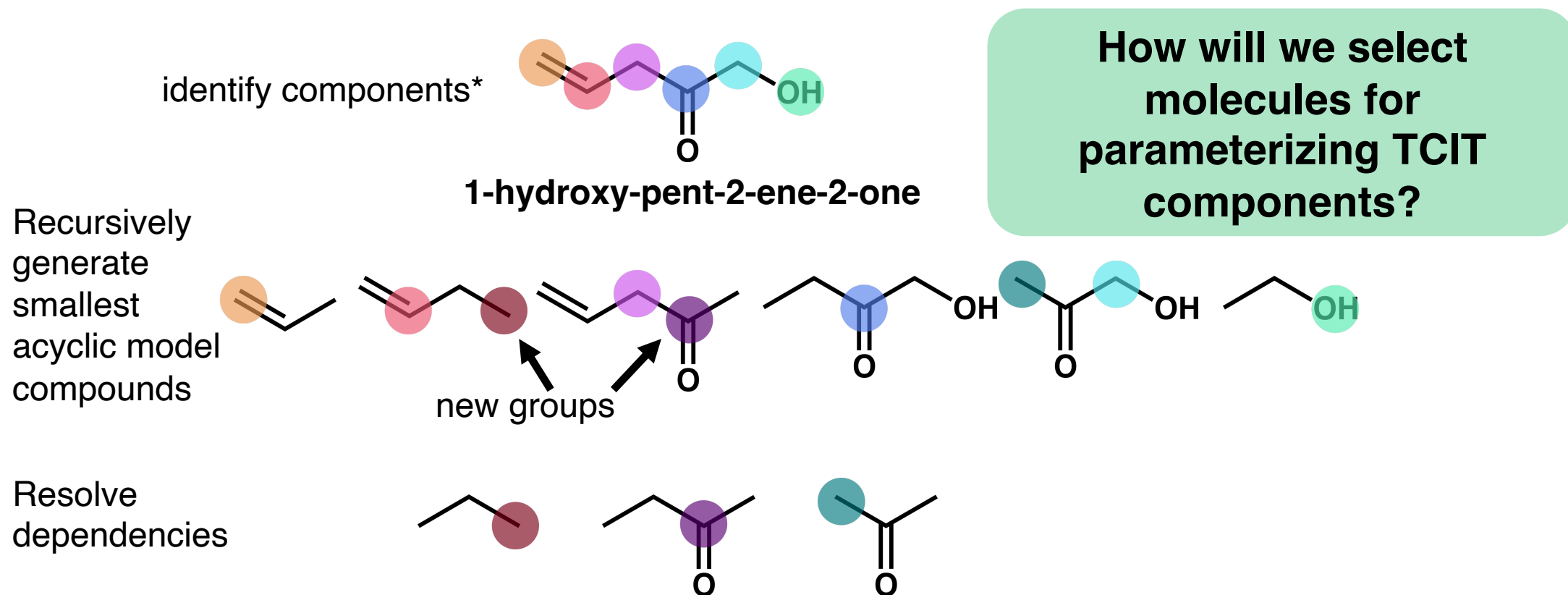
Graphical Decomposition of Model Compounds



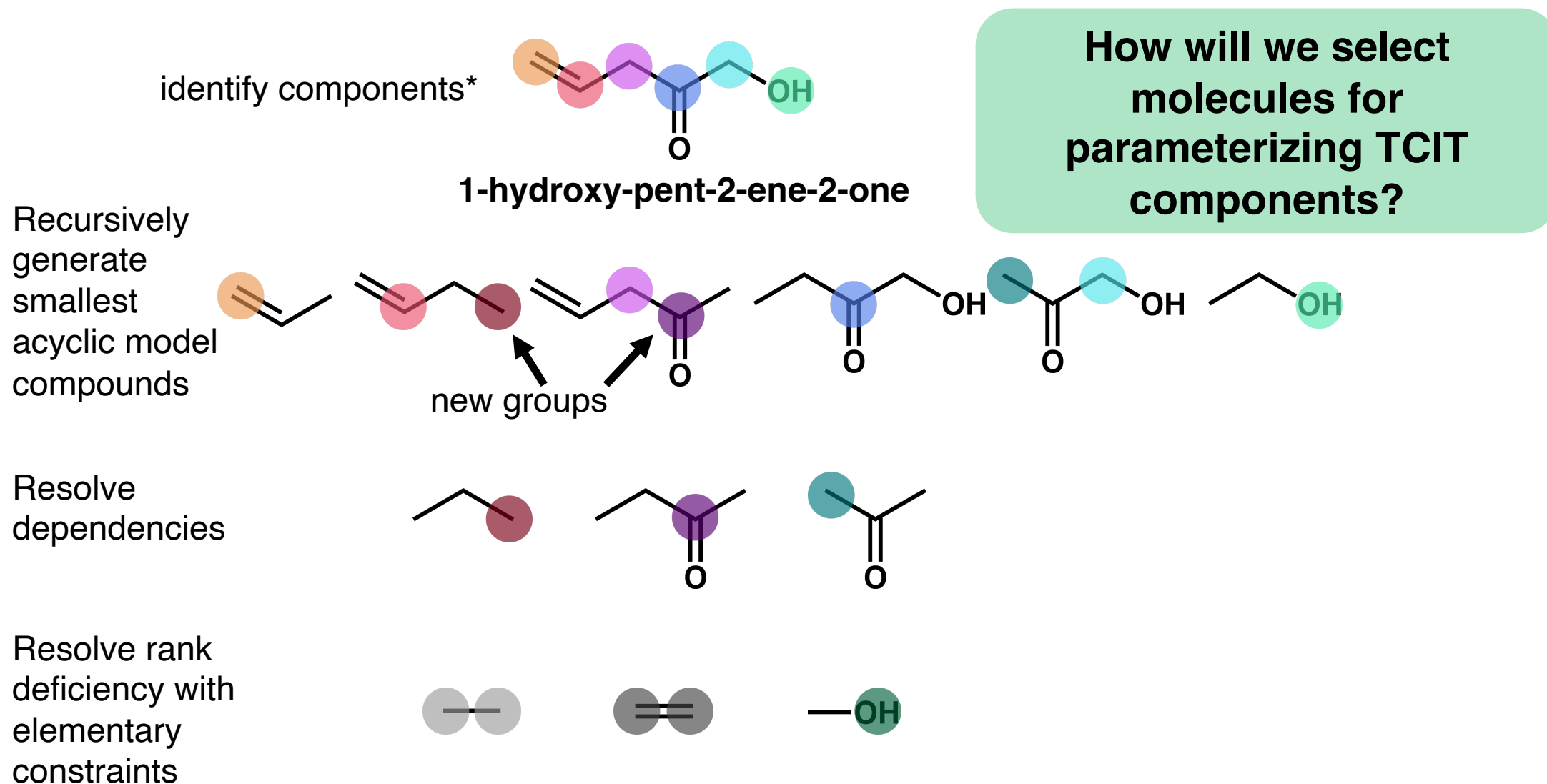
Graphical Decomposition of Model Compounds



Graphical Decomposition of Model Compounds

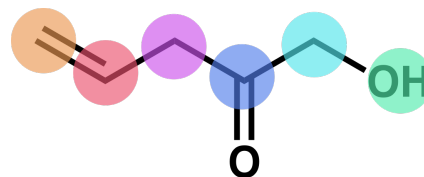


Graphical Decomposition of Model Compounds



Graphical Decomposition of Model Compounds

Prediction target:



1-hydroxy-pent-2-ene-2-one

$$\Delta H_{f,G4} = -259.9 \text{ kJ/mol}$$

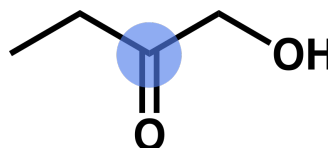
$$\Delta H_{f,TCIT} = -259.3 \text{ kJ/mol}$$

no experimental data

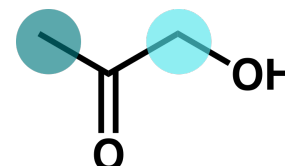
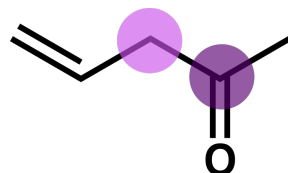
Topologically
sort
dependency
graph

(Automatically
handled by
TCIT software)

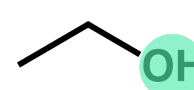
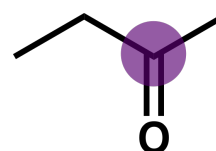
Gen 4:



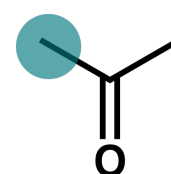
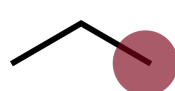
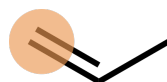
Gen 3:



Gen 2:



Gen 1:



Gen 0:



Model compounds
are small enough to
perform the highest
quality quantum
chemistry
calculations (G4
throughout)

Graphical Decomposition of Model Compounds

Have we solved the specificity problem?

All components are unique out to a graph depth of two,
no exceptions.

Have we solved the provenance problem?

All ΔH_f data is calculated at the G4 composite level,
no exceptions.

Have we solved the extensibility problem?

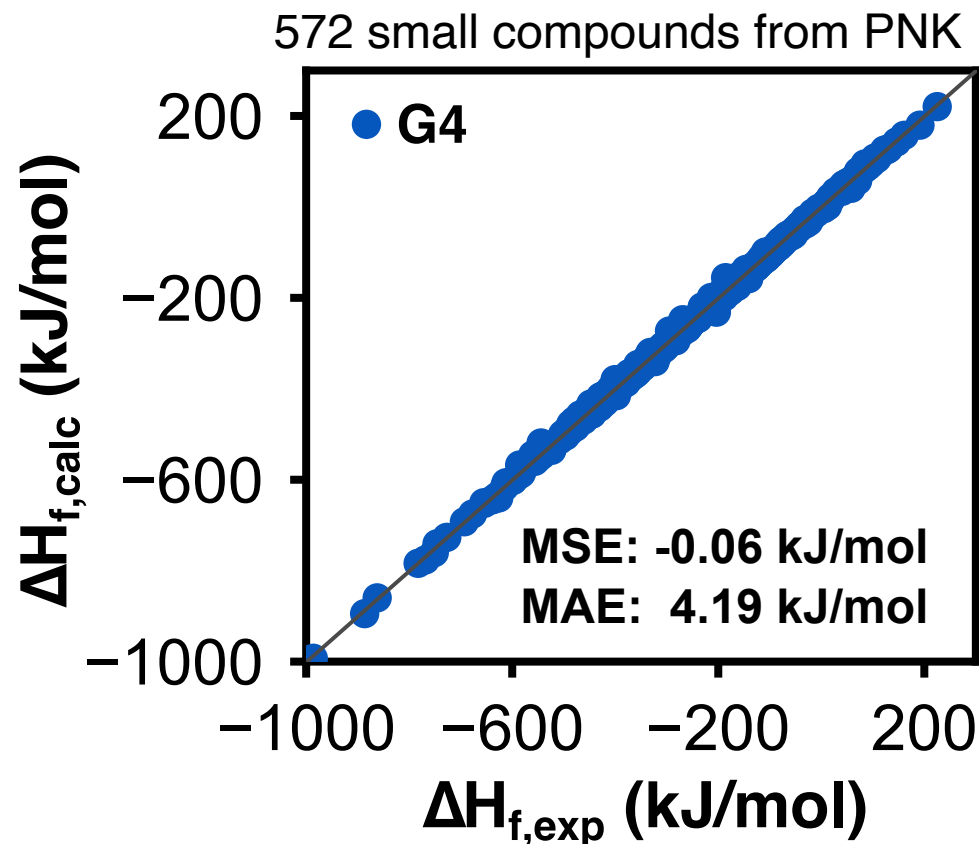
Model compounds exist for all conceivable components,
no exceptions.

Benchmarking $\Delta H_{f,gas}$ Predictions Against the PNK Dataset

- Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK¹

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2nd ed. 1986

- PNK is a core dataset for fitting Benson groups
- ~600 PNK compounds are small enough for G4 calculations and comparison with experiment.



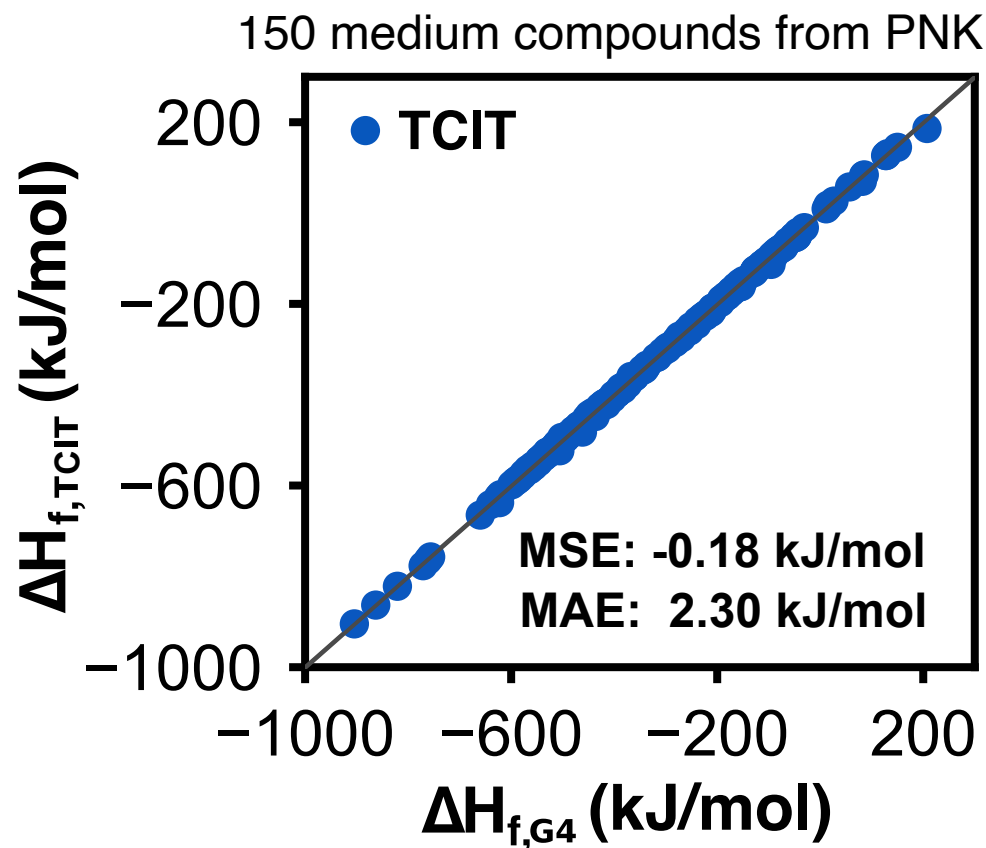
Zhao, Q.; Savoie, B. M.; Enthalpy of Formation Prediction via a Fully Self-Consistent Component Increment Theory. *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Benchmarking $\Delta H_{f,gas}$ Predictions Against the PNK Dataset

- Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK¹

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2nd ed. 1986

- PNK is a core dataset for fitting Benson groups
- ~600 PNK compounds are small enough for G4 calculations and comparison with experiment.
- ~150 PNK compounds are large enough for direct G4 calculation and comparison with TCIT.



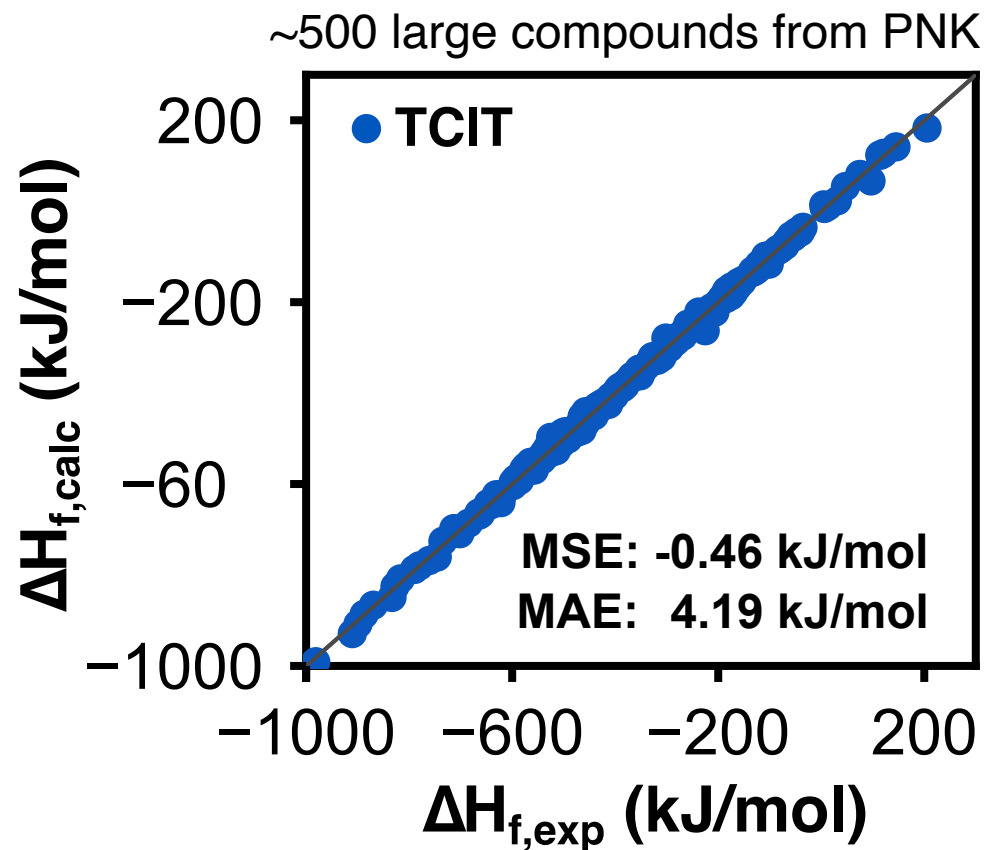
Zhao, Q.; Savoie, B. M.; Enthalpy of Formation Prediction via a Fully Self-Consistent Component Increment Theory. *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Benchmarking $\Delta H_{f,gas}$ Predictions Against the PNK Dataset

- Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK¹

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2nd ed. 1986

- PNK is a core dataset for fitting Benson groups
- ~600 PNK compounds are small enough for G4 calculations and comparison with experiment.
- ~150 PNK compounds are large enough for direct G4 calculation and comparison with TCIT.
- ~500 PNK compounds are large enough to evaluate the predictive accuracy of the increment theories.



