# _New Computational Tools for Predicting Reactivity_

**Brett M. Savoie**

Davidson Associate Professor of Chemical Engineering,
Purdue University

**Students: Qiyuan Zhao, Tyler Pasut, Michael Woulfe**

P2SAC Fall Conference, Purdue University, 12/5/23

## A → B

• To safely plan a known reaction, we need access to solid thermodynamic data (e.g., $\Delta H_f$, $S°$, $C_v$) to understand and classify risks.

• This is a **"known unknown"** in that we know the reaction, A → B, but we need values for a few unknown variables.

## A → ? → B ; A → B + ? ; A → ?

• **A → ? → 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 → B + ?** (unknown side-reaction) and **A → ?** (unknown main product), problems have similar **"unknown unknown"** characteristics.

**A → B**

- For a reaction, we need access to solid thermodynamic data (e.g., $\Delta H_f$, $S°$, $C_v$) to understand and classify risks.

- This is a "known unknown" in that we know the reaction, A → B, but we need values for a few unknown variables.

# TAFFI Component Increment Theory (TCIT)

**A → ? → B ; A → B + ? ; A → ?**

- **A → ? → 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 → B + ?** (unknown side-reaction) and **A → ?** (unknown main product), problems have similar "unknown unknown" characteristics.

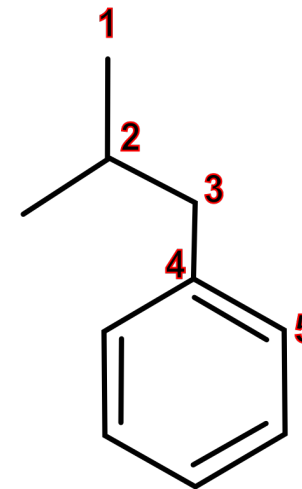# Yet Another Reaction Program (YARP)

## Benson Group Theory:

• The idea is to decompose molecular properties ($\Delta H_f$, $S^o$, $C_v$) 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 $\Delta H_f$:** -5.15 +/- 0.34 kcal/mol

## Benson Group Theory:

• The idea is to decompose molecular properties ($\Delta H_f$, $S^\circ$, $C_v$) 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.
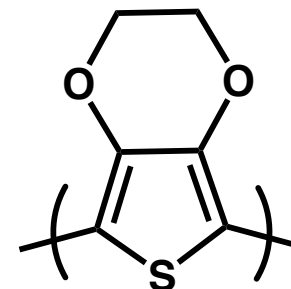
**$\Delta H_f$ from modern quantum chemistry**



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

**Benson Group Theory:**

• The idea is to decompose molecular properties ($\Delta H_f$, $S°$, $C_v$) as the sum of "group" contributions.

• Grou... on trus... data, a...

**ΔH_f from modern quantum chemistry**

200

(mol)

200

572 small molecules

G4

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

**Prob...**

• **Spec...** 200
formal...

• **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

## The fundamental idea

**TCIT is a component theory (2-bond specific)**

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

**T**opology **A**utomated **F**orce **F**ield **I**nteractions
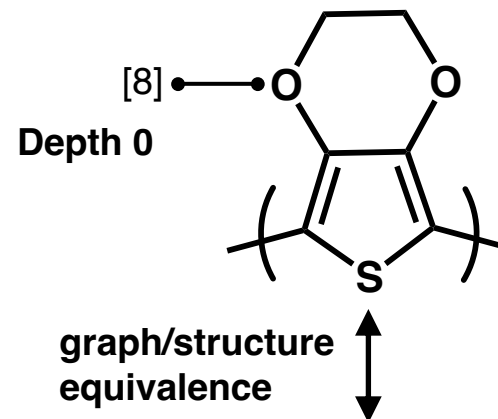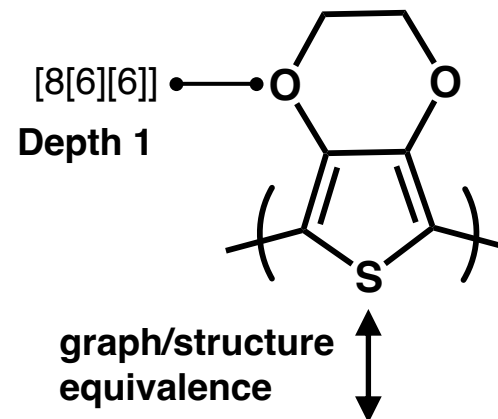
graph/structure equivalence

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**

$$
\begin{array}{c}
\begin{array}{l}S\\C\\C\\C\\C\\O\\C\\C\\H\\H\\O\\H\\H\end{array}
\begin{bmatrix}
0&1&0&0&1&0&0&0&0&0&0&0&0&0\\
1&0&1&0&0&0&0&0&0&0&0&0&0&0\\
0&1&0&1&0&1&0&0&0&0&0&0&0&0\\
0&0&1&0&1&0&0&0&0&0&0&1&0&0\\
1&0&0&1&0&0&0&0&0&0&0&0&0&0\\
0&0&1&0&0&0&1&0&0&0&0&0&0&0\\
0&0&0&0&0&1&0&1&1&1&0&0&0&0\\
0&0&0&0&0&0&1&0&0&0&1&1&1\\
0&0&0&0&0&0&1&0&0&0&0&0&0&0\\
0&0&0&0&0&0&1&0&0&0&0&0&0&0\\
0&0&0&1&0&0&0&1&0&0&0&0&0&0\\
0&0&0&0&0&0&0&1&0&0&0&0&0&0\\
0&0&0&0&0&0&0&1&0&0&0&0&0&0
\end{bmatrix}
\end{array}
$$

**Adjacency matrix for PEDOT monomer**

**The fundamental idea**

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

**TCIT is a <u>component</u> theory (2-bond specific)**



[8] •——•O
**Depth 0**

**T**opology **A**utomated **F**orce **F**ield **I**nteractions

graph/structure equivalence

$$\begin{array}{c} S \\ C \\ C \\ C \\ C \\ O \\ C \\ C \\ H \\ H \\ O \\ H \\ H \end{array} \begin{bmatrix} 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$
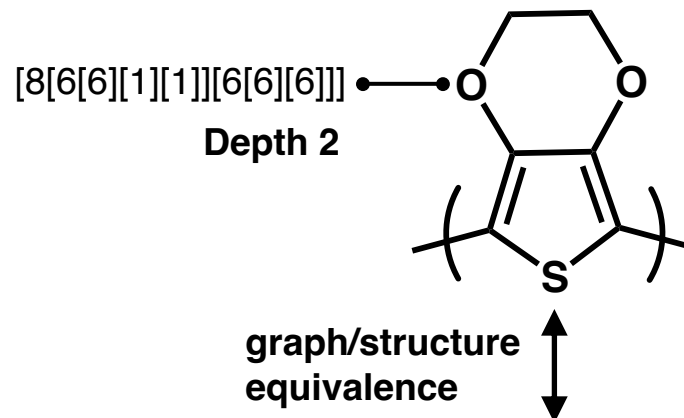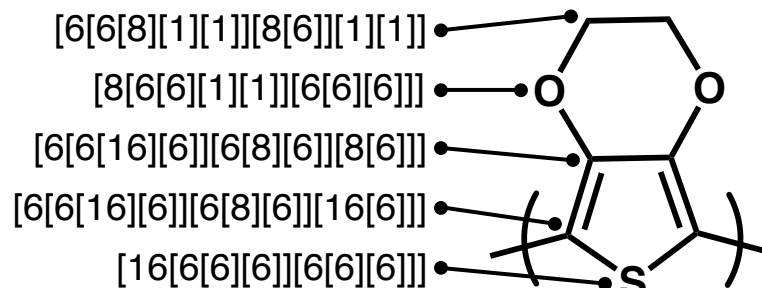
**Adjacency matrix for PEDOT monomer**