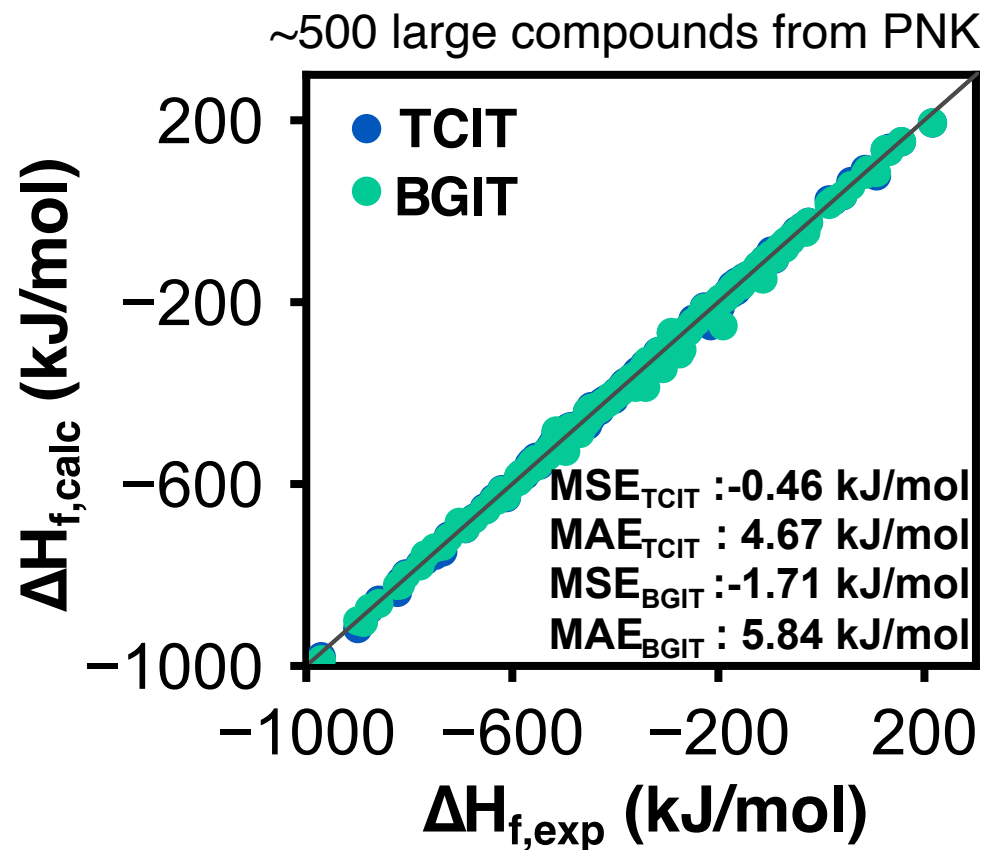
Zhao, Q.; Savoie, B. M.; Enthalpy of Formation Prediction via a Fully Self-Consistent Component Increment Theory. *J. Chem. Info. Model.* **2020**, 60, 2199-2207

Benchmarking $\Delta H_{f,gas}$ Predictions Against the PNK Dataset

- Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK¹

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2nd ed. 1986

- PNK is a core dataset for fitting Benson groups
- ~600 PNK compounds are small enough for G4 calculations and comparison with experiment.
- ~150 PNK compounds are large enough for direct G4 calculation and comparison with TCIT.
- ~500 PNK compounds are large enough to evaluate the predictive accuracy of the increment theories.



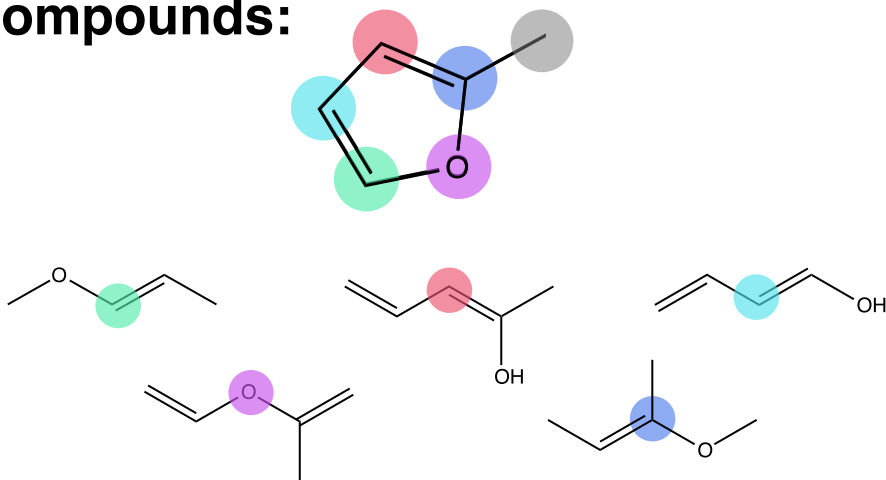
Zhao, Q.; Savoie, B. M.; Enthalpy of Formation Prediction via a Fully Self-Consistent Component Increment Theory. *J. Chem. Info. Model.* **2020**, 60, 2199-2207

TCIT shows comparable performance to BGIT/CHETAH but is derived exclusively from extensible G4 data.

Extension to Ring-Containing Molecules

- Ring-containing molecules have additional strain and/or conjugation corrections that exacerbate the extensibility issues of Benson Theory.
- In TCIT we are addressing this through chemically specific ring corrections that account for differences in substitution pattern and topology:

1. Decompose ring into acyclic model compounds:



2. Add ring correction (RC) to final prediction:

$$RC = H_f(\text{ring}) - H_f(\text{red}) - H_f(\text{cyan}) - H_f(\text{green}) - H_f(\text{purple}) - H_f(\text{blue}) - H_f(\text{grey})$$

Technical Developments

RC Model Compounds

RC₀ RC₁ RC₂

○ : Depth 0 ● : Depth 1 ● : Depth 2

Method 1: Use RC₁ based model parameterized to G4 data.

Method 2: Use graph-NN to predict RC₀-RC₂

Benchmarking Ring-Correction Performance

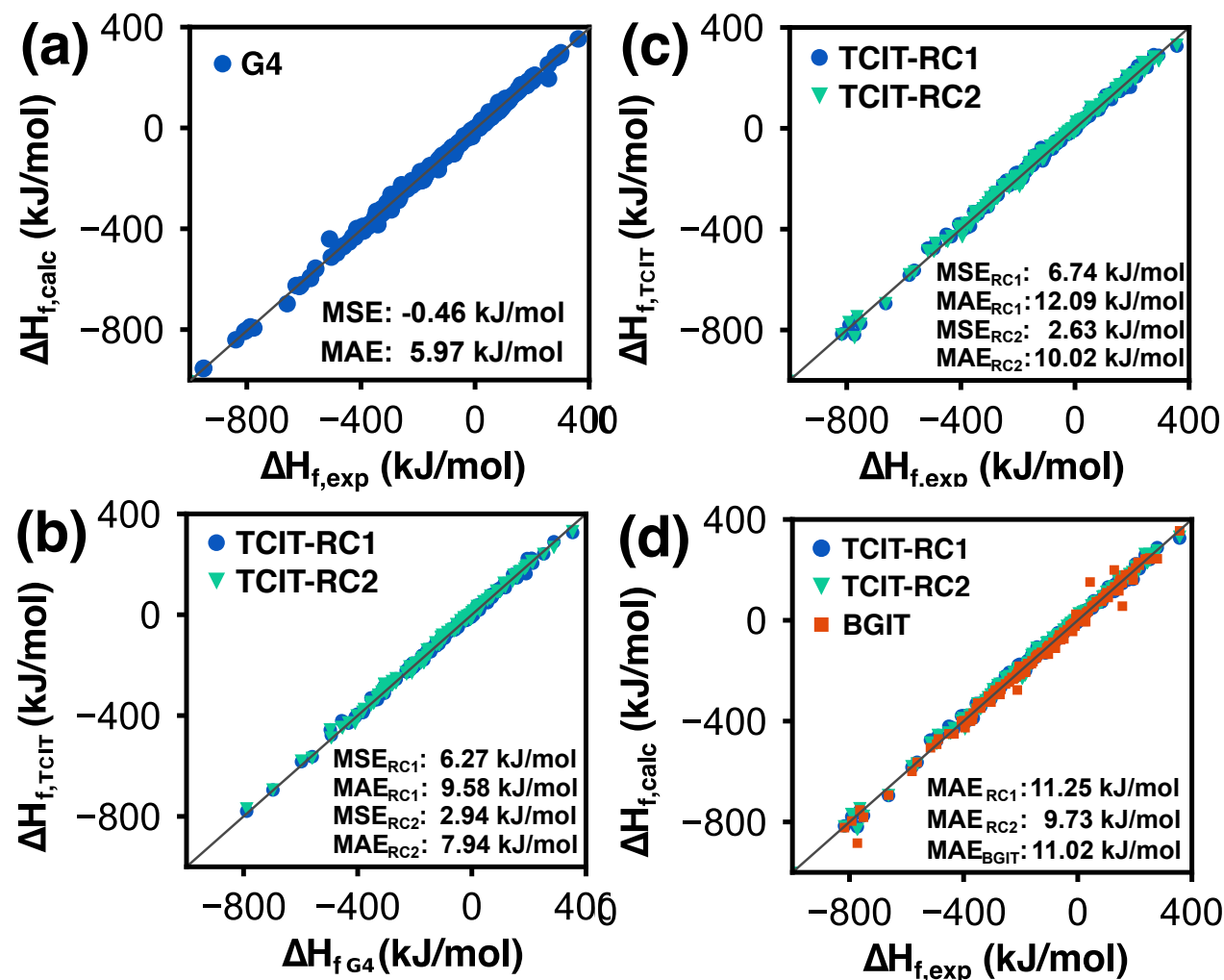
(a) G4 errors are marginally larger for ring-containing compounds but still very accurate

(b) The neural-network based ring-correction exhibits excellent reproduction of the G4 predictions (MSE: ~ 3 kJ/mol; MAE: ~ 8 kJ/mol).

(c) TCIT is completely transferable to new testing compounds that are experimentally characterized. Errors are consistent with G4 comparison

(d) The TCIT-R2 model outperforms BGIT on the large molecule benchmark while being extensible. Significantly, these compounds are within BGIT's training data.

~ 120 ring-containing compounds from PNK (excluding training)



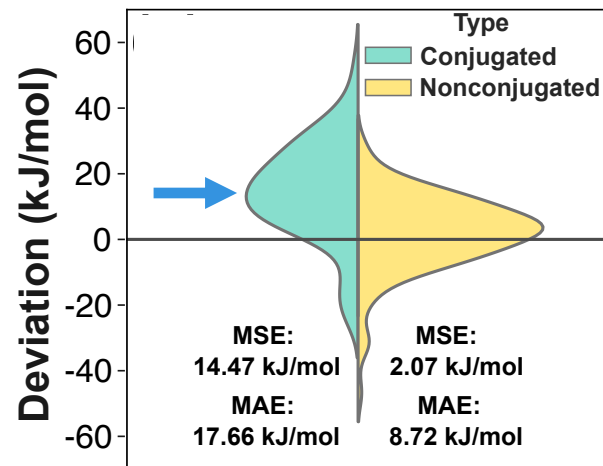
BGIT cannot make predictions for $\sim 2\%$ of PNK compounds

Benchmarking Ring-Correction Performance

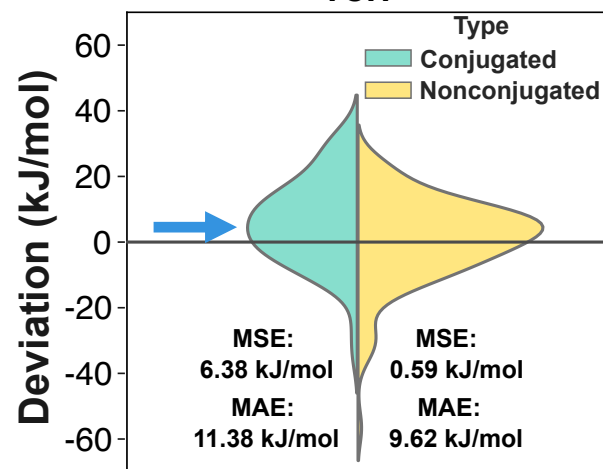
Breaking down distinct subsets:

- Conjugated systems are challenging to accurately predict with a local ring correction.
- BGIT has excellent performance on benzene rings due to the prevalence of experimental data, but poor performance on novel rings.
- The ML ring-correction shows the strongest overall performance. This strategy could also be used to generically correct for long-range conjugation effects.

~120 ring-containing compounds from PNK (excluding training)



TCIT



TCIT-ML

BGIT cannot make predictions for ~2% of PNK compounds

TCIT Extension to Other Properties and Phases

Condensed Phases: The condensed-phase and gas-phase standard enthalpies of formation differ by the heats of sublimation and vaporization^[1]:

$$\Delta_f H_{(s)}^\circ = \Delta_f H_{(g)}^\circ - \Delta_{\text{sub}} H^\circ$$

$$\Delta_f H_{(\ell)}^\circ = \Delta_f H_{(g)}^\circ - \Delta_{\text{vap}} H^\circ$$

We have implemented group contribution models for heat of vaporization^[2] and sublimation^[3], respectively. The group assignments and group values associated with these models have been automated within the context of TCIT.

Standard Molar Entropy (S°) and heat capacity (C_v): The molar entropies and constant volume heat capacities are accessible from quantum chemistry using the harmonic oscillator approximation for the molecular partition function and corrections based on the number of rotatable bonds (N_{rot}) and molecular symmetry:

$$S^\circ = \langle S_{\text{harm}}^\circ \rangle + RN_{\text{rot}} + R \log \sigma \quad C_v = \langle C_{v,\text{harm}} \rangle + \alpha N_{\text{rot}} + \beta$$

[1] Murray, J.S., Brinck, T. and Politzer, P., **1996**. *Chemical physics*, 204, 289-299.

[2] Pankow, J.F. and Asher, W.E., **2008**. *Atmospheric Chemistry and Physics*.

[3] Bagheri, M.; Gandomi, A. H.; Golbraikh, A. **2012**, *Thermochim. Acta*, 543, 96–106

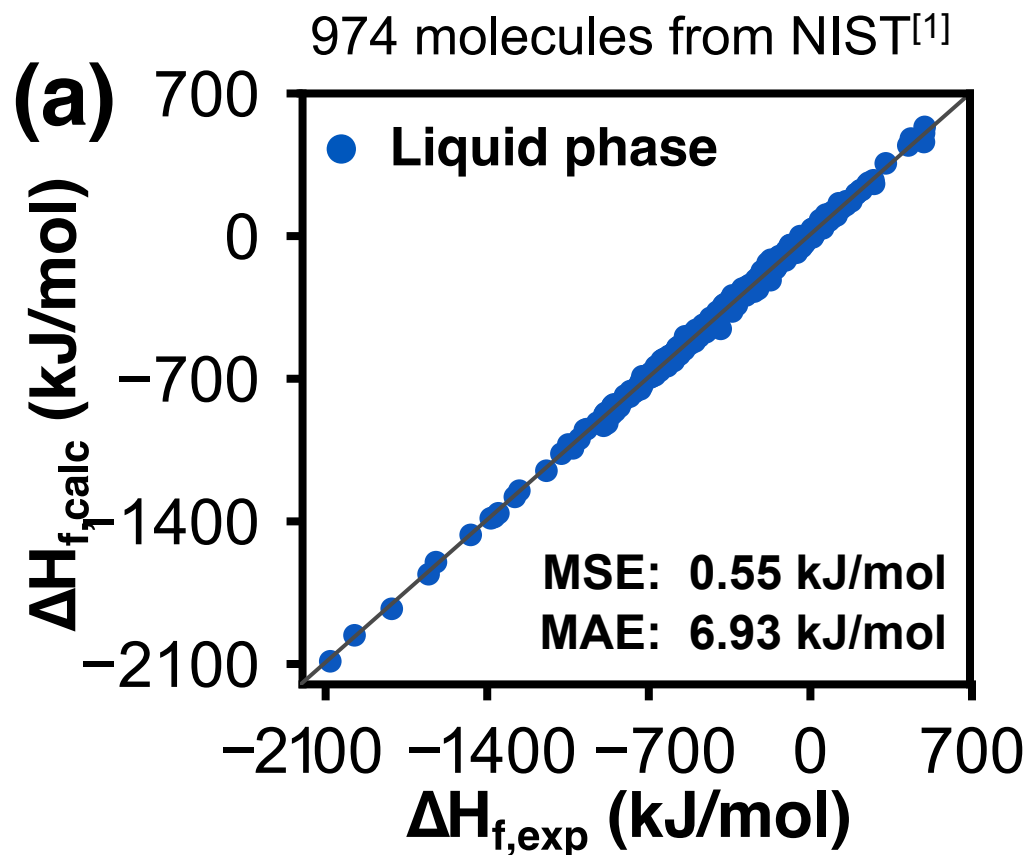
<.> Indicates conformational averaging

R: ideal gas constant

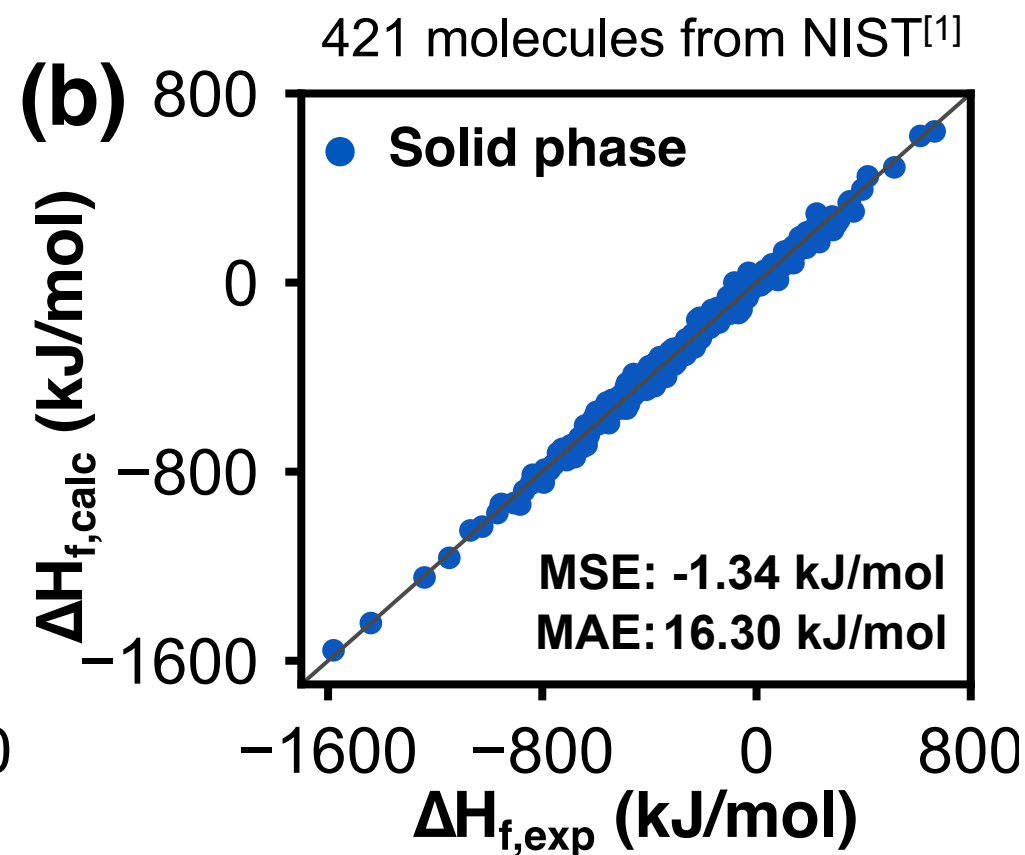
σ : symmetry number

α , β : regressed constants

Benchmarking Condensed Phase ΔH_f Predictions



- Testing set includes both linear and cyclic compounds with number of heavy atoms varying from 1 to 30.



- Low MSE indicates no systematic bias, larger absolute errors result from the quality of the ΔH_{vap} and ΔH_{sub} models.

Benchmarking TCIT S° and C_v Predictions

(a) G4/TCIT S° comparison for 314 medium sized molecules.

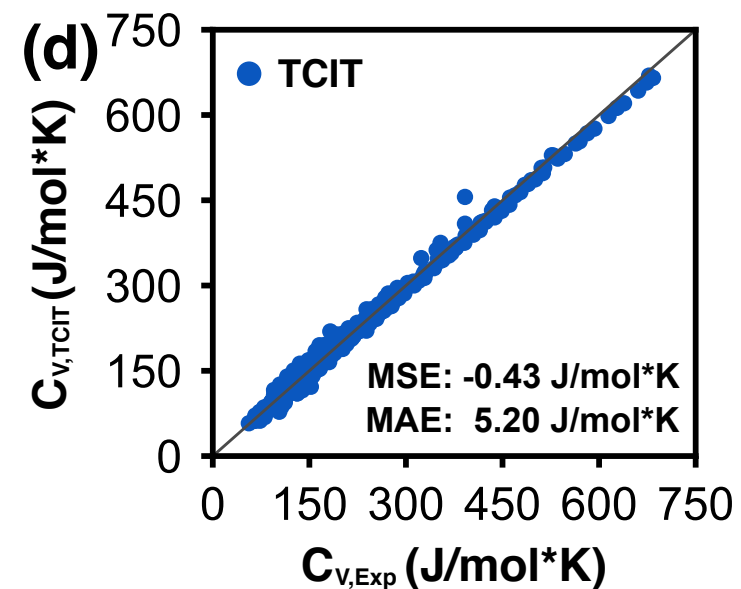
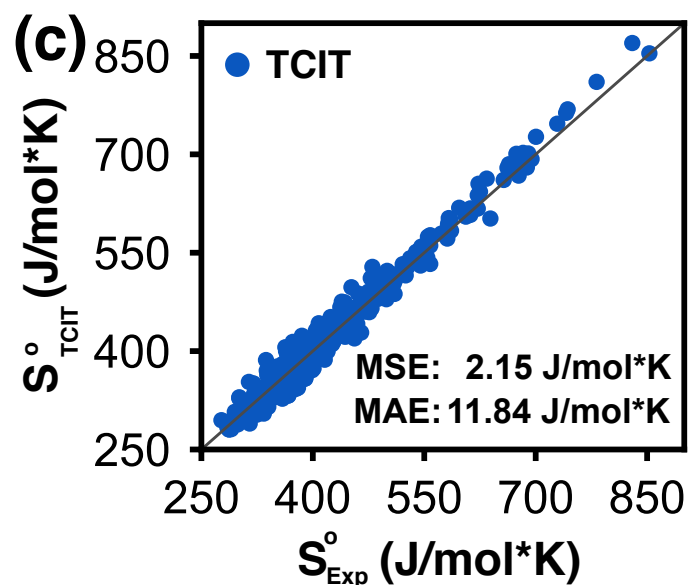
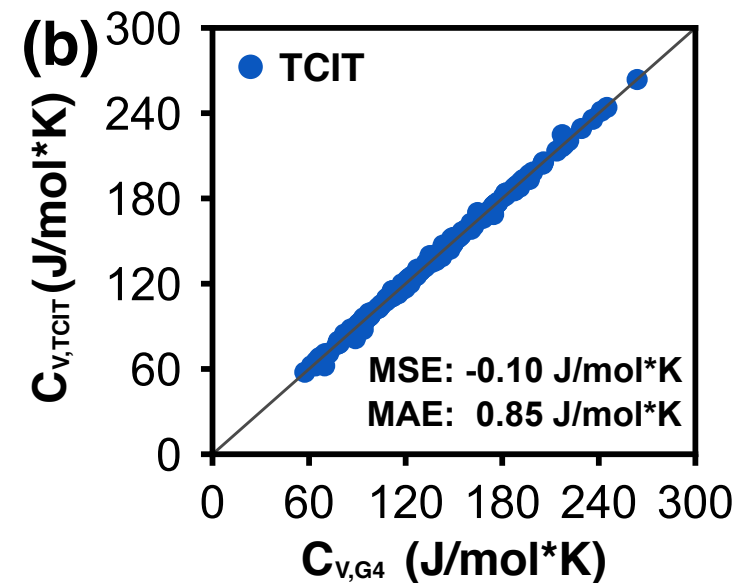
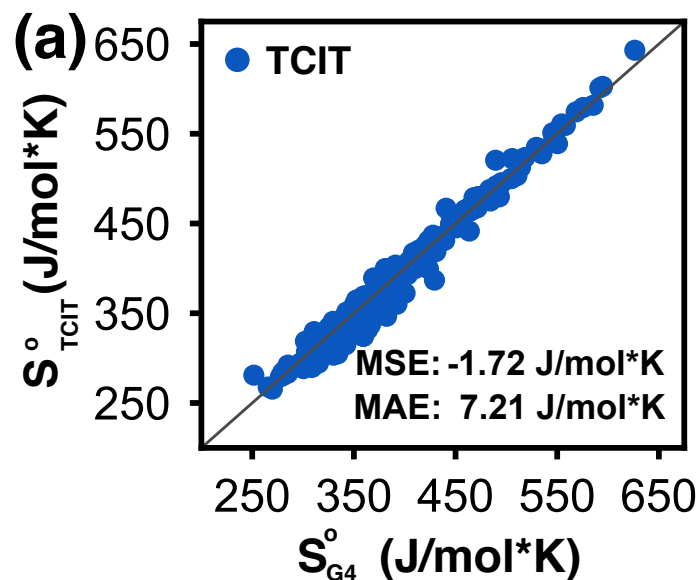
(b) G4/TCIT S° comparison for 314 medium sized molecules.

(c) TCIT S° comparison for 439 large molecules from NIST^[1]

(d) TCIT heat capacity comparison for 904 large molecules from NIST^[1]

- The TCIT errors are consistent with error propagation of G4:exp and TCIT:G4 errors.

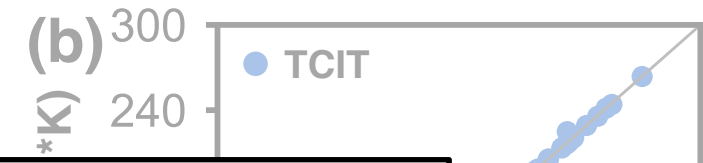
TCIT now supports S° and C_v predictions with accuracies comparable to G4 model chemistry.



[1] Linstrom, P.J. and Mallard, W.G., 2001. *Journal of Chemical & Engineering Data*, 46(5), pp.1059-1063.

Benchmarking TCIT S° and C_v Predictions

(a) G4/TCIT S° comparison for 314 medium sized molecules.



(b) G4/TCIT S° comparison for 314 medium sized molecules.

(c) TCIT S° comparison for 314 medium sized molecules.

(d) TCIT S° comparison for 904 large molecules.

• The TCIT error for S° is comparable to G4 model chemistry.

P2SAC funding over the past 3 years has allowed us to establish a component increment theory that addresses the major gaps in Benson group theory.

All of the promised functionality from that original proposal has now been delivered.

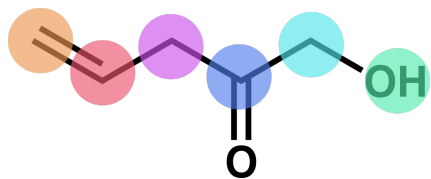
-0.10 J/mol*K
0.85 J/mol*K
180 240 300
(J/mol*K)

-0.43 J/mol*K
5.20 J/mol*K
450 600 750
(J/mol*K)

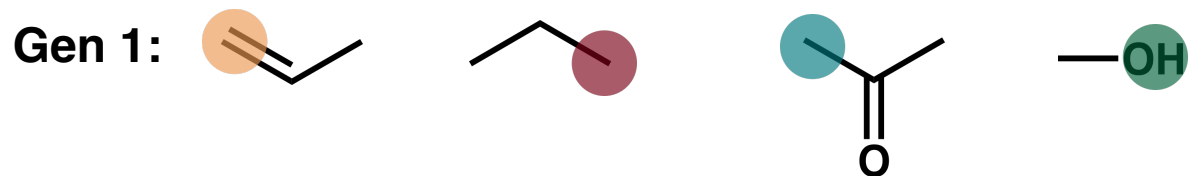
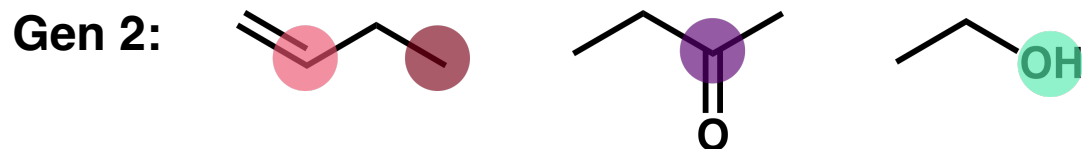
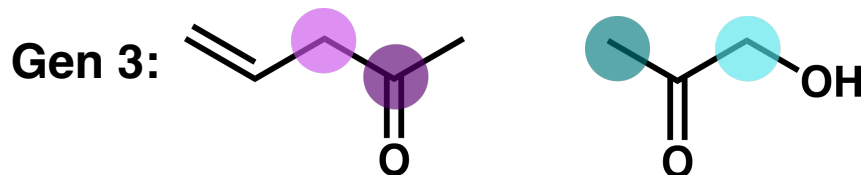
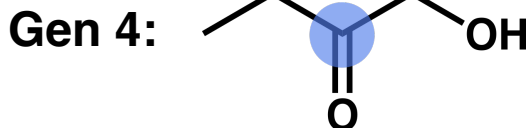
TCIT predictions with accuracies comparable to G4 model chemistry.

How Many Components are Possible?

Prediction target:



1-hydroxy-pent-2-ene-2-one



We database all model compounds and components for reuse.

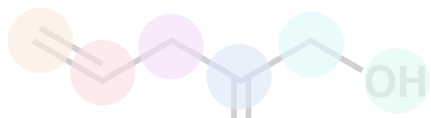
Over the past two years, we have parameterized new components in response to distinct project needs (many from P2SAC Pharma Members)

Current Database:

- ~35k distinct components for ΔH_f relevant to organic chemistry
- ~35k distinct G4 calculations on organic molecules.
- ~450 distinct ring corrections

How Many Components are Possible?

Prediction target:



We database all model compounds

How many components are required to predict the ΔH_f of **all** (physically relevant) organic molecules?

Gen 3:



response to distinct project needs
(many from P2SAC Pharma Members)

How many P2SAC funding periods would it take to make a “complete”/gapless component theory?

Current Database:

- ~450 distinct ring corrections

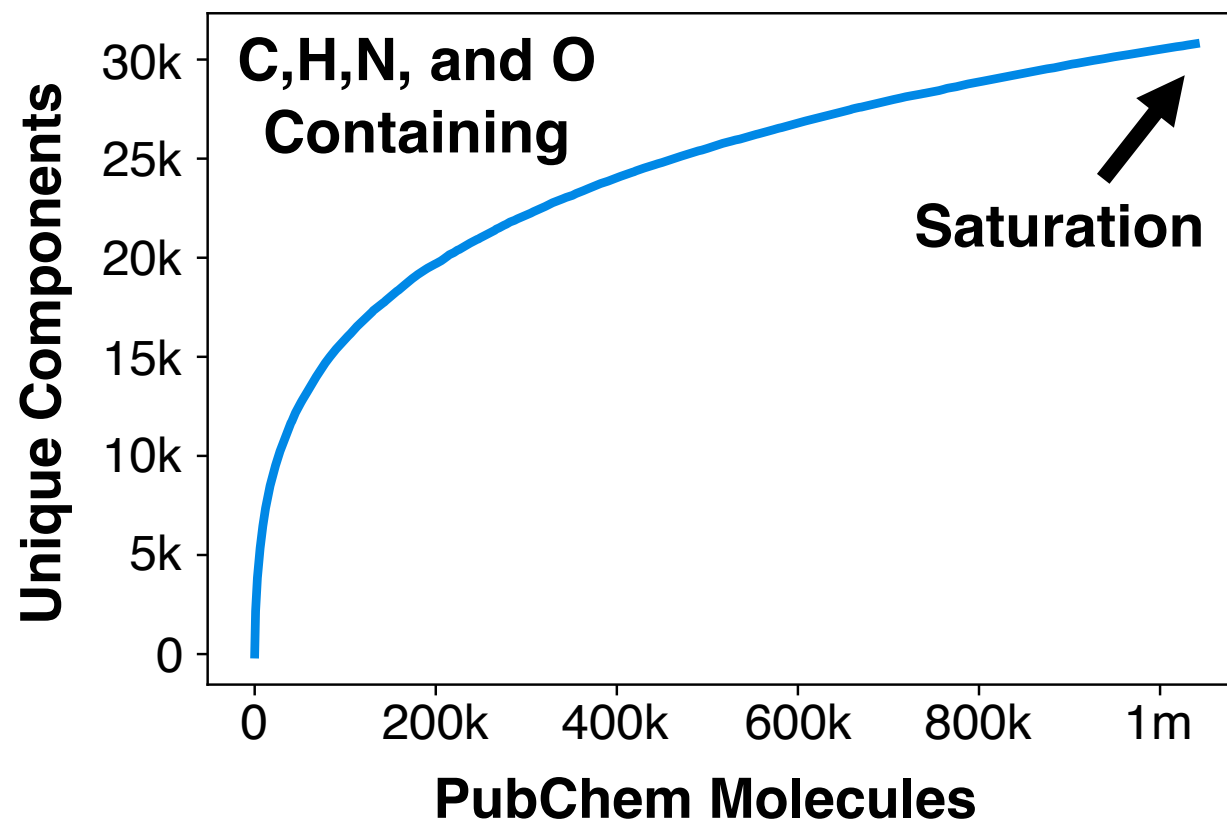
Gen 0:



Treating PubChem as a Model of Organic Chemical Space

PubChem is a repository of chemical properties that contains many millions of organic species ranging from small molecules to oligonucleotides.

We recently started mining PubChem's H,C,N, and O containing molecules for distinct components and the model compounds necessary to predict ΔH_f

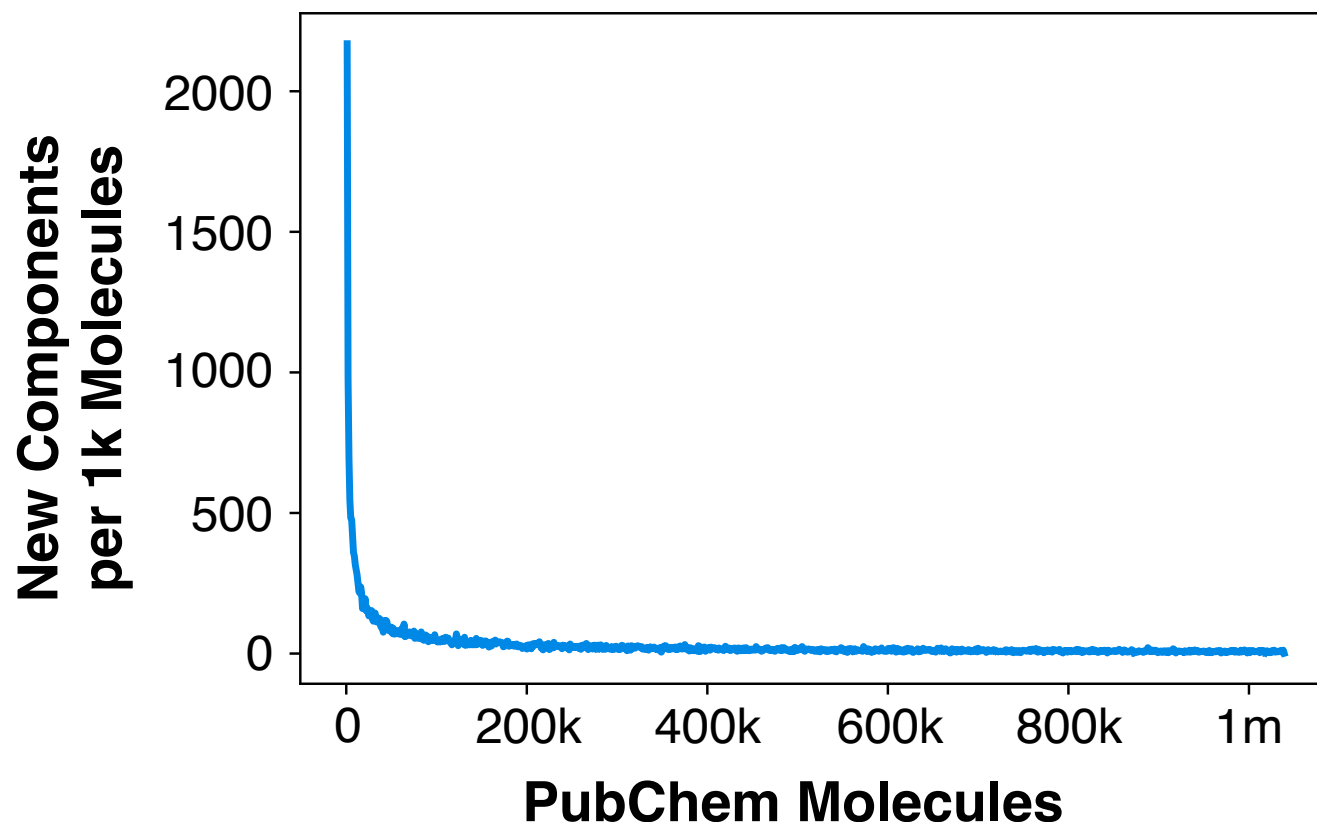


Treating PubChem as a Model of Organic Chemical Space

PubChem is a repository of chemical properties that contains many millions of organic species ranging from small molecules to oligonucleotides.

We recently started mining PubChem's H,C,N, and O containing molecules for distinct components and the model compounds necessary to predict ΔH_f

The derivative plot shows that TCIT initially generates ~2 new components per molecule, but by the end of the sampling ~100 molecules need to be sampled to find a new component.