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**

## The fundamental idea

**TCIT is a <u>component</u> theory (2-bond specific)**

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

**[8[6][6]]**

**Depth 1**

**T**opology **A**utomated **F**orce **F**ield **I**nteractions

TAFFI

**graph/structure equivalence**

$$\begin{array}{c} S \\ C \\ C \\ C \\ C \\ O \\ C \\ C \\ H \\ H \\ O \\ H \\ H \end{array} \begin{bmatrix} 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{bmatrix}$$

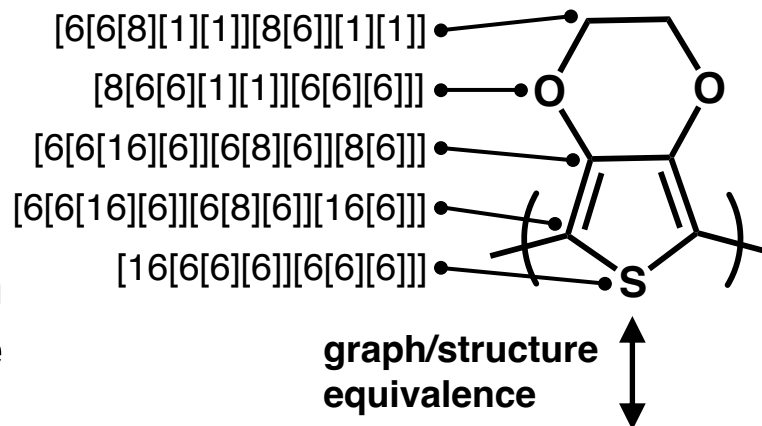**Adjacency matrix for PEDOT monomer**

**P2SAC Publications**

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.

**The fundamental idea**

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

**TCIT is a component theory (2-bond specific)**

[8[6[6][1][1]][6[6][6]]] ●——● O

**Depth 2**

**T**opology **A**utomated **F**orce **F**ield **I**nteractions

graph/structure equivalence

$$
\begin{array}{c}
S \\ C \\ C \\ C \\ C \\ O \\ C \\ C \\ H \\ H \\ O \\ H \\ H
\end{array}
\begin{bmatrix}
0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\
1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 1 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\
0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0
\end{bmatrix}
$$

**Adjacency matrix for PEDOT monomer**

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**

## The fundamental idea

**TCIT is a <u>component</u> theory (2-bond specific)**

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

[6[6[8][1][1]][8[6]][1][1]]
[8[6[6][1][1]][6[6][6]]]
[6[6[16][6]][6[8][6]][8[6]]]
[6[6[16][6]][6[8][6]][16[6]]]
[16[6[6][6]][6[6][6]]]

**T**opology **A**utomated **F**orce **F**ield **I**nteractions

TAFFI

graph/structure equivalence

**Adjacency matrix for PEDOT monomer**

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**

## The fundamental idea

**TCIT is a <u>component </u>theory (2-bond specific)**



[6[6[8][1][1]][8[6]][1][1]]
[8[6][6][1][1]][6[6][6]]]
[6[6[16][6]][6[8][6]][8[6]]]
[6[6[16][6]][6[8][6]][16[6]]]
[16[6][6]][6[6][6]]]

**T**opology **A**utomated
**F**orce **F**ield **I**nteractions

• 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**.

**graph/structure equivalence**

$$\begin{array}{c} S \\ C \\ C \\ C \\ C \\ O \\ C \\ C \\ H \\ H \\ O \\ H \\ H \end{array}\left[\begin{array}{ccccccccccccc} 0 & 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 1 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 \\ 1 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 1 & 0 & 1 & 1 & 1 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 1 & 1 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 1 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \\ 0 & 0 & 0 & 0 & 0 & 0 & 0 & 1 & 0 & 0 & 0 & 0 & 0 \end{array}\right]$$
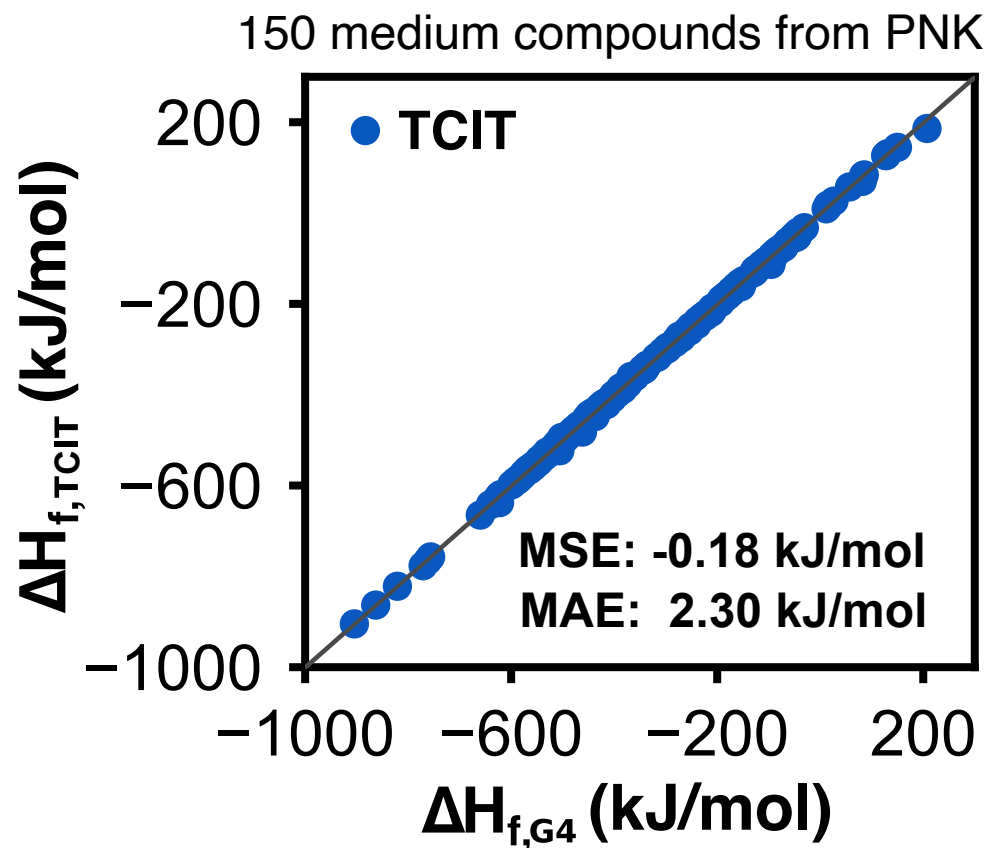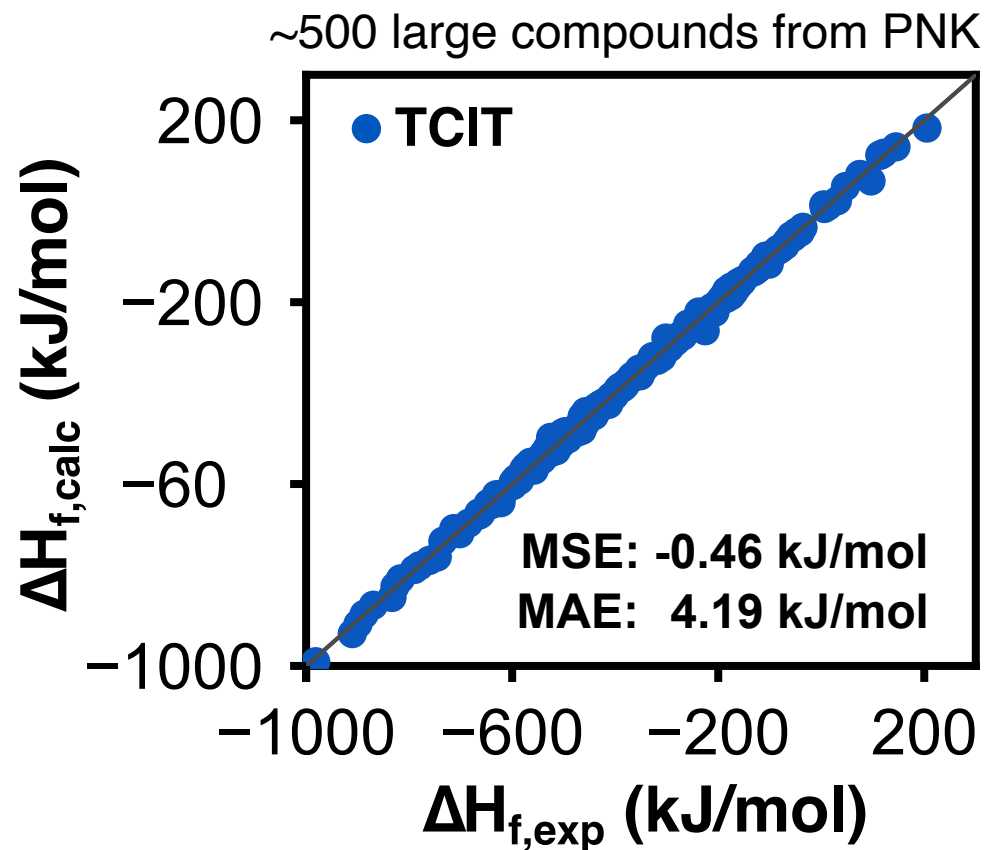
**Adjacency matrix for PEDOT monomer**

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**

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

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

identify components*



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

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

identify components*

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

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

Recursively generate smallest acyclic model compounds

identify components*

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

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

Recursively generate smallest acyclic model compounds

identify components*

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

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

Recursively generate smallest acyclic model compounds

new groups

Savoie Research Group  |

Zhao, Q.; Savoie, B. M. *J. Chem. Info. Model.* **2020**, 60, 2199-2207.

| 17

identify components*

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

How will we select molecules for parameterizing TCIT components?

Recursively generate smallest acyclic model compounds

new groups

Resolve dependencies

Savoie Research Group  l

Zhao, Q.; Savoie, B. M. *J. Chem. Info. Model.* **2020, 60, 2199-2207.**

l 18

identify components*

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

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

Recursively generate smallest acyclic model compounds

new groups

Resolve dependencies

Resolve rank deficiency with elementary constraints

**Prediction target:**

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

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

$\Delta H_{f,TCIT}$ = -259.3 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)**

## 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 $\Delta H_f$ 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.

• Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK[1]

(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.

572 small compounds from PNK



MSE: -0.06 kJ/mol
MAE:  4.19 kJ/mol

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

• Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK[1]

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2$^{nd}$ 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.

150 medium compounds from PNK



- TCIT

MSE: -0.18 kJ/mol
MAE:  2.30 kJ/mol

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

• Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK[1]

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2$^{nd}$ 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.

~500 large compounds from PNK



MSE: -0.46 kJ/mol
MAE:  4.19 kJ/mol

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

• Initial benchmarking set consists of ~1100 **linear** C,H, and O containing compounds from PNK[1]

(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.

~500 large compounds from PNK



$MSE_{TCIT}$ :-0.46 kJ/mol
$MAE_{TCIT}$ : 4.67 kJ/mol
$MSE_{BGIT}$ :-1.71 kJ/mol
$MAE_{BGIT}$ : 5.84 kJ/mol

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.**

**Prediction target:**

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

**Gen 4:**

**Gen 3:**

**Gen 2:**

**Gen 1:**

**Gen 0:**

We database all model compounds and components for reuse.

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

**Current Database:**
- ~35k distinct components for $\Delta H_f$ relevant to organic chemistry
- ~35k distinct G4 calculations on organic molecules.
- ~450 distinct ring corrections

How many components are required to predict the $\Delta H_f$ of **all** (physically relevant) organic molecules?

How many P2SAC funding periods would it take to make a "complete" or "gapless" component theory?

**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 $\Delta H_f$



**C,H,N, and O Containing Molecules**

**Saturation**

Unique Components

30k
25k
20k
15k
10k
5k
0

0    200k    400k    600k    800k    1m

**PubChem Molecules**

**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 $\Delta H_f$

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.

**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 $\Delta H_f$

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**

PubChem is a repository of chemical properties that contains many millions of organic

TCIT now contains all CAVs necessary to predict $\Delta H_f$ of 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.

4000

3000

2

4

6

8

10

12

14

Heavy Atoms

A recurring question is when will TCIT support predictions on **radicals** and **ions**?

**TCIT** already covers neutral close-shell species, so these extensions require us only to predict the difference between the target and the **nearest closed-shell neutral**.

**This amounts to developing models to predict IP/EA/+H⁺/-H⁺**

A recurring question is when will TCIT support predictions on **radicals** and **ions**?

**TCIT** already covers neutral close-shell species, so these extensions require us only to predict the difference between the target and the **nearest closed-shell neutral**.

**This amounts to developing models to predict IP/EA/+H⁺/-H⁺**



Existing Radical/Ion G4 Database

**A → B**

• To safely plan a known reaction, we need access to solid thermodynamic data (e.g., $\Delta H_f$, $S°$, $C_v$) to understand and classify risks.

• This is a **"known unknown"** in that we know the reaction, A → B, but we need values for a few unknown variables.

**A → ? → B ; A → B + ? ; A → ?**

• **A → ? → 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 → B + ?** (unknown side-reaction) and **A → ?** (unknown main product), problems have similar **"unknown unknown"** characteristics.

**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.

**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.



**Thermal, pH, h$v$, O$_2$, other stressors**

**?**

**3-hydroperoxypropanal**

**Idea:** Turn the **A→?** problem into tractable (and parallelizable) **A→B** problems.

**Observations:**

• Product enumeration is easier than transition state enumeration.

• Transition state algorithms for A→B problems are mature. Let the TS algorithm identify physical reactions.

• Recent developments in semi-empirical models and ML create opportunities.

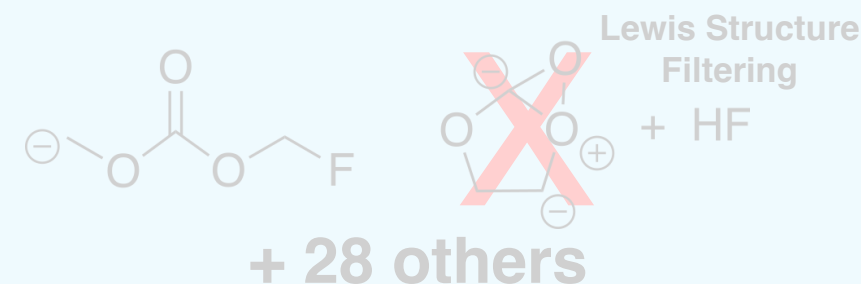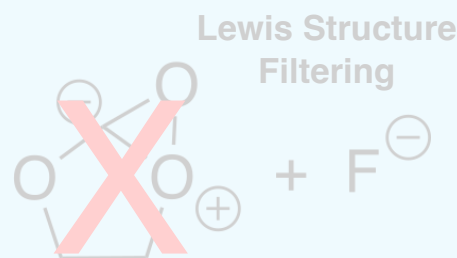• Solving the **A→?** problem is the prerequisite for reaction network prediction.