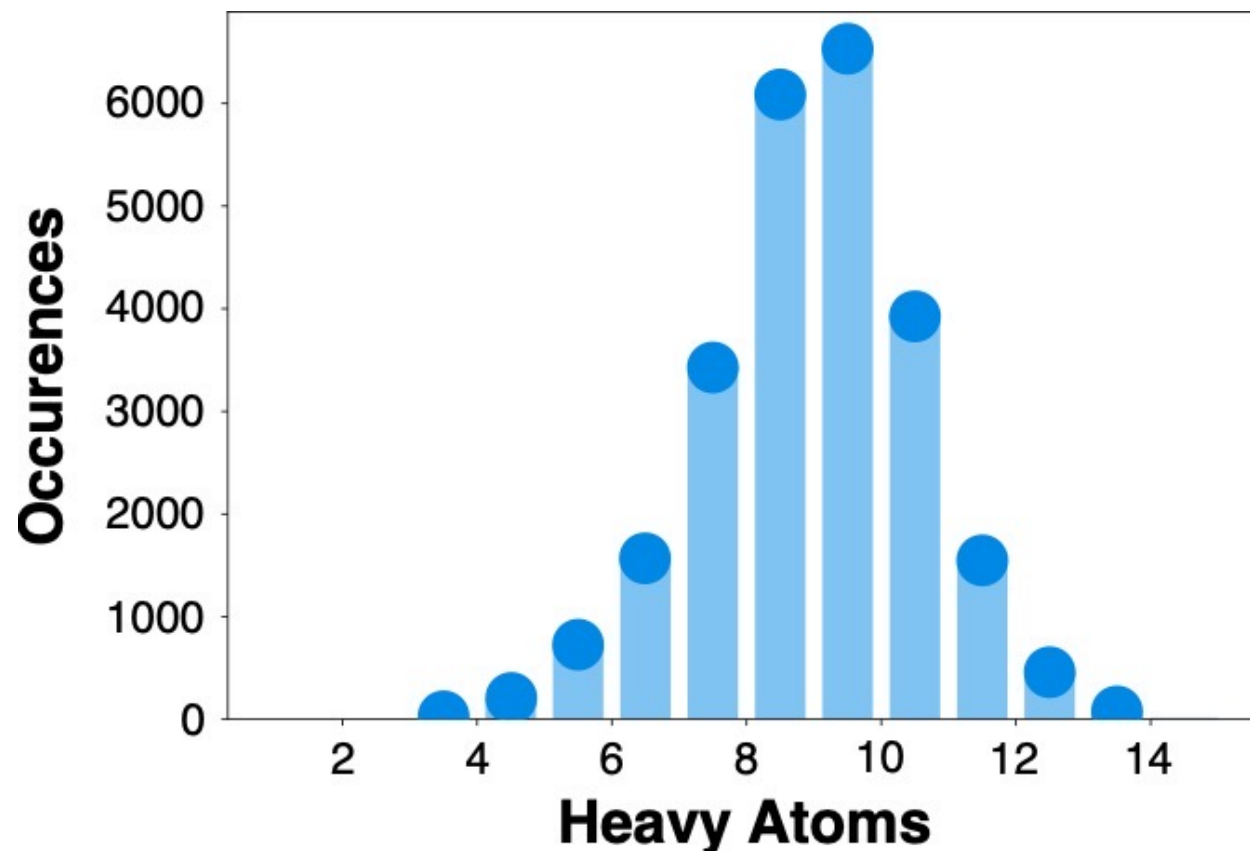
Treating PubChem as a Model of Organic Chemical Space

PubChem is a repository of chemical properties that contains many millions of organic species ranging from small molecules to oligonucleotides.

We recently started mining PubChem's H,C,N, and O containing molecules for distinct components and the model compounds necessary to predict ΔH_f

The derivative plot shows that TCIT initially generates ~2 new components per molecule, but by the end of the sampling ~100 molecules need to be sampled to find a new component.

New model compounds



Treating PubChem as a Model of Organic Chemical Space

In the past six months we have generated all training data necessary to make predictions on all N, H, O, and C-containing molecules in pubchem. **This is the largest repository of G4 calculations on large molecules in the world.**

It is foreseeable that we could complete all B, F, Cl, S, and P containing structures over the next few years.

“Known Unknowns” and “Unknown Unknowns”

$A \rightarrow B$

- To safely plan a known reaction, we need access to solid thermodynamic data (e.g., ΔH_f , S° , C_v) to understand and classify risks.
- This is a “known unknown” in that we know the reaction, $A \rightarrow B$, but we need values for a few unknown variables.

$A \rightarrow ? \rightarrow B$; $A \rightarrow B + ?$; $A \rightarrow ?$

- $A \rightarrow ? \rightarrow B$, means that we know the net reaction, but there may be a consequential (e.g., potentially reactive) intermediate. Even if we have accurate thermodynamic data on A/B, neglecting the intermediate could be disastrous.
- The $A \rightarrow B + ?$ (unknown side-reaction) and $A \rightarrow ?$ (unknown main product), problems have similar “unknown unknown” characteristics.

The Reaction Prediction Problem

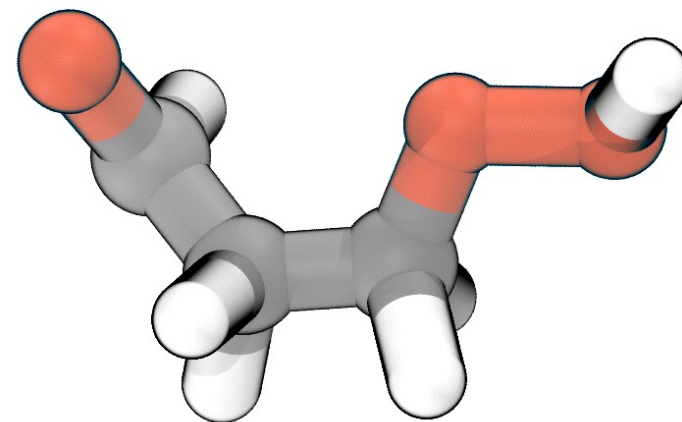
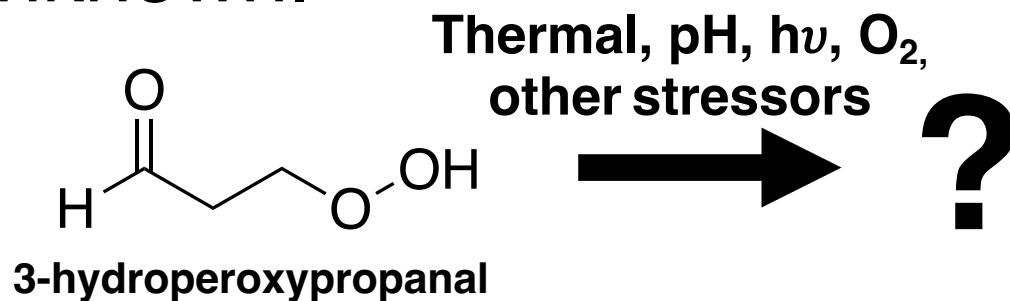
A → B : When we know the reactants and products, mature quantum chemistry tools exist to characterize transition states and establish pathways

A → ? : For degradation reactions, plausible reactions are often unknown.

The Reaction Prediction Problem

A → B : When we know the reactants and products, mature quantum chemistry tools exist to characterize transition states and establish pathways

A → ? : For degradation reactions, plausible reactions are often unknown.

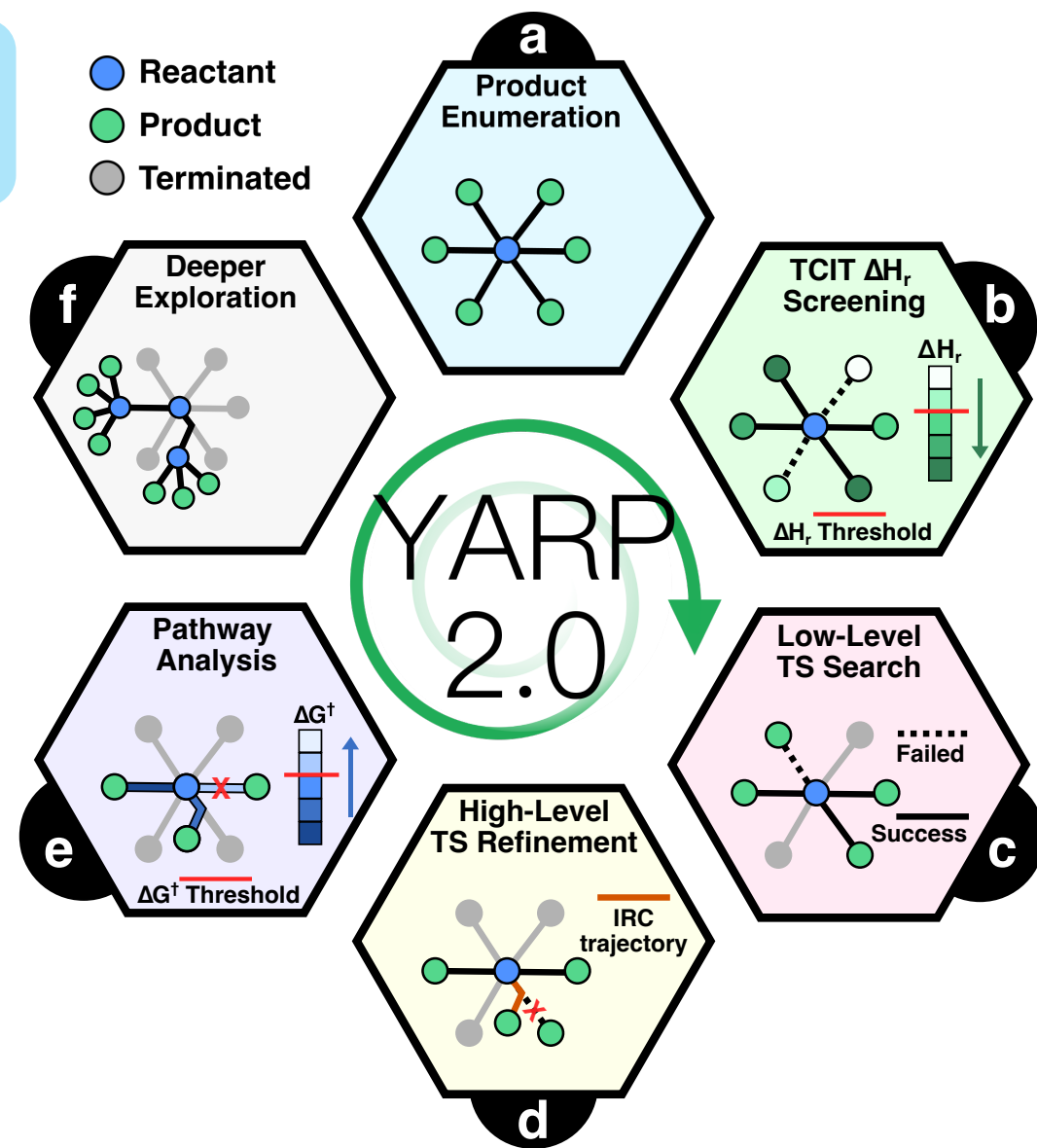


Yet Another Reaction Program (YARP)

Idea: Turn the $A \rightarrow ?$ problem into tractable (and parallelizable) $A \rightarrow B$ problems.

Observations:

- Product enumeration is easier than transition state enumeration.
- Transition state algorithms for $A \rightarrow B$ problems are mature. Let the TS algorithm identify physical reactions.
- Recent developments in semi-empirical models and ML create opportunities.
- If you are fast enough, you can brute force the $A \rightarrow ?$ problem.



Converting Reactions into a Machine-Readable Grammar

Bond-Electron Matrix Formalism: matrix representation of molecules with bond order indicated in off-diagonal elements and lone electrons along the diagonal.

Ugi, I. et al. "New Applications of Computers in Chemistry." *Angew. Chem.* **1979**, 18 (2), 111–123.



$$\begin{matrix} \text{O} \\ \text{C} \\ \text{C} \\ \text{H} \\ \text{H} \\ \text{O} \\ \text{H} \\ \text{H} \end{matrix} \begin{bmatrix} 4 & 1 & 0 & 0 & 0 & 0 & 1 & 0 \\ 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 2 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 2 & 0 & 0 & 4 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

$$\begin{matrix} \text{O} \\ \text{C} \\ \text{C} \\ \text{H} \\ \text{H} \\ \text{O} \\ \text{H} \\ \text{H} \end{matrix} \begin{bmatrix} 4 & 0 & 1 & 0 & 0 & 0 & 1 & 0 \\ 0 & 0 & 1 & 1 & 1 & 1 & 0 & 0 \\ 1 & 1 & 0 & 0 & 0 & 1 & 0 & 1 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 1 & 0 & 0 & 4 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

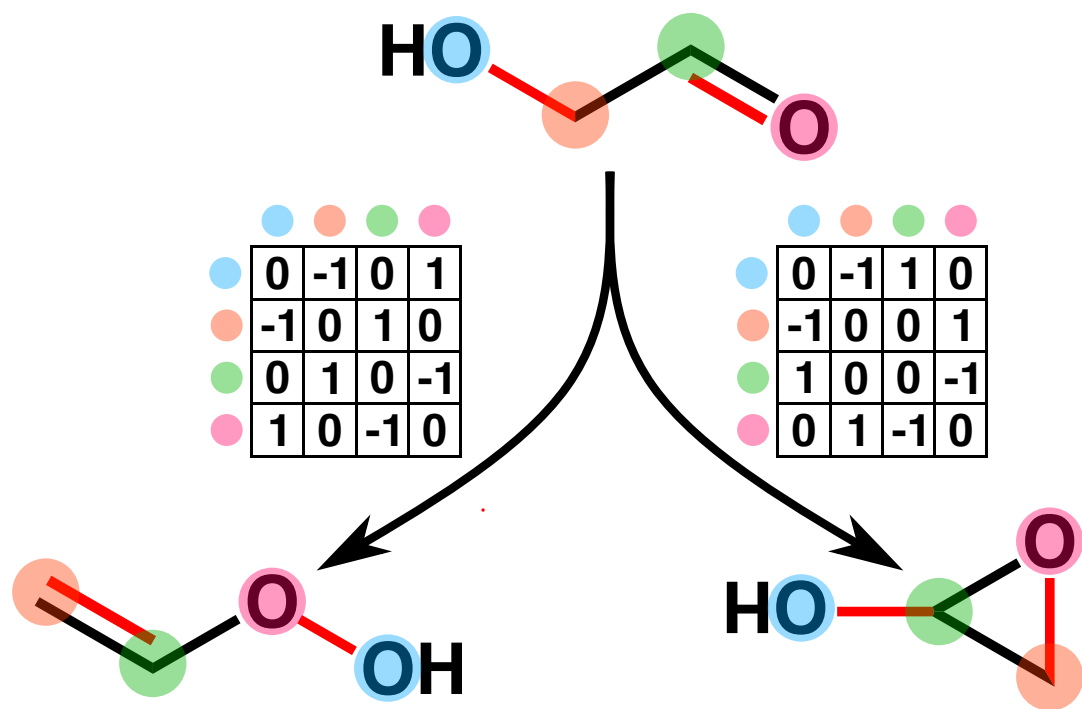
$$R = B - A$$

$$\begin{matrix} \text{O} \\ \text{C} \\ \text{C} \\ \text{H} \\ \text{H} \\ \text{O} \\ \text{H} \\ \text{H} \end{matrix} \begin{bmatrix} 0 & -1 & 1 & 0 & 0 & 0 & 0 & 0 \\ -1 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 0 & 0 & -1 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & -1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

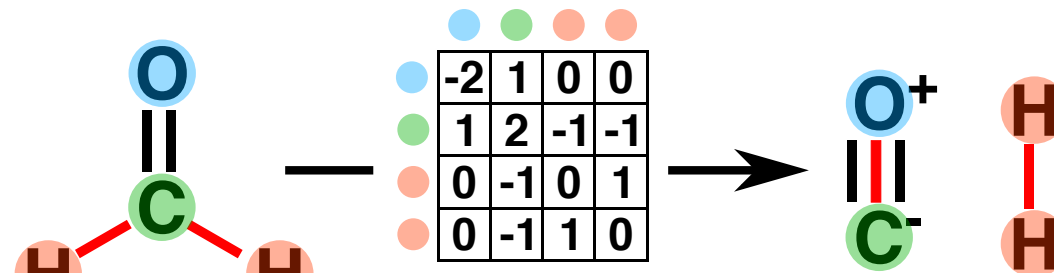
This is essentially a way of converting arrow pushing into something we can automate and interpret by a program

YARP: Elementary Reaction Step(s)

For full-octet neutral organic molecules, “break 2 bonds form 2 bonds” (**b2f2**) is the simplest ERS that yields non-trivial closed-shell neutral products.



Example with concomitant e-transfer



Note: for each pair of broken bonds two distinct reactions are possible

YARP supports a more general suite of ERS(s), but **b2f2** strikes a useful balance that scales as N^2 with reactant size and generates **b1** and **b2f1** type products for free.

YARP: Elementary Reaction Step(s)

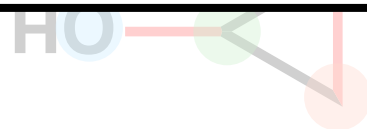
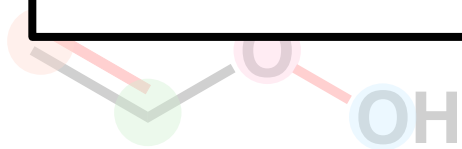
For full-octet neutral organic molecules, “break 2 bonds form 2 bonds” (**b2f2**) is the simplest ERS that yields non-trivial closed-shell neutral products.