Polar and pericyclic organic reactions are decomposed into elementary electron donor and acceptor reactions with concomitant σ-bond breaks

**bnfn**
will refer to σ-bond changes, π-bonds are allowed to arbitrarily rearrange.
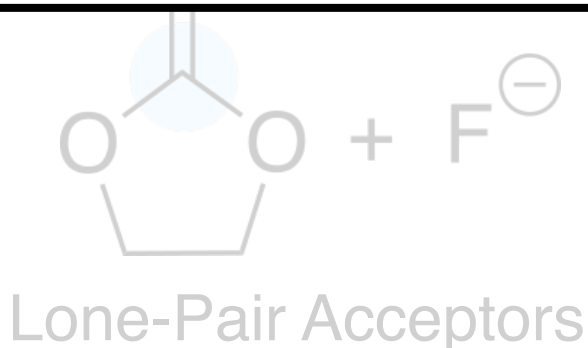


Lone-Pair Donors



Lone-Pair Acceptors

Polar and pericyclic organic reactions are decomposed into elementary electron donor and acceptor reactions with concomitant σ-bond breaks
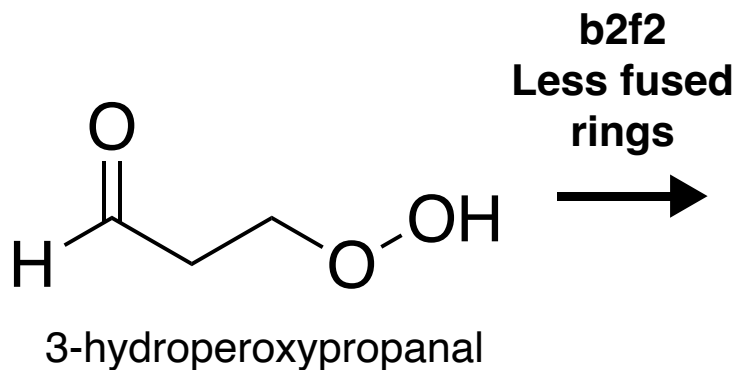
**bnfn**
will refer to σ-bond changes, π-bonds are allowed to arbitrarily rearrange.

Lone-Pair Donors

Lone-Pair Acceptors

**Form 1 Products**

Lewis Structure Filtering

Polar and pericyclic organic reactions are decomposed into elementary electron donor and acceptor reactions with concomitant σ-bond breaks



**bnfn** will refer to σ-bond changes, π-bonds are allowed to arbitrarily rearrange.

Lone-Pair Donors

Lone-Pair Acceptors

**Form 1 Products**

Lewis Structure Filtering

**Break 1 Form 1 Products**

Lewis Structure Filtering

**+ 28 others**

Polar and pericyclic organic reactions are decomposed into elementary electron donor and acceptor reactions with concomitant σ-bond breaks

**Form 1 Products**

**Break 1 Form 1 Products**

**All bnfn products are b(n-1)f(n-1) decomposable**

This means that using only "break 1 bond form 1 bond" (b1f1) for radicals and ions 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)
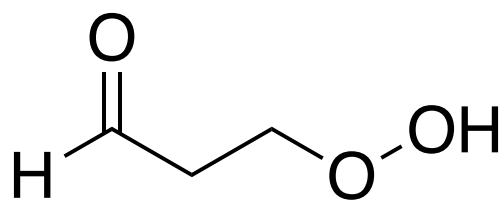
Lone-Pair Acceptors

Lewis Structure Filtering

Lewis Structure Filtering

+ HF

**+ 28 others**

**b2f2
Less fused
rings**

3-hydroperoxypropanal

**3-hydroperoxypropanal**

**b2f2
Less fused
rings**

3-hydroperoxypropanal
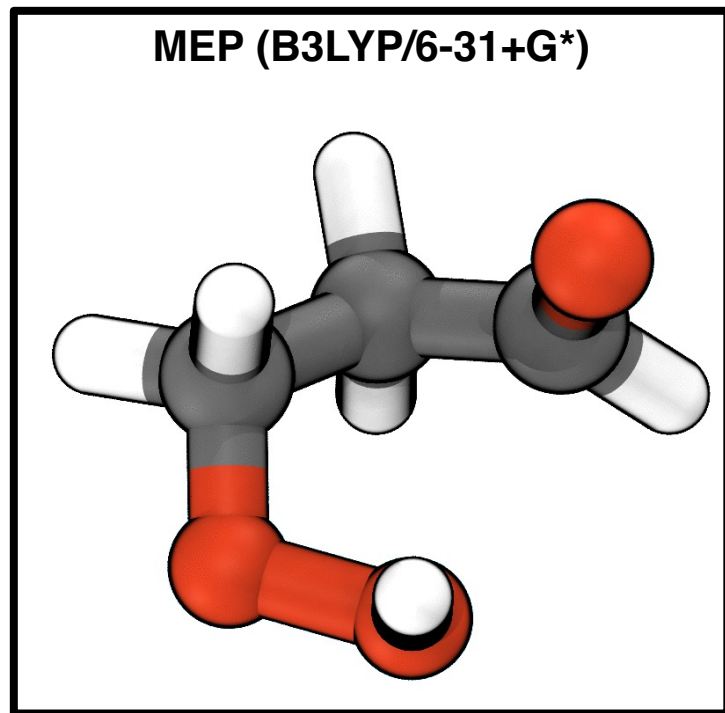
b2f2
Less fused
rings

Filtering
3 and 4
membered
rings
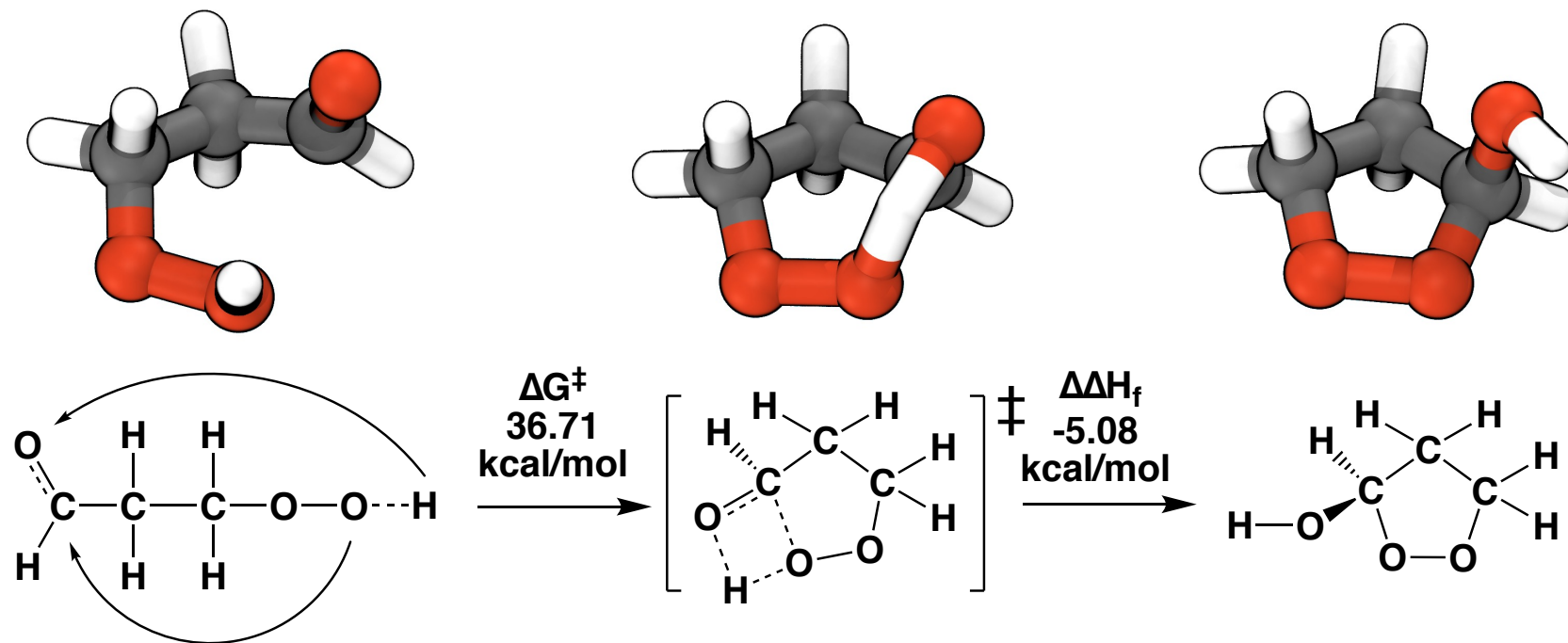
Jensen, R. K.; Korcek, S.; Mahoney, L. R.; Zinbo, M. *JACS* **1979**, 101, 7574

## **The Korcek Mechanism**

According to YARP, this is the lowest barrier unimolecular reaction.

**MEP (B3LYP/6-31+G*)**



$\Delta G^{\ddagger}$
**36.71
kcal/mol**

$\Delta\Delta H_f$
**-5.08
kcal/mol**

Fully resolved (along with subsequent ROOH and R=O formation) 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.
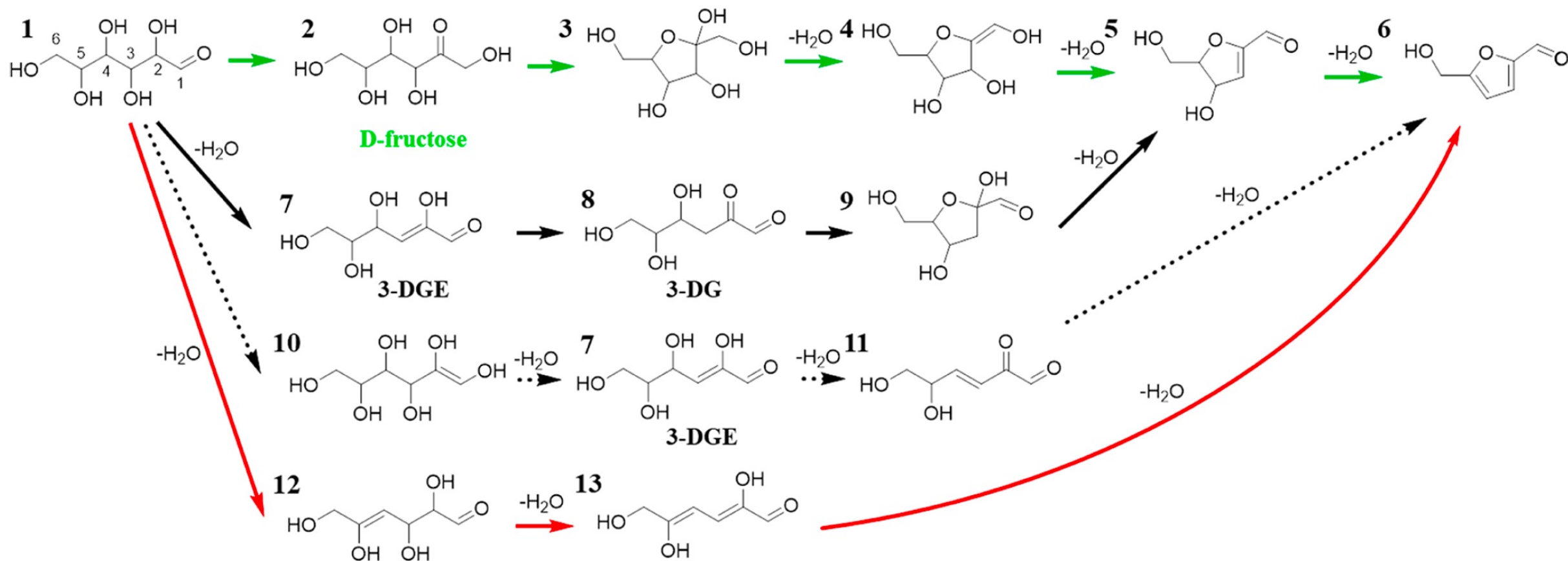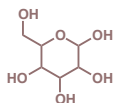
**Figure 1.** Proposed pathways in literature from glucose to HMF, namely the fructose path (green), 3-DG paths (black and black dotted), and direct path (red). The molecules are indicated by numbers and some key molecules are named as follows: **1**. D-glucose; **2**. D-fructose; **3**. D-fructofuranose; **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).

# β-D-Glucose Pyrolysis Network Exploration

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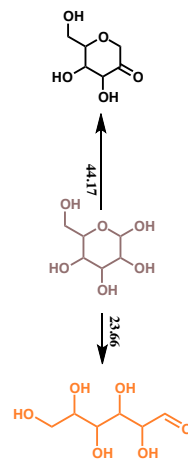
At each iteration:

**(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
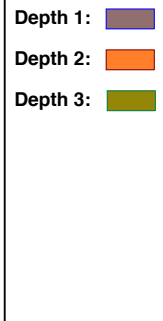
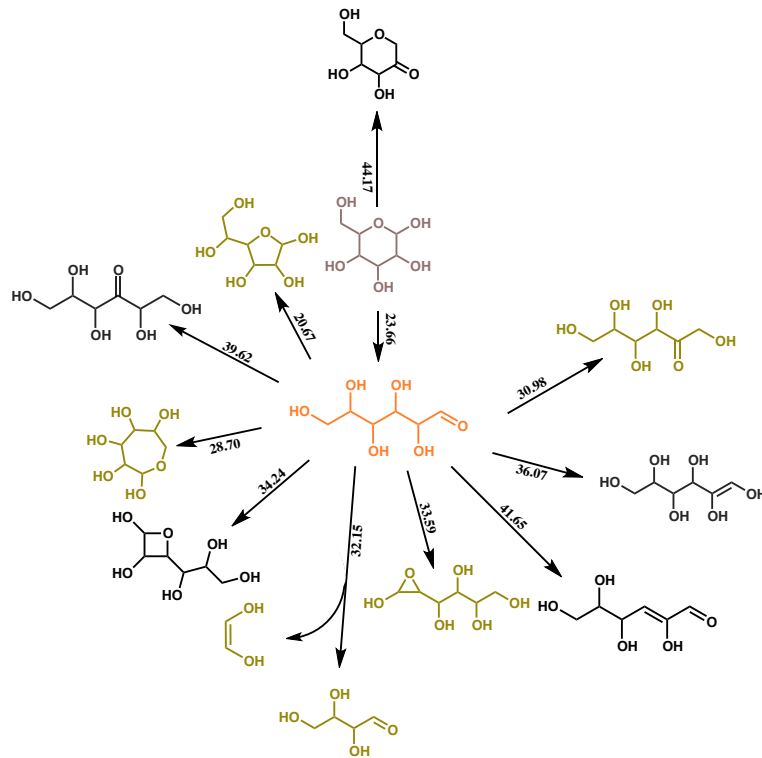**Depth 1:**