Example with concomitant e-transfer

All b3f3 and b4f4 products are b2f2 decomposable

This means that using only b2f2 won't miss any products, but it will potentially miss important transition states (i.e., by predicting a sequential mechanism when a concerted mechanism is favored)

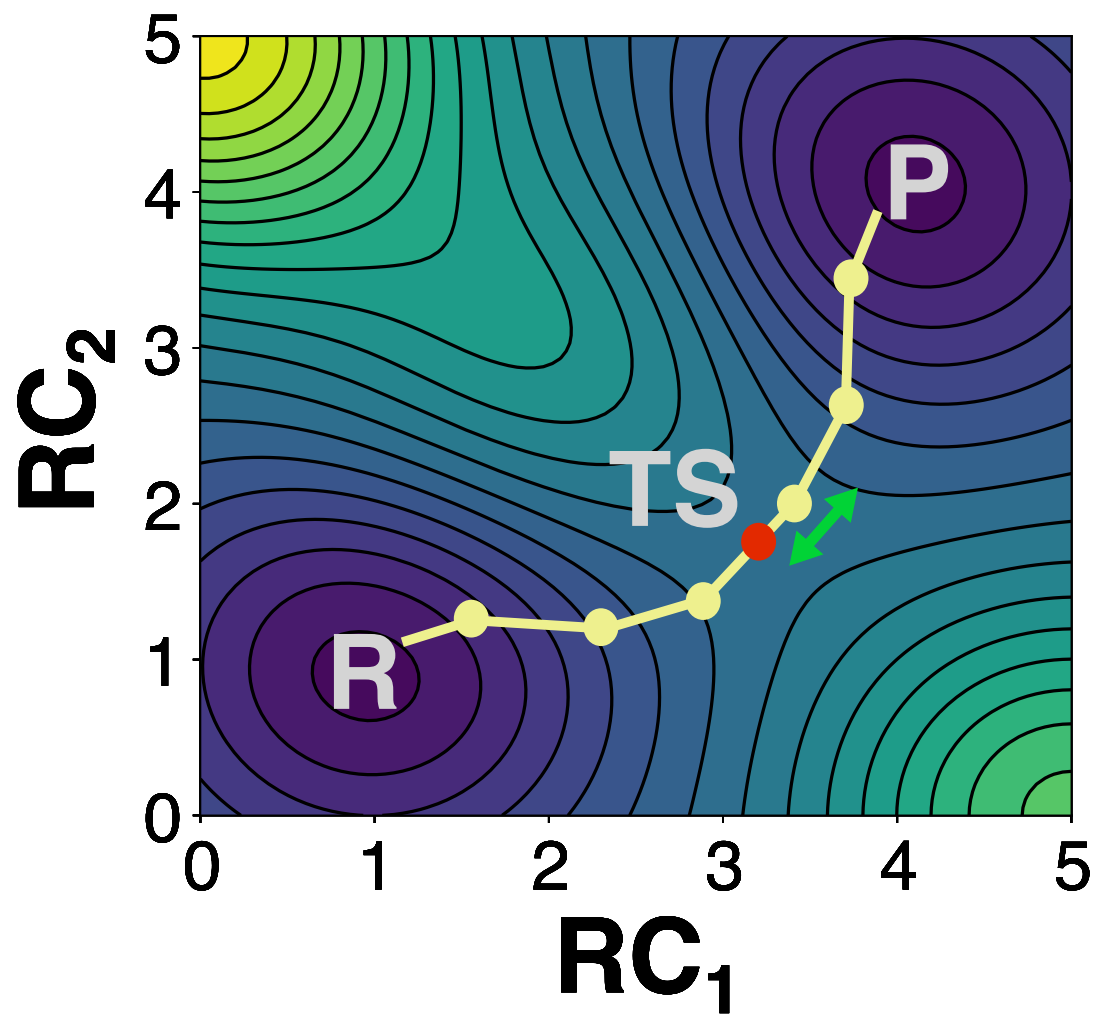


Note: for each pair of broken bonds two distinct reactions are possible

YARP supports a more general suite of ERS(s), but **b2f2** strikes a useful balance that scales as N^2 with reactant size and generates **b1** and **b2f1** type products for free.

YARP: Pseudo-1D Transition State Searches

Searching for saddle points in $3N-6$ space is expensive.



Double-Ended Searches: Using the product and reactant geometries as a constraint, the search can be recast as an effective 1-D search about the connecting coordinate.

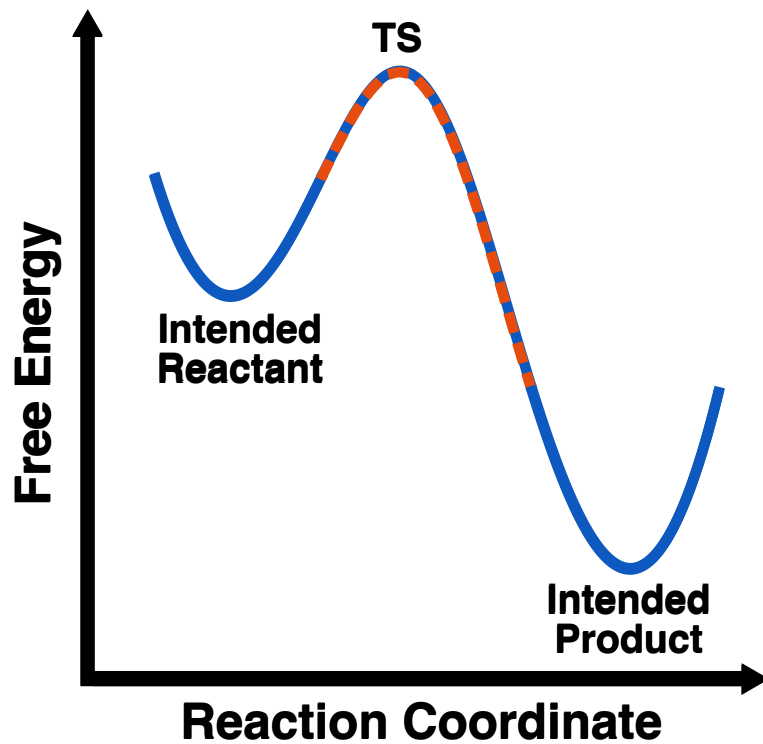
Many Flavors:

- (Growing, Freezing) String methods
- Nudged-Elastic band methods

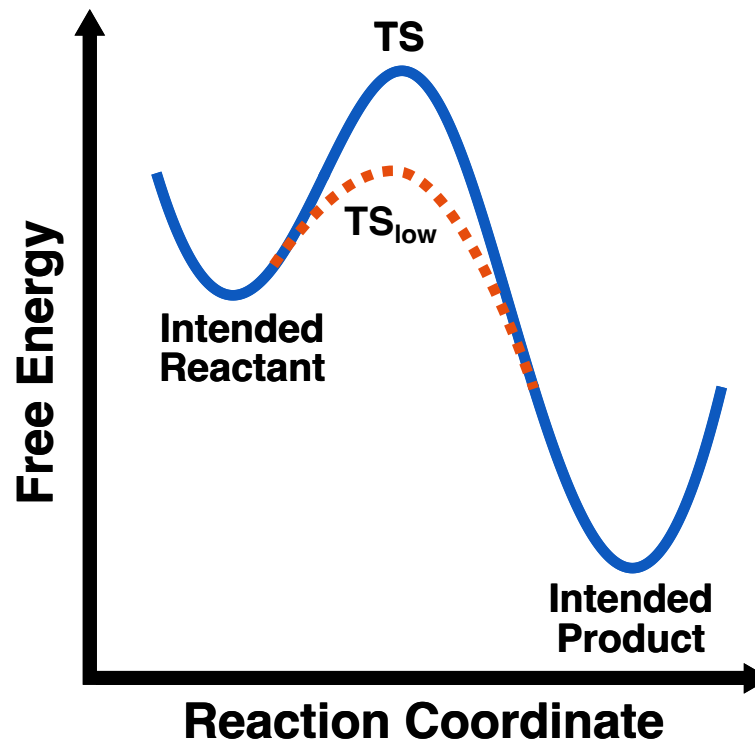
Peters, B.; Heyden, A.; Bell, A. T.; Chakraborty, A. *J. Chem. Phys.* **2004**, 120 (17), 7877–7886.

We like strings because they avoid the bad initial pathway problems of bands.

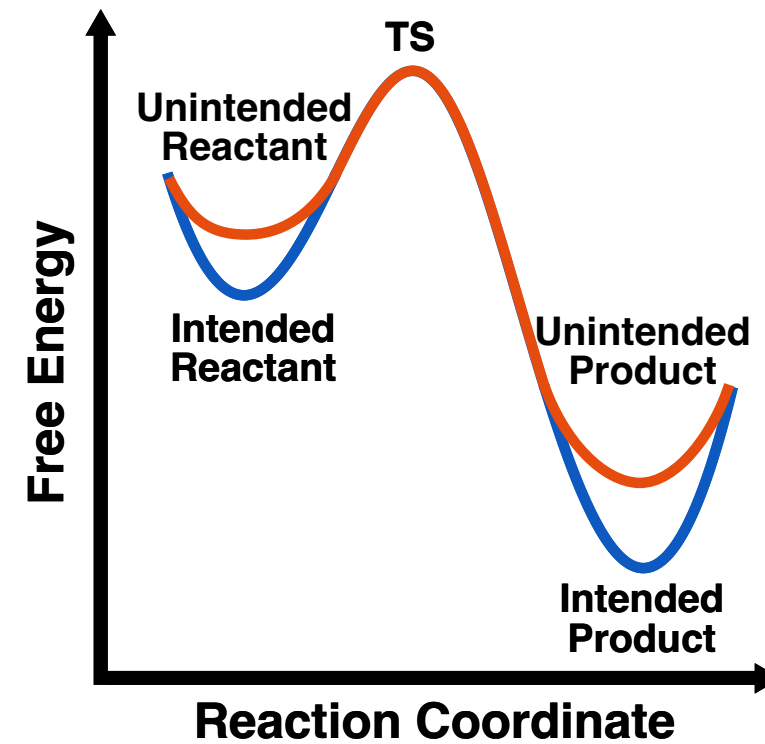
Three Sources of Error in TS Searches



Failure to localize a transition state for a given $A \rightarrow B$ reaction

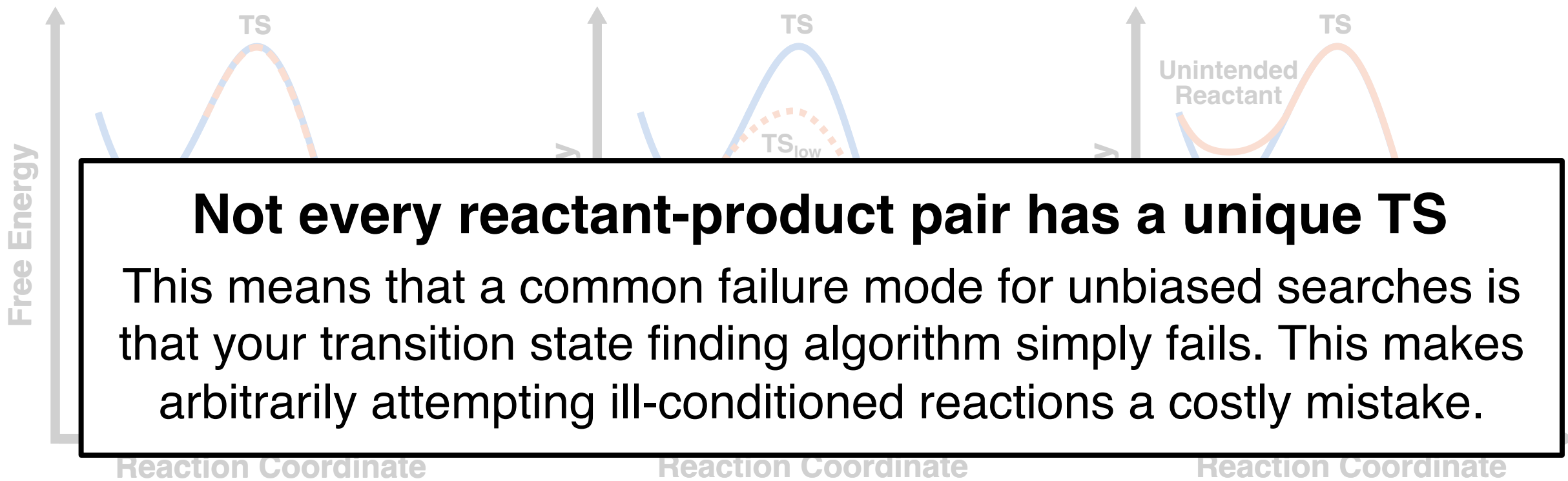


Failure to localize the kinetically relevant TS for a given reaction



Failure to localize an intended TS for a given reaction

Three Sources of Error in TS Searches

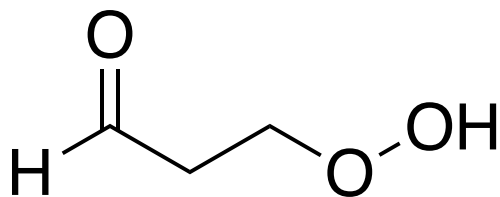


Failure to localize a transition state for a given $A \rightarrow B$ reaction

Failure to localize the kinetically relevant TS for a given reaction

Failure to localize an intended TS for a given reaction

Testing YARP on a Unimolecular Decomposition Problem

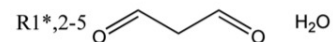


The 3-hydroperoxypropanal reaction network out to b4f4 was recently published as a benchmark for 5 reaction discovery methods.

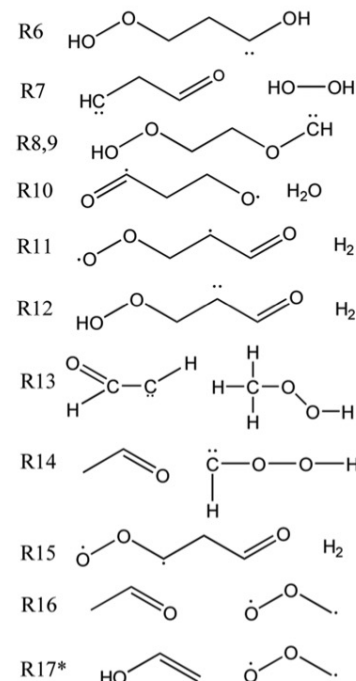
Grambow, C. A, Suleimanov, Y. V. et al. *J. Am. Chem. Soc.* **2018**, 140 (3), 1035–1048.



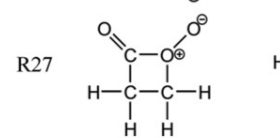
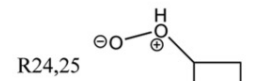
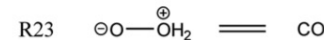
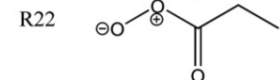
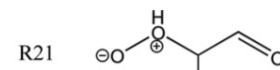
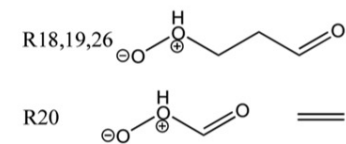
H₂O + malondialdehyde channels:



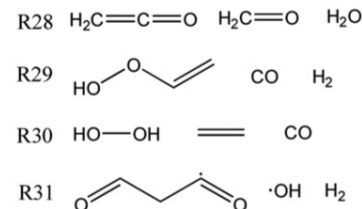
Biradical products including carbenes and the Criegee intermediates:



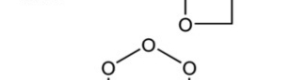
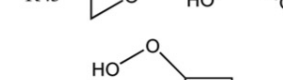
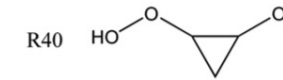
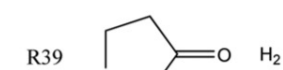
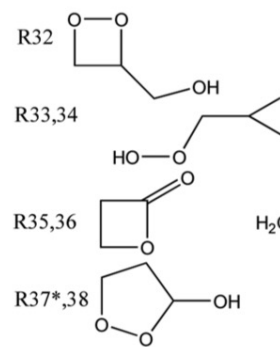
Zwitterionic structures:



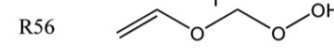
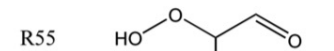
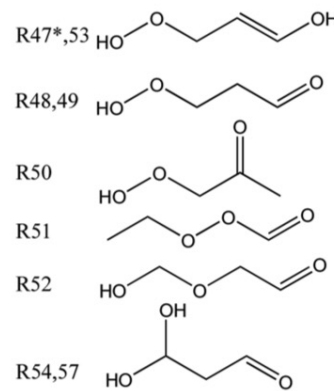
Channels with three products except zwitterionic structures:



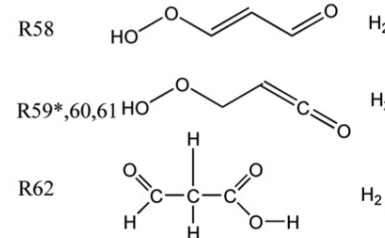
Channels with cyclic products:



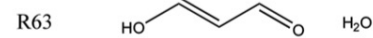
Stable (not radical or zwitterionic) unimolecular noncyclic channels:



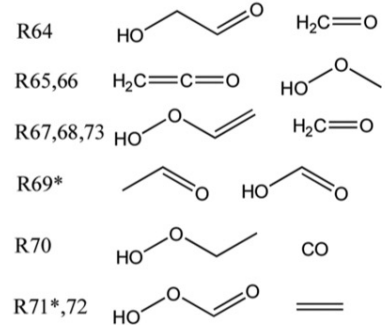
H₂ elimination channels:



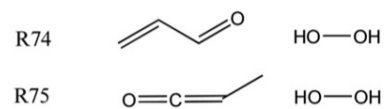
Non-malondialdehyde H₂O elimination channel:



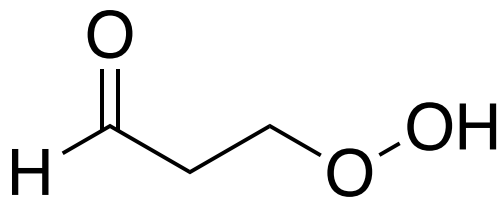
CH₂-CH₂ or CH₂-CHO bond breaking and forming two non-cyclic products:



HOOH elimination channels:

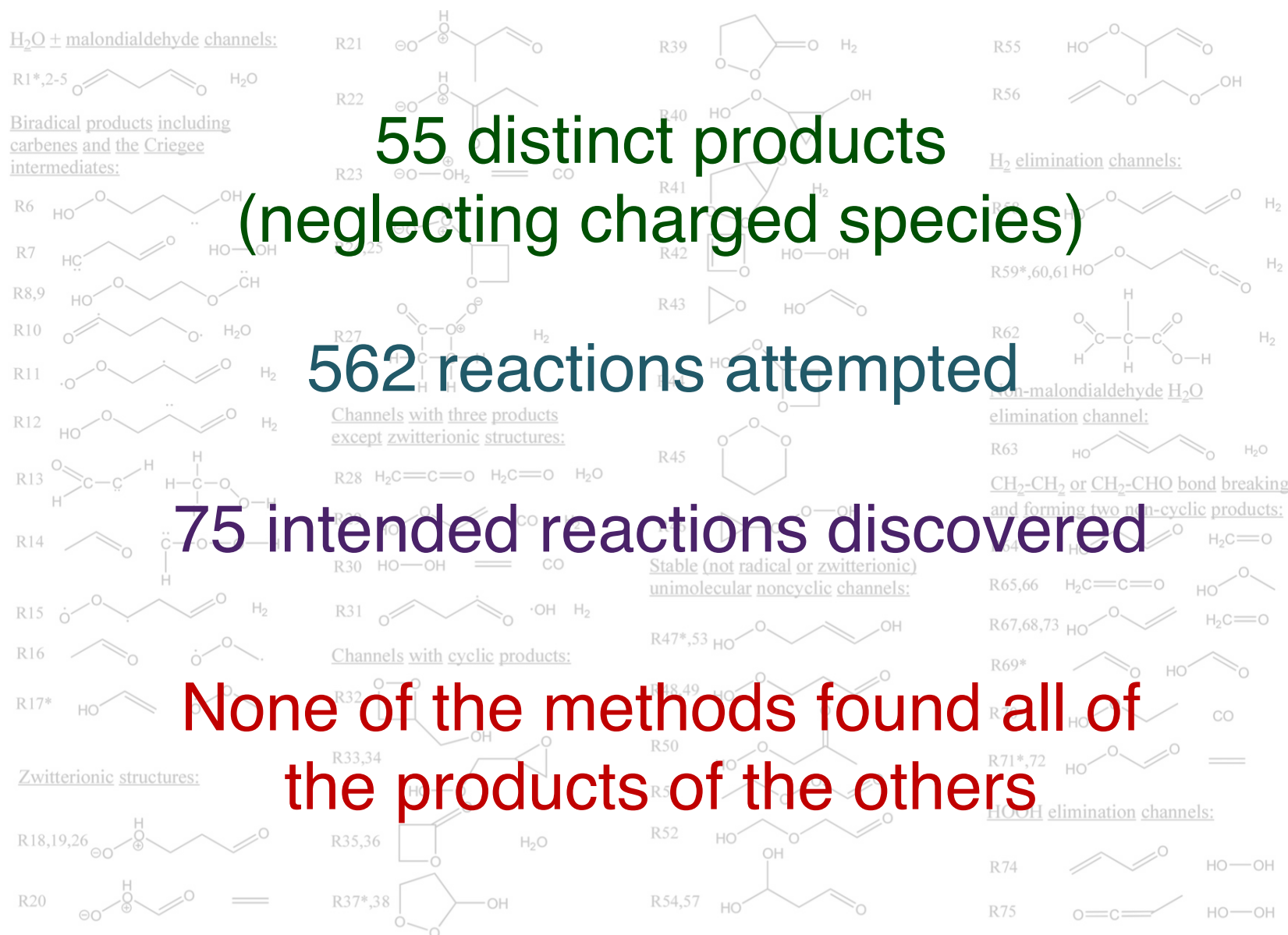


Testing YARP on a Unimolecular Decomposition Problem



The 3-hydroperoxypropanal reaction network out to b4f4 was recently published as a benchmark for 5 reaction discovery methods.