44.17

23.66

Depth 1:
Depth 2:

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

At each iteration:

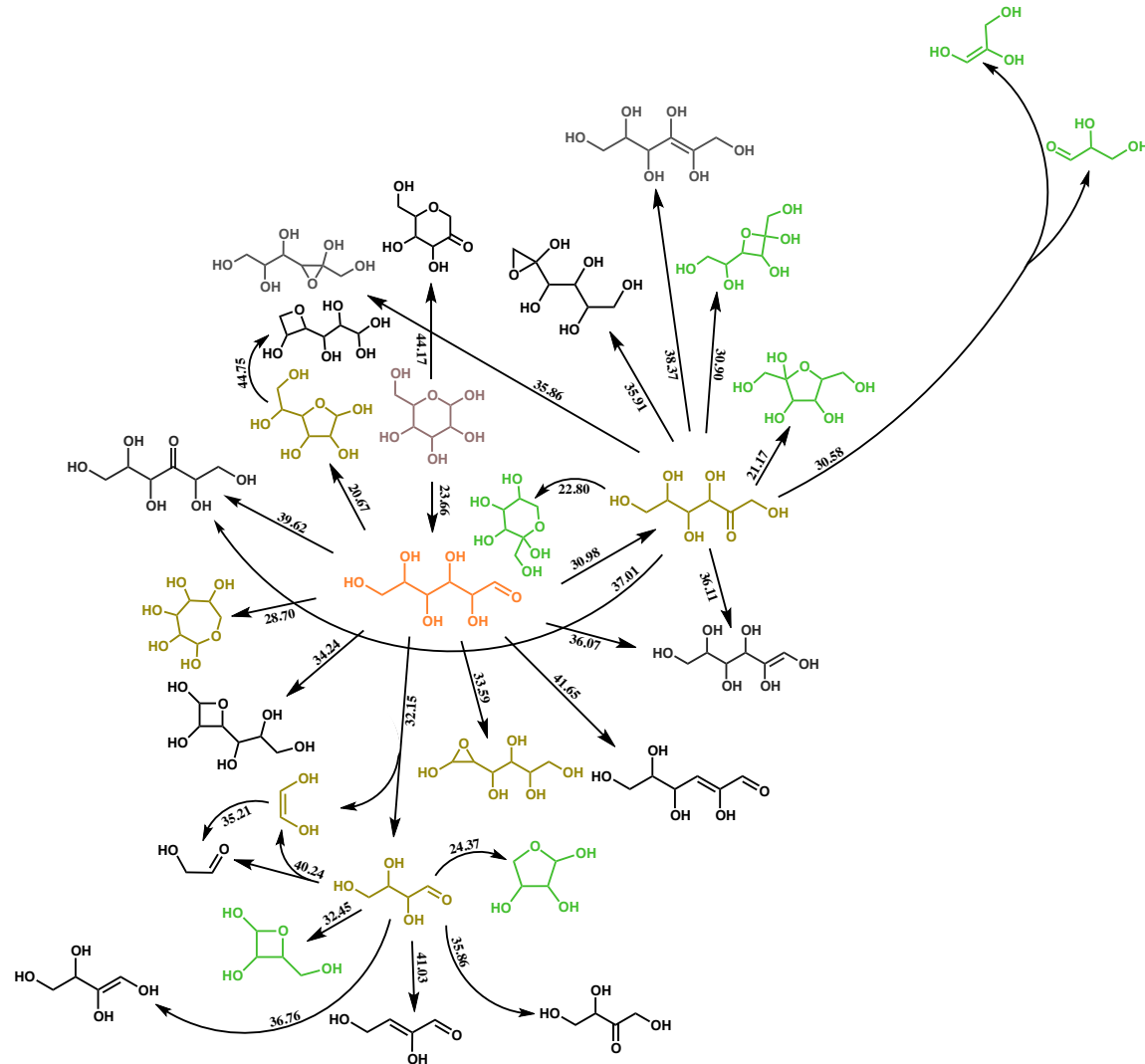**(1)** all b2f2 reactions are explored for active nodes.

**(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

# β-D-Glucose Pyrolysis Network 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
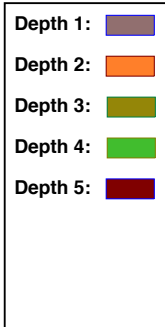
# β-D-Glucose Pyrolysis Network Exploration

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:
Depth 2:
Depth 3:
Depth 4:

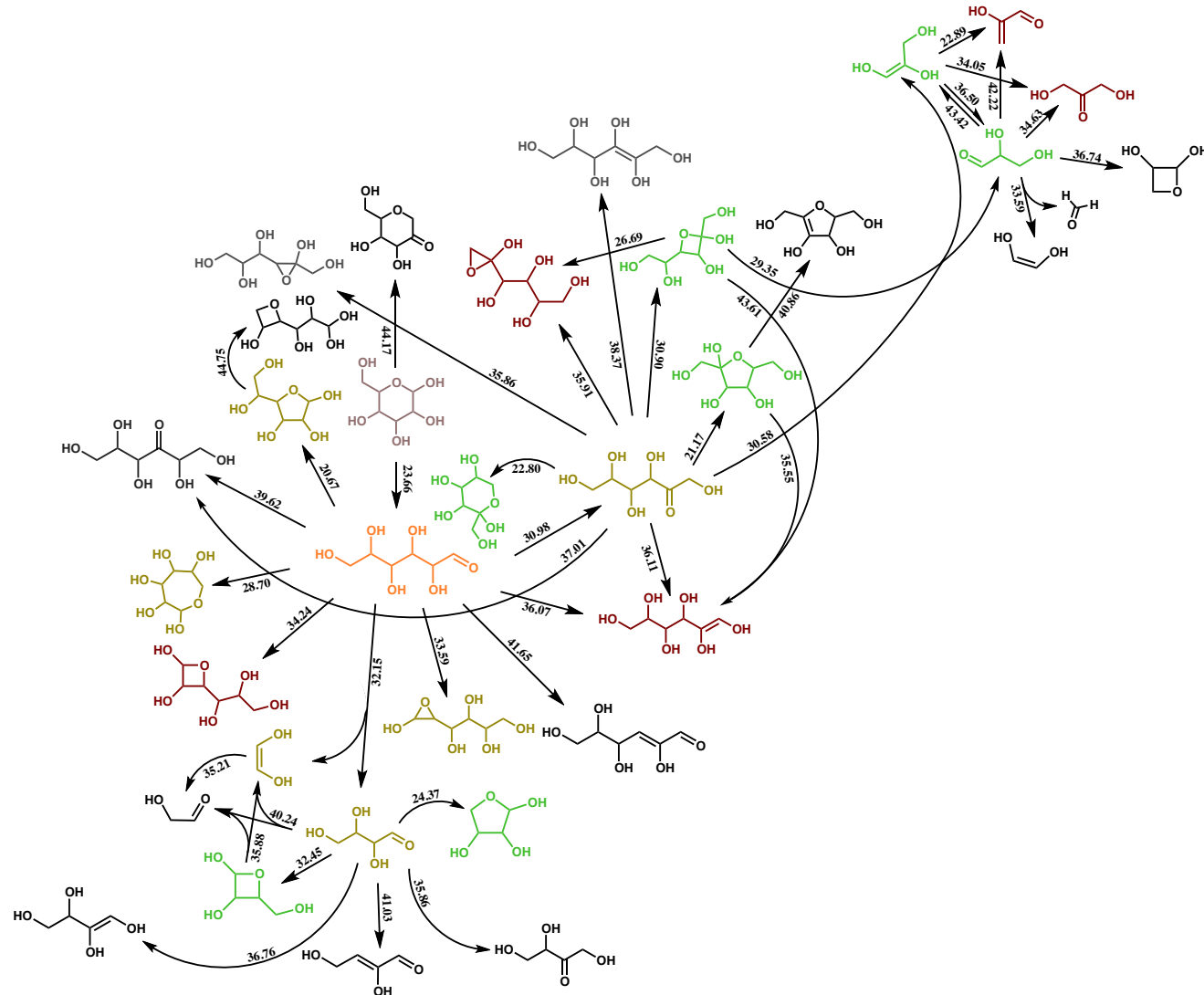# β-D-Glucose Pyrolysis Network Exploration