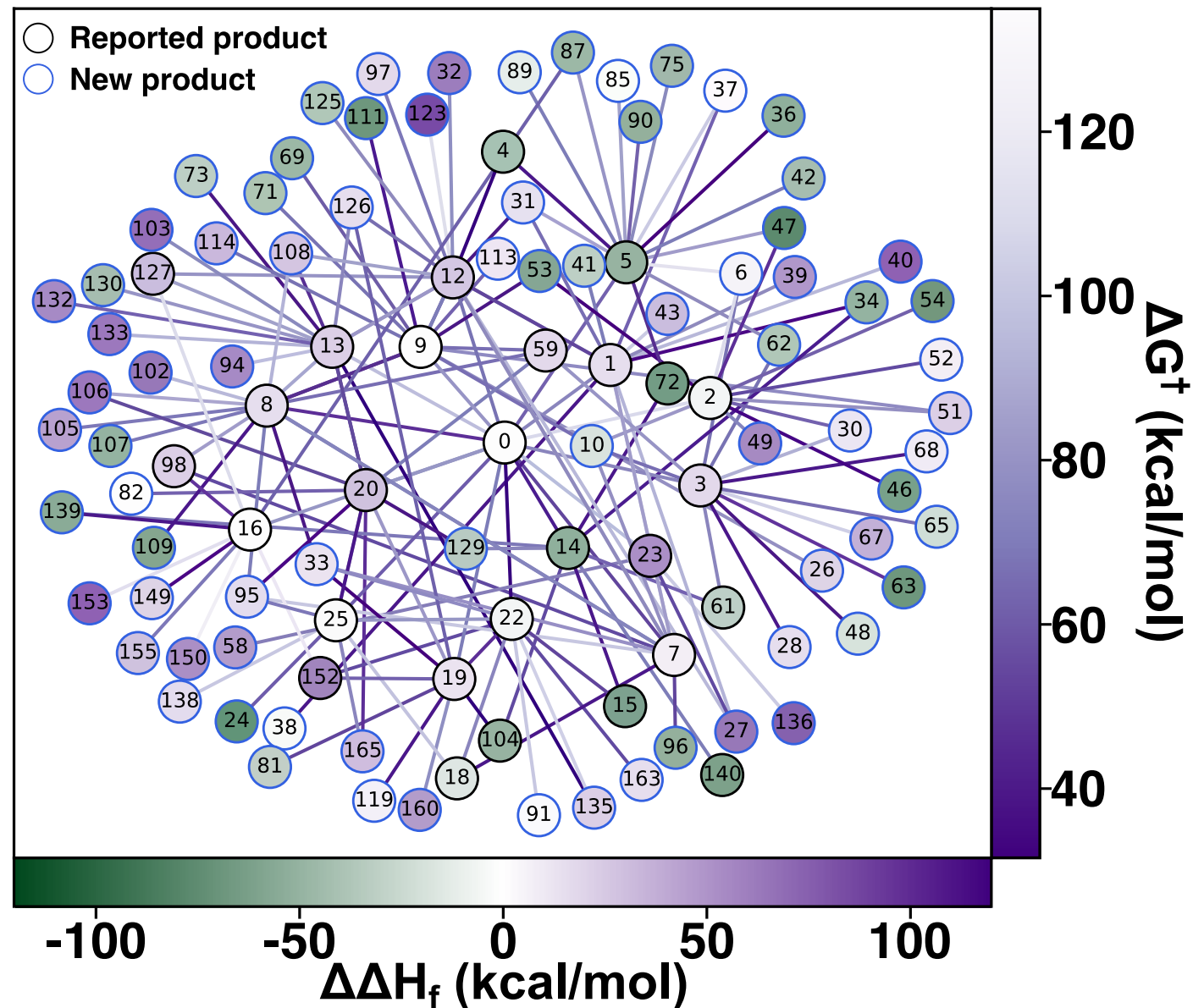
Grambow, C. A, Suleimanov, Y. V. et al. *J. Am. Chem. Soc.* **2018**, 140 (3), 1035–1048.



3-Hydroperoxypropanal - Reaction Network

We used YARP to recursively elucidate the 3-hydroperoxypropanal unimolecular thermal degradation network for comparison with Grambow et al.

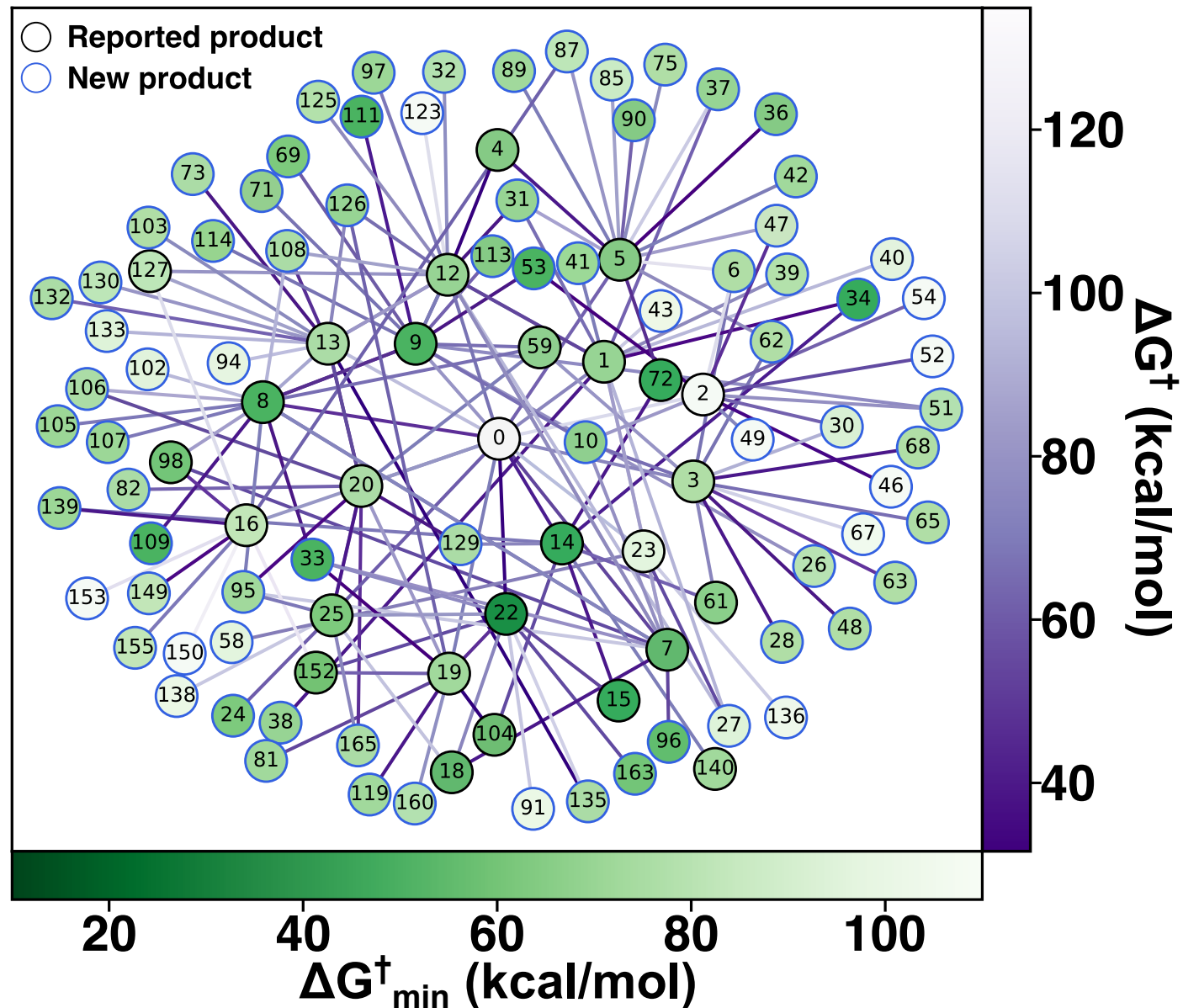
YARP finds **all known products** of this thermal decomposition network, as well as new products (77), and new reactions (157).



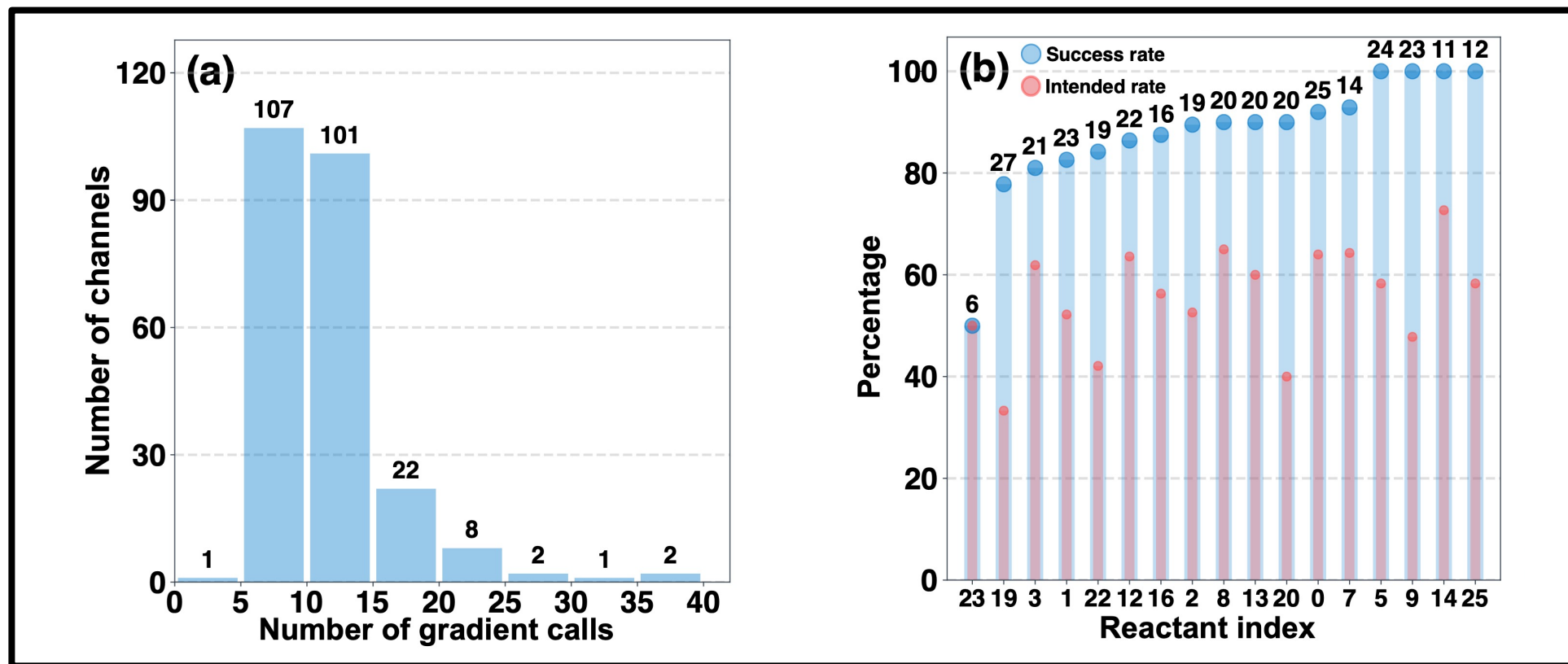
3-Hydroperoxypropanal - Reaction Network

We used YARP to recursively elucidate the 3-hydroperoxypropanal unimolecular thermal degradation network for comparison with Grambow et al.

YARP finds **all known products** of this thermal decomposition network, as well as new products (77), and new reactions (157).



Predicting More (Reactions) with Less (Cost)



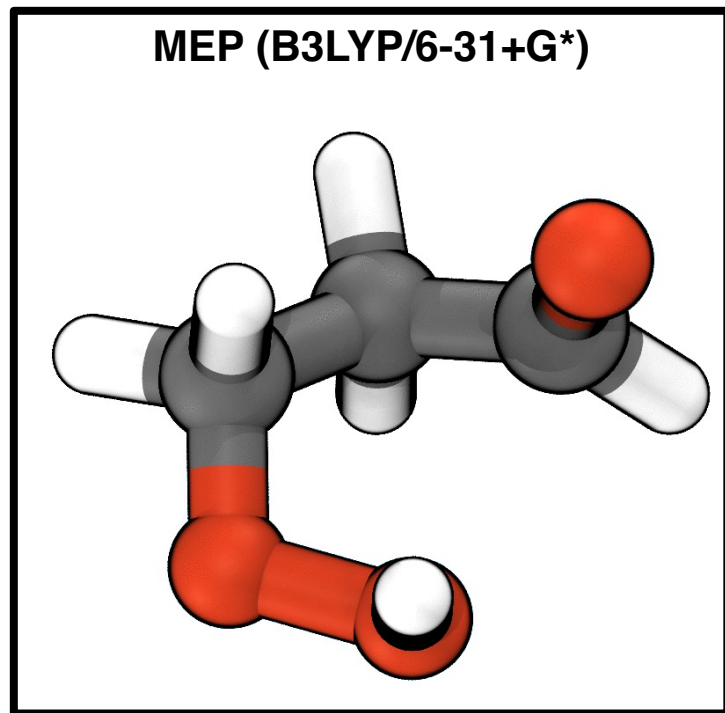
Constructing the whole network required **8364** DFT gradient calls for YARP compared with **756,227** for the earlier benchmark (**100-fold reduction**)

Average success and intended rates for YARP are **81.4%** and **41.1%**, respectively, compared with **38%** and **4%**, in the earlier benchmark.

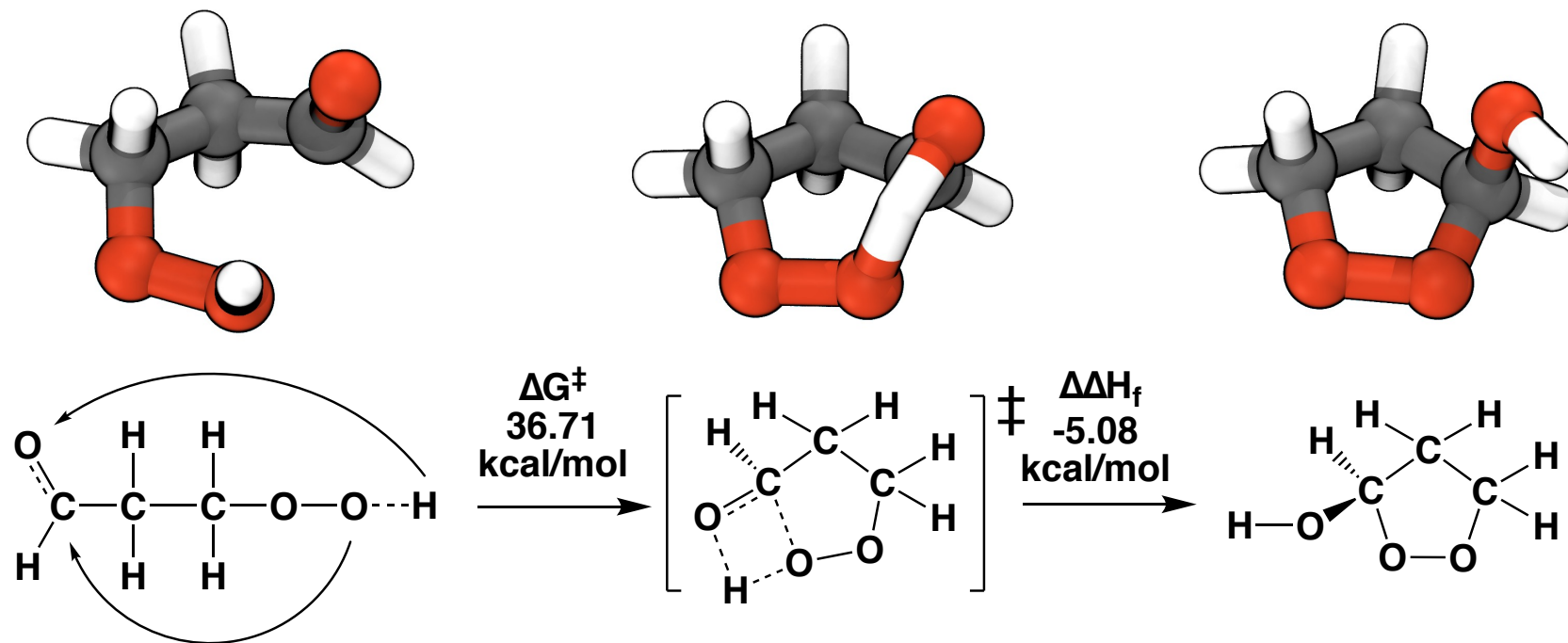
What Happens First?

Jensen, R. K.; Korcek, S.; Mahoney, L. R.; Zinbo, M. *JACS* **1979**, 101, 7574

The Korcek Mechanism



According to YARP, this is the lowest barrier degradation product.



Validated 30 years later by Green and Truhlar:

Jalan, A.; Alecu, I. M.; Meana-Pañeda, R.; Aguilera-Iparraguirre, J.; Yang, K. R.; Merchant, S. S.; Truhlar, D. G.; Green, W. H. *JACS* **2013**, 135 (30), 11100–11114.

Recent A → ? and A → ? → B Case Studies

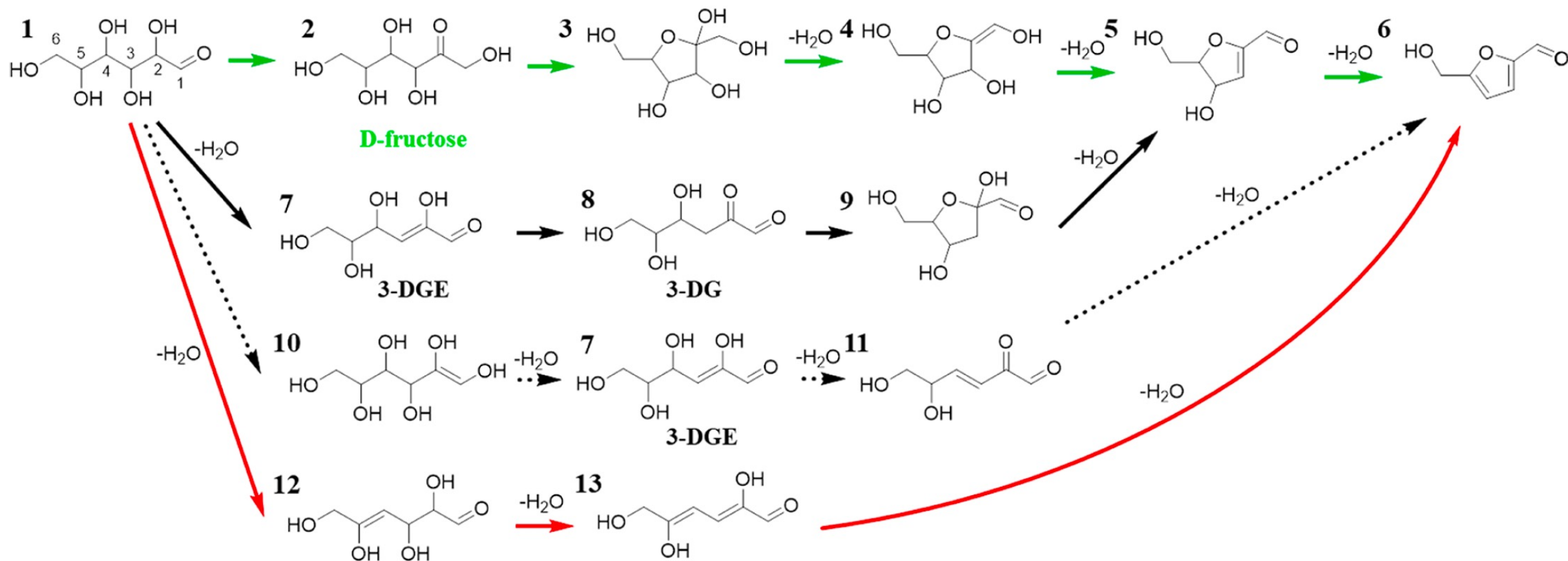
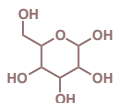


Figure 1. Proposed pathways in literature from glucose to HMF, namely the fructose path (green), 3-DG paths (black and black dotted), and direct path (red). The molecules are indicated by numbers and some key molecules are named as follows: **1.** D-glucose; **2.** D-fructose; **3.** D-fructofuranose; **6.** 5-hydroxymethylfurfural (5-HMF); **7.** 3-deoxyglucos-2-ene (3-DGE); **8.** 3-deoxyglucosone (3-DG); and **10.** hex-1-ene-1,2,3,4,5,6-hexaol (enol form of glucose).

Glucose Pyrolysis Network from YARP Exploration




To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

At each iteration:

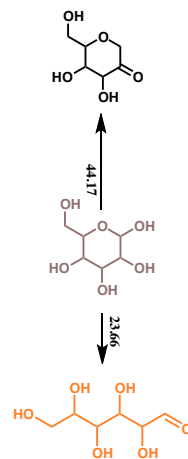
(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are determined by the minimum barrier to a given product (with a window)

(3) Water catalyzed reactions are considered for all H-transfers

Depth 1: 

Glucose Pyrolysis Network from YARP Exploration



Depth 1: 
Depth 2: 

To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

At each iteration:

- (1) all b2f2 reactions are explored for active nodes.
- (2) Active nodes are determined by the minimum barrier to a given product (with a window)
- (3) Water catalyzed reactions are considered for all H-transfers

Glucose Pyrolysis Network from YARP Exploration

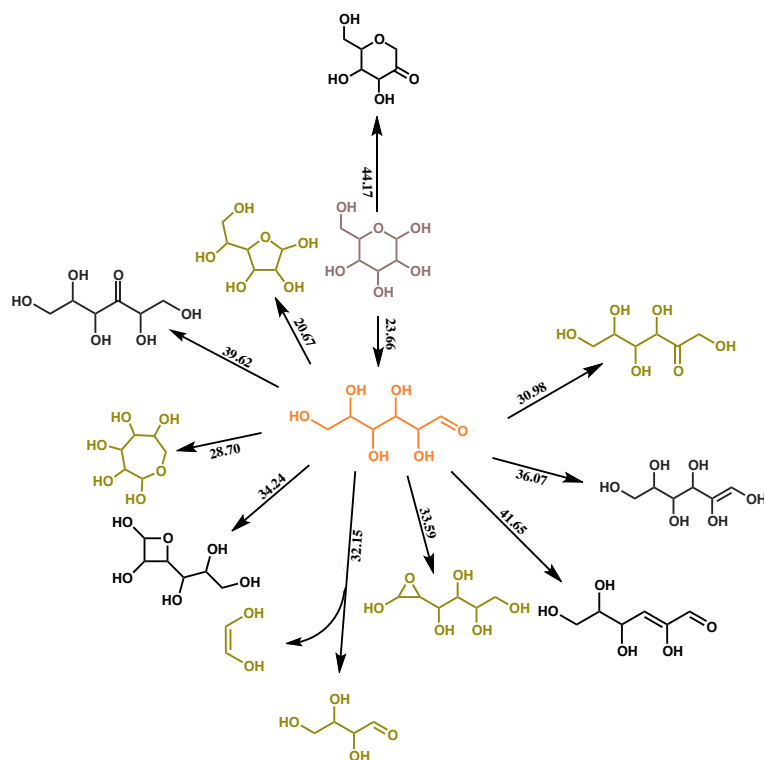
To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

At each iteration:

(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are determined by the minimum barrier to a given product (with a window)

(3) Water catalyzed reactions are considered for all H-transfers



Glucose Pyrolysis Network from YARP Exploration

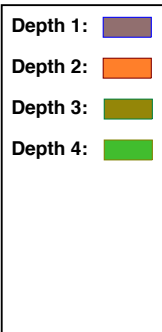
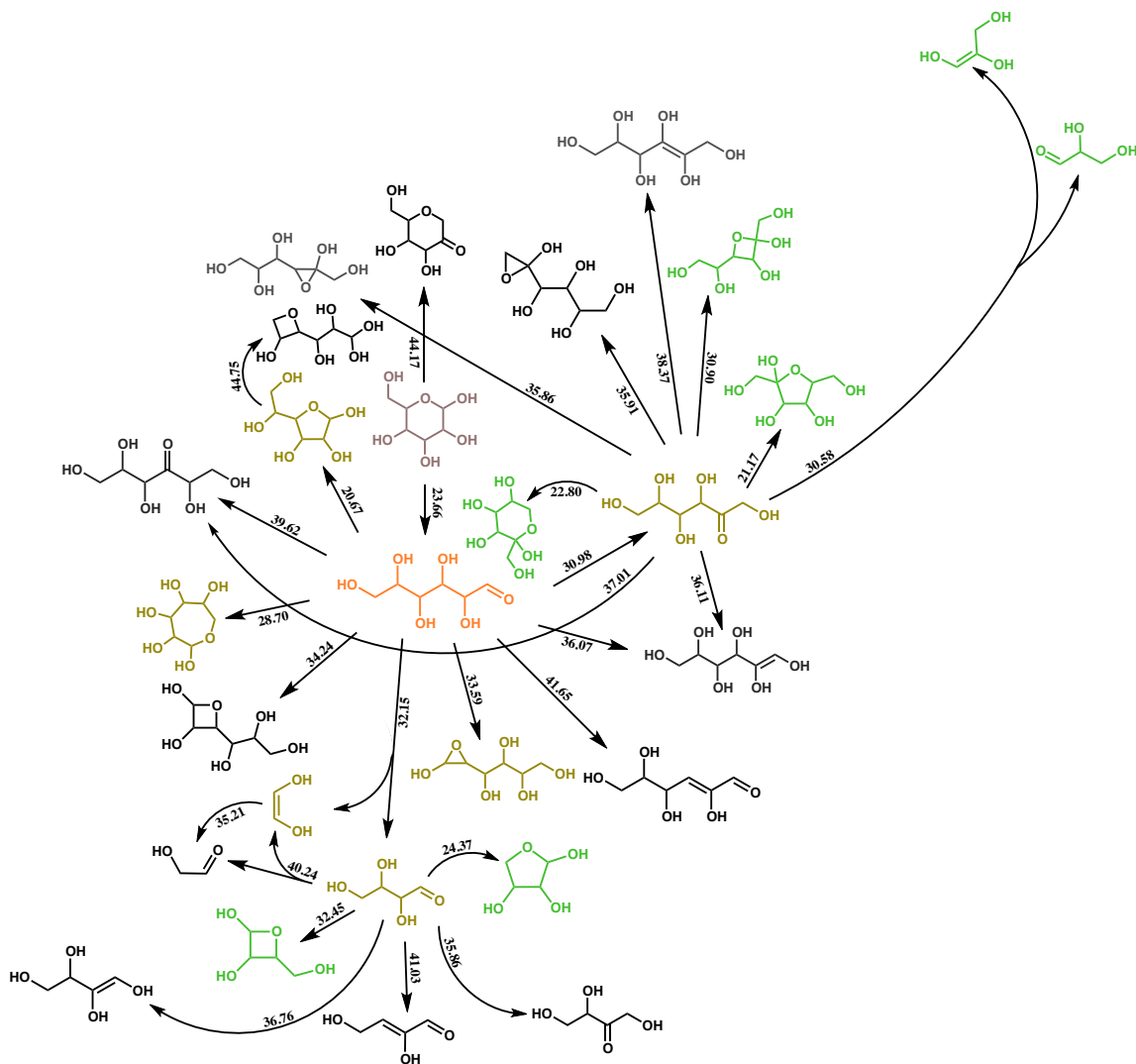
To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

At each iteration:

(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are determined by the minimum barrier to a given product (with a window)

(3) Water catalyzed reactions are considered for all H-transfers



Glucose Pyrolysis Network from YARP Exploration

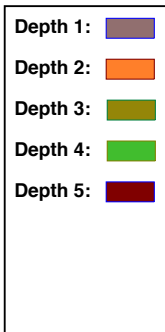
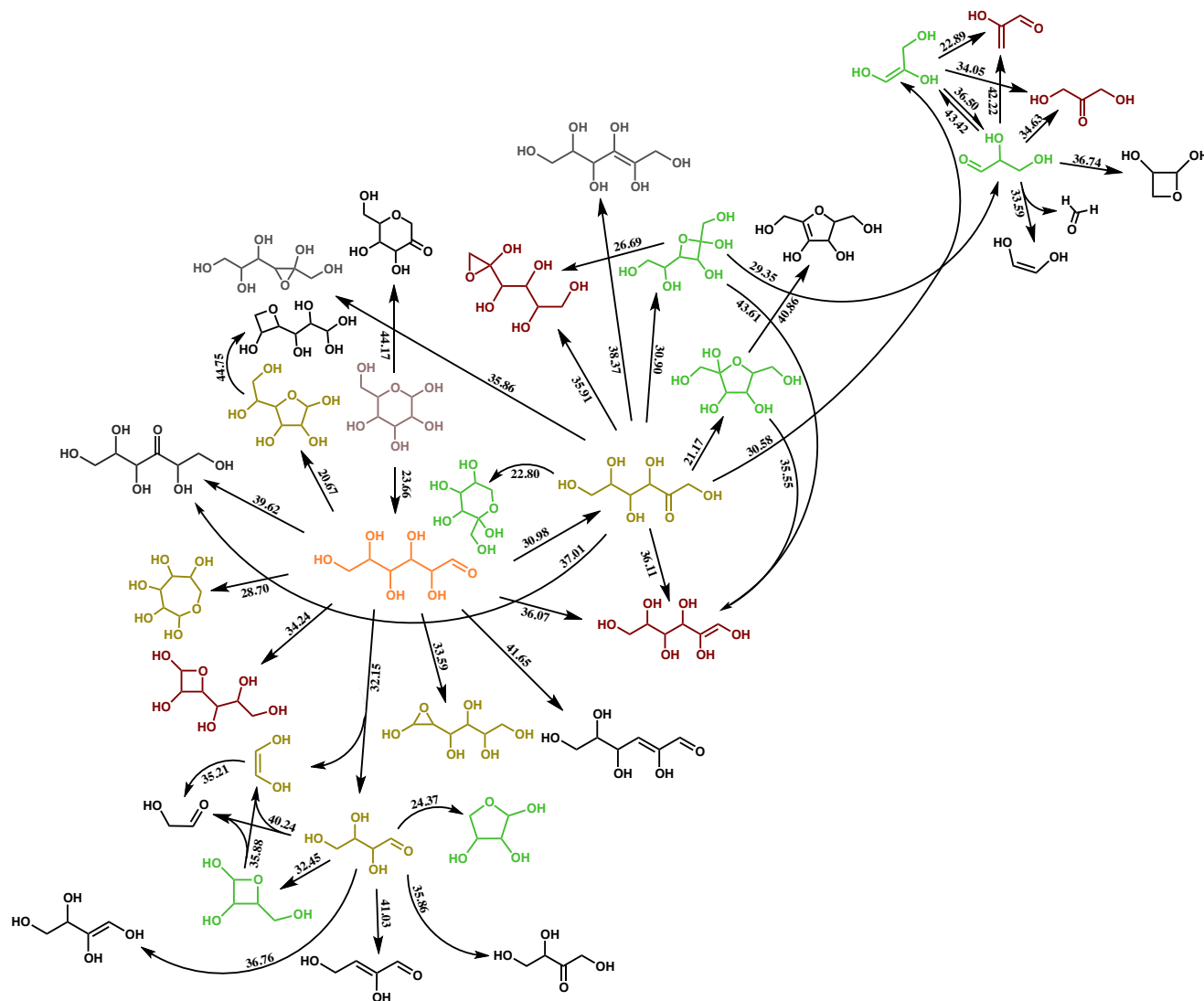
To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

At each iteration:

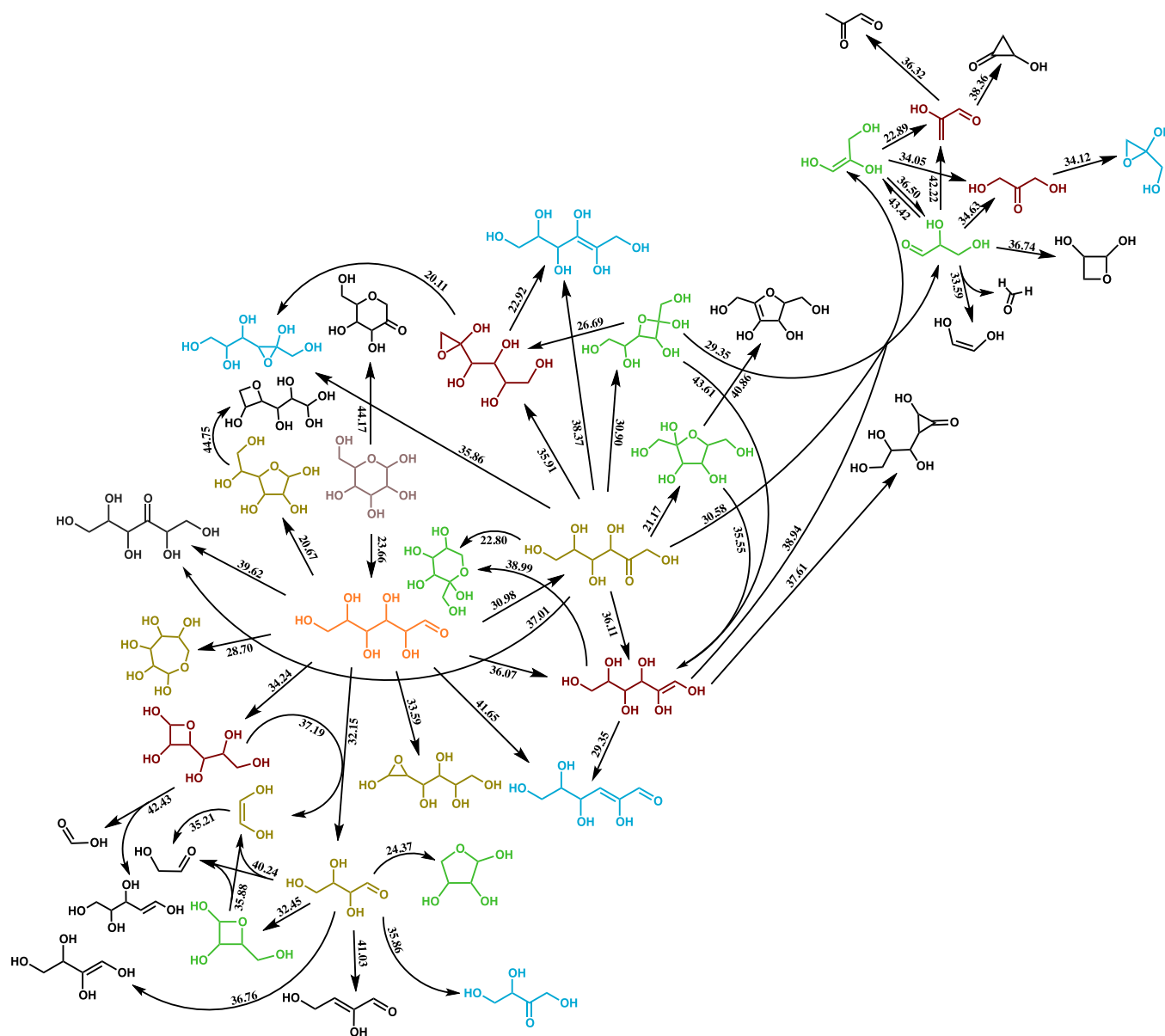
(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are determined by the minimum barrier to a given product (with a window)

(3) Water catalyzed reactions are considered for all H-transfers



Glucose Pyrolysis Network from YARP Exploration



To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

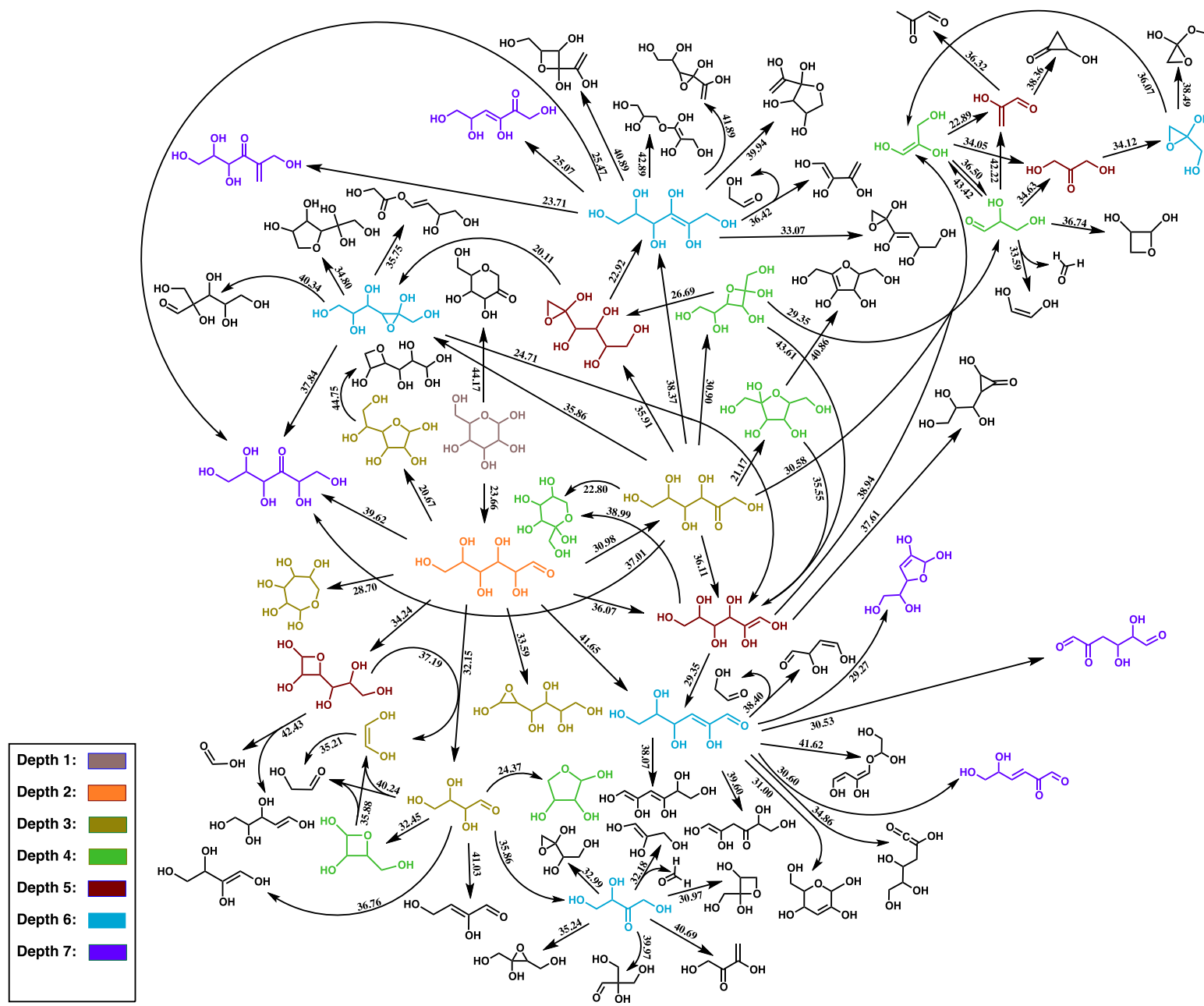
At each iteration:

(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are determined by the minimum barrier to a given product (with a window)

(3) Water catalyzed reactions are considered for all H-transfers

Glucose Pyrolysis Network from YARP Exploration



To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

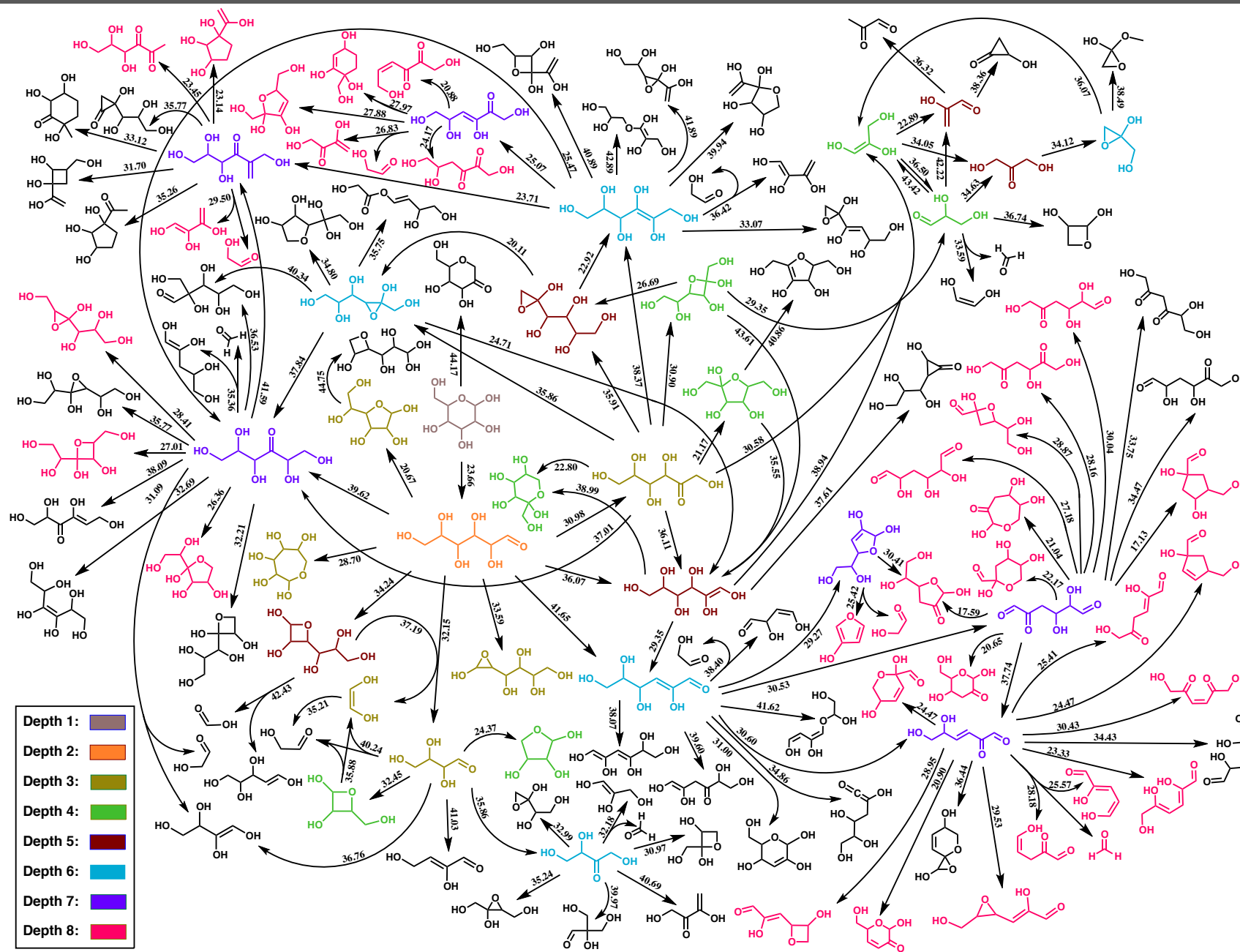
At each iteration:

(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are determined by the minimum barrier to a given product (with a window)

(3) Water catalyzed reactions are considered for all H-transfers

Glucose Pyrolysis Network from YARP Exploration



To perform a deep network exploration, we've implemented a modified version of Dijkstra's algorithm

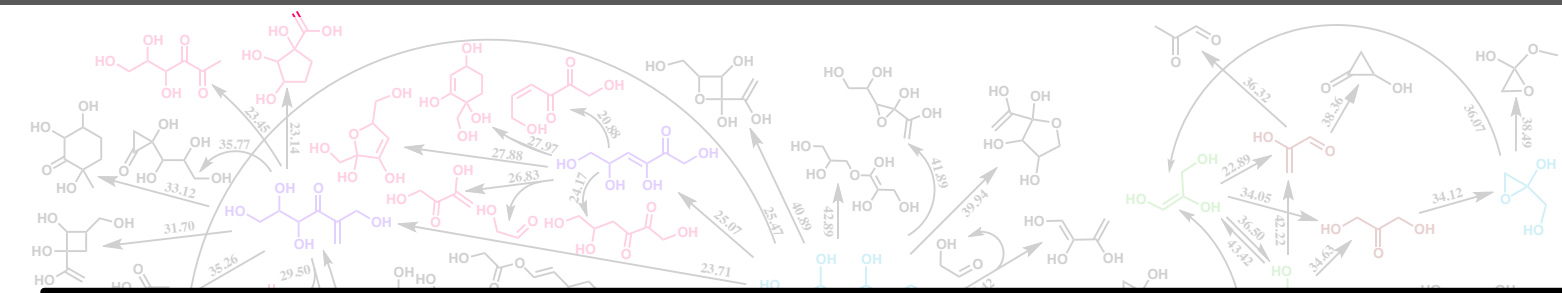
At each iteration:

(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are determined by the minimum barrier to a given product (with a window)

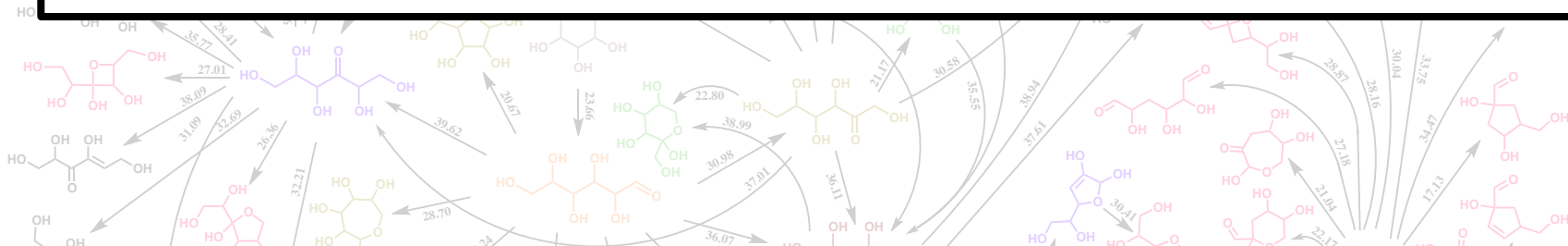
(3) Water catalyzed reactions are considered for all H-transfers

Glucose Pyrolysis Network from YARP Exploration



To perform a deep network exploration, we've implemented a modified version of Dijkstra's

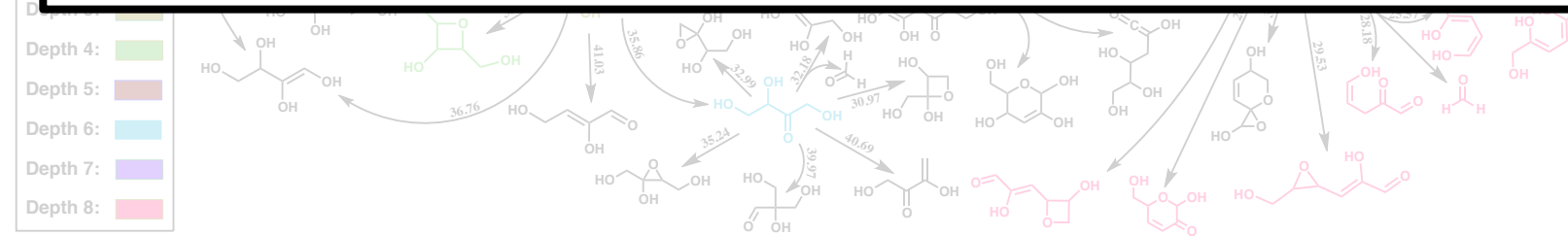
By following the kinetically favorable pathways, YARP spontaneously discovers all major pyrolysis products.



(1) all b2f2 reactions are explored for active nodes.

(2) Active nodes are

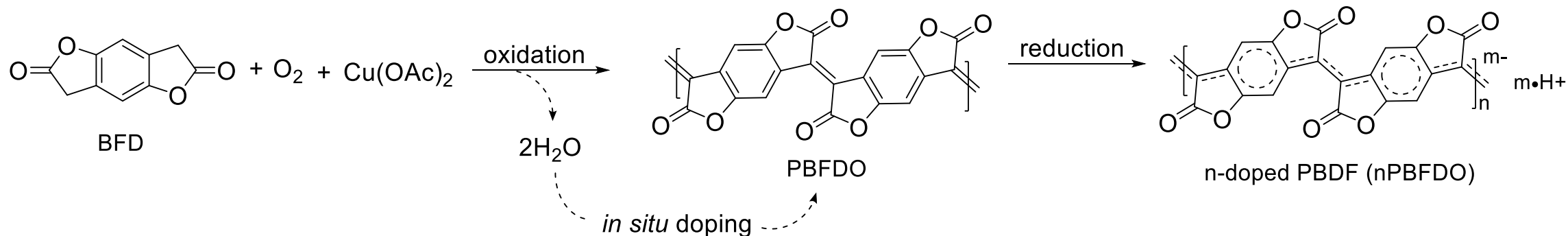
Pathways to minor products can also be found using backward searches.



(3) Water catalyzed reactions are considered for all H-transfers

A Recent $A \rightarrow ? \rightarrow B$ Case Study

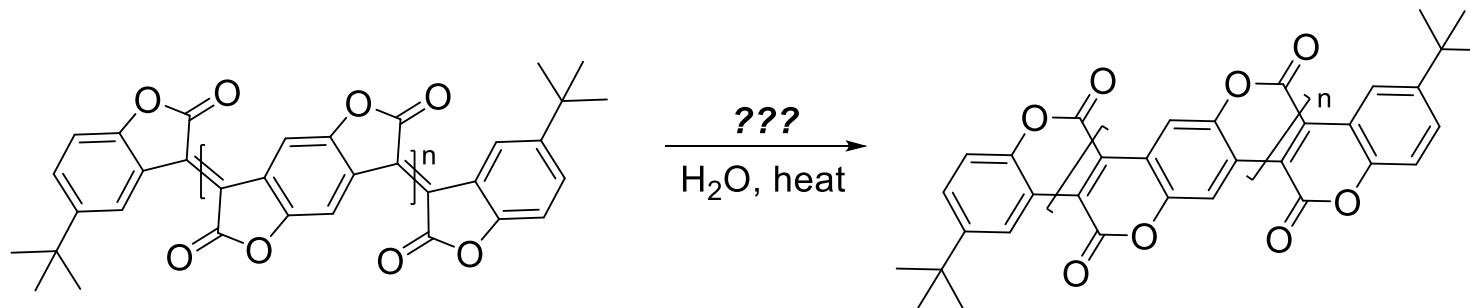
PBDF is a new record breaking n-type semi-conducting plastic



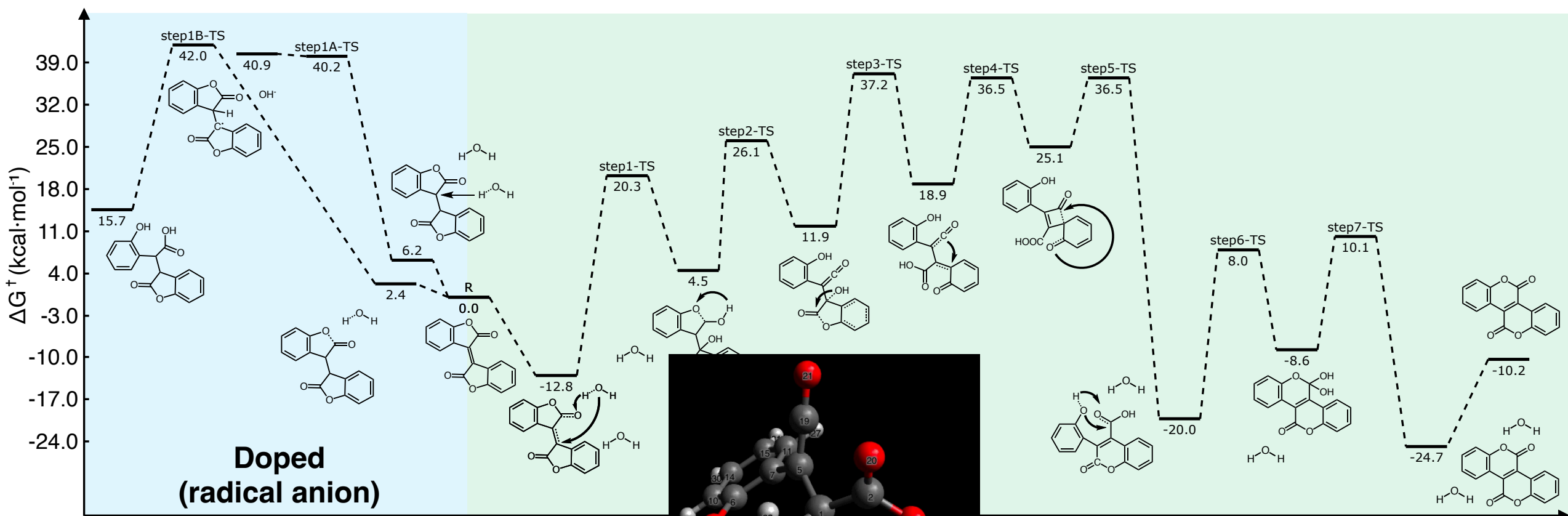
Small molecule studies show a significant synthetic side product



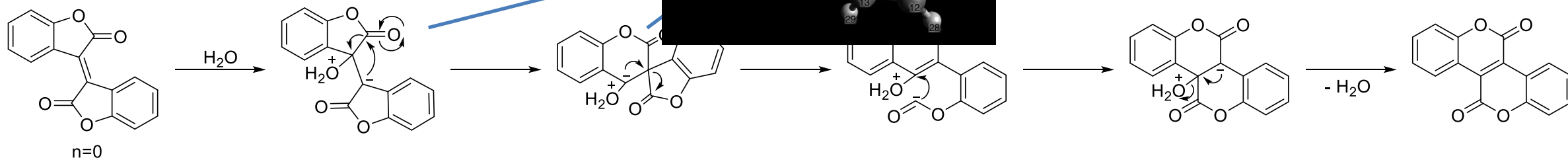
To what extent is this occurring in the polymer? or why don't we see it in the polymer?



Pathways identified by YARP



- ✓ Proposed reaction pathways



Outlook and Acknowledgements

Students: Qiyuan Zhao, Tyler Pasut

State-of-the-art:

- The accurate calculation of thermodynamic properties has become routine in many scenarios. Major opportunities lie in automation, systemization, and low-cost models.
- Practical solutions to the $A \rightarrow ? \rightarrow B$, $A \rightarrow B + ?$, and $A \rightarrow ?$ problems are now available. We envision black-box tools for non-experts in the near future that will assist in hypothesis generation and potentially reactivity screening.



- P2SAC and ONR for funding.
- Ray Mentzer (Purdue)
- Katherine Young (Purdue UG)