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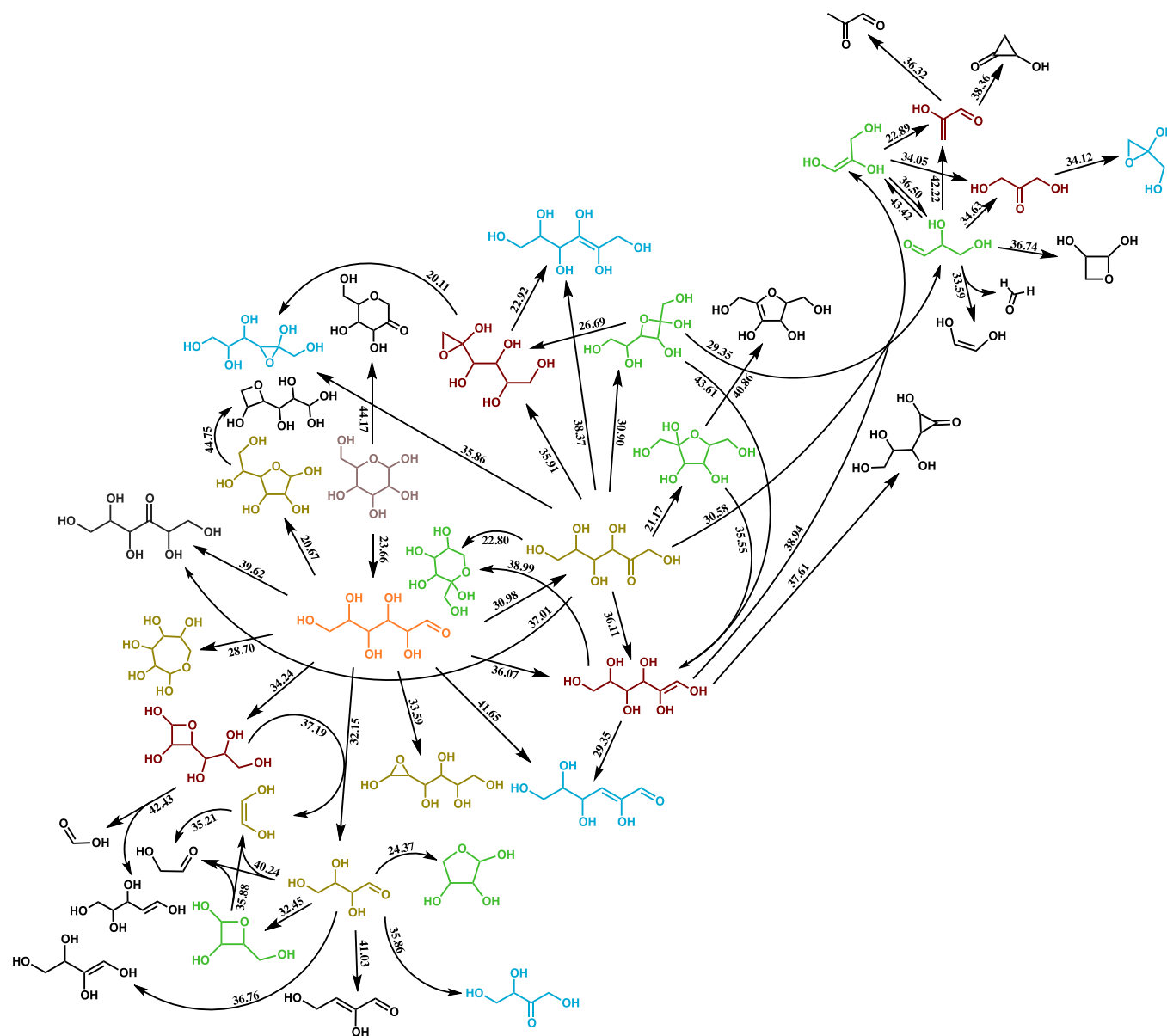
At each iteration:

**(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:
Depth 2:
Depth 3:
Depth 4:
Depth 5:

# β-D-Glucose Pyrolysis Network Exploration

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

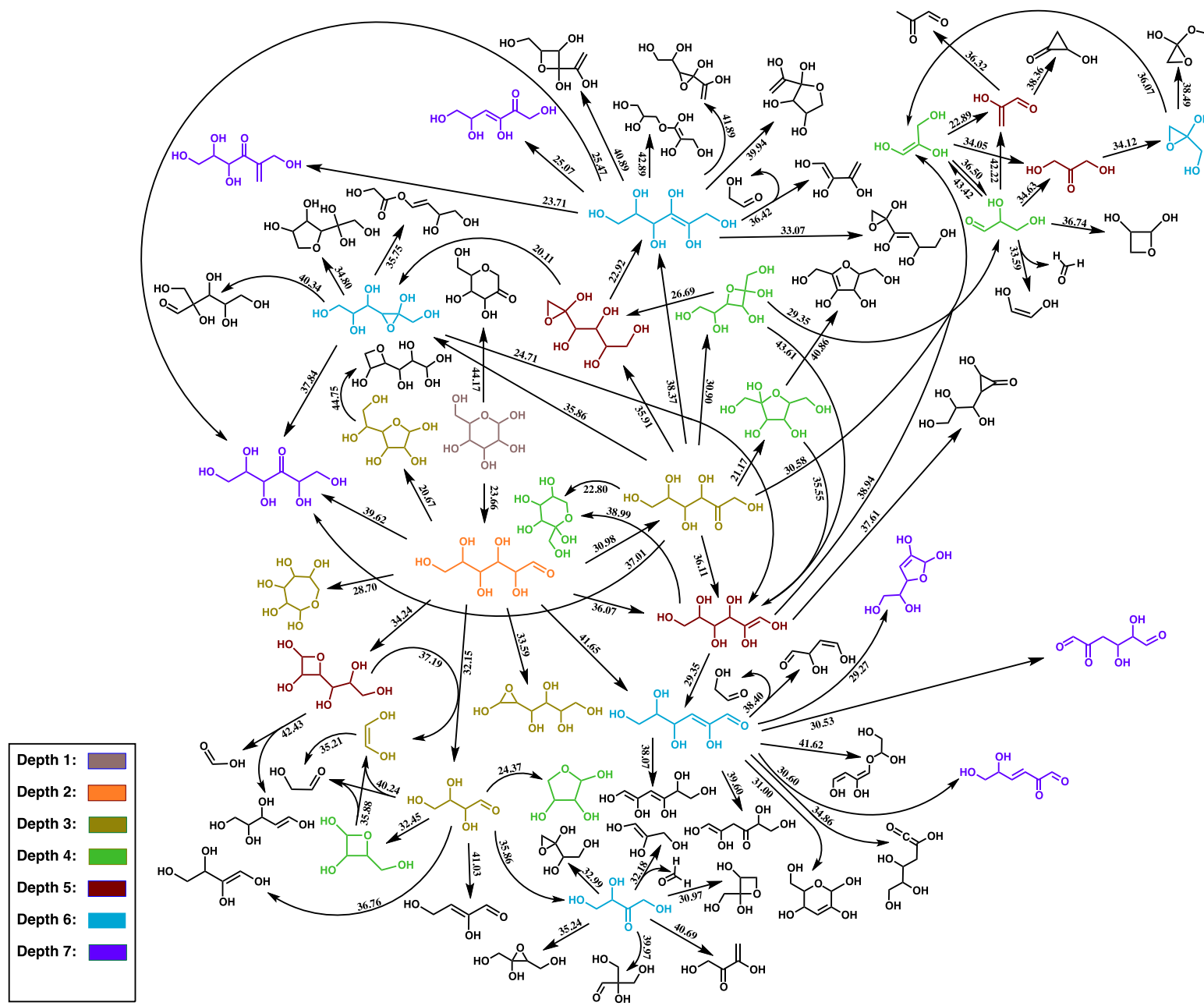At each iteration:

**(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

# β-D-Glucose Pyrolysis Network Exploration

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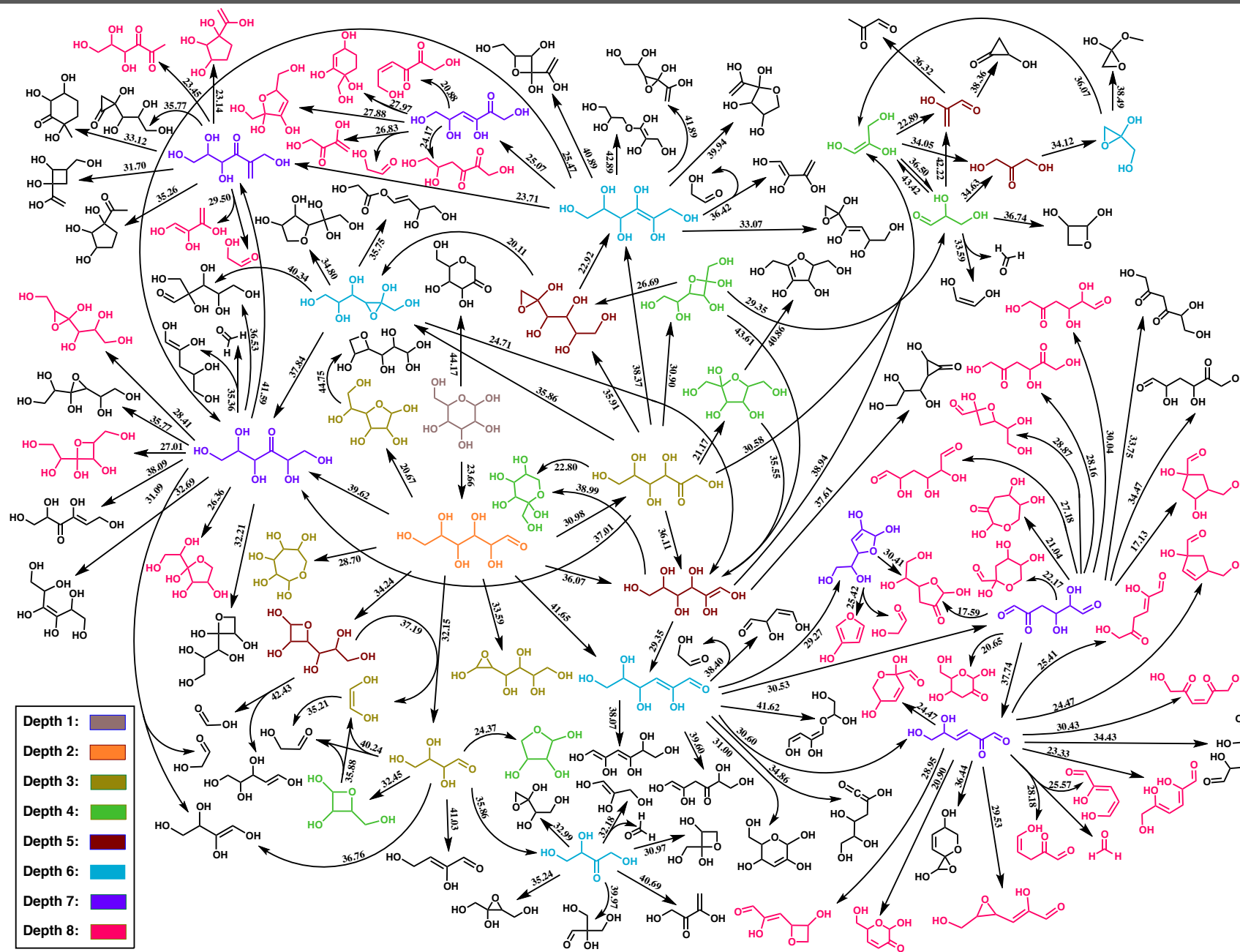
At each iteration:

**(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:
Depth 2:
Depth 3:
Depth 4:
Depth 5:
Depth 6:
Depth 7:

# β-D-Glucose Pyrolysis Network 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.
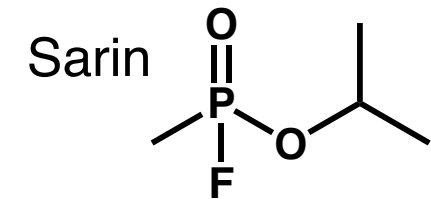
**(2)** Active nodes are determined by the minimum barrier to a given product (with a window)

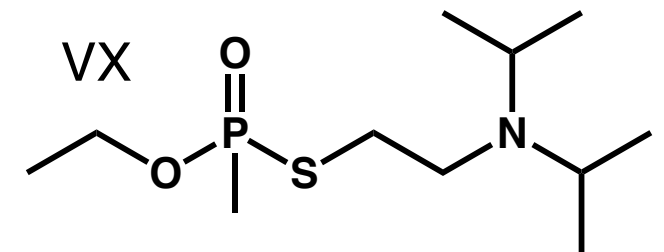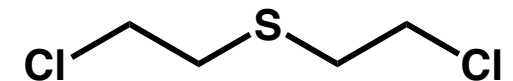**(3)** Water catalyzed reactions are considered for all H-transfers

Depth 1:
Depth 2:
Depth 3:
Depth 4:
Depth 5:
Depth 6:
Depth 7:
Depth 8:

Degradation products are often the only evidence of CWA use or existence. Establishing mechanistic pathways provides evidentiary value to investigators.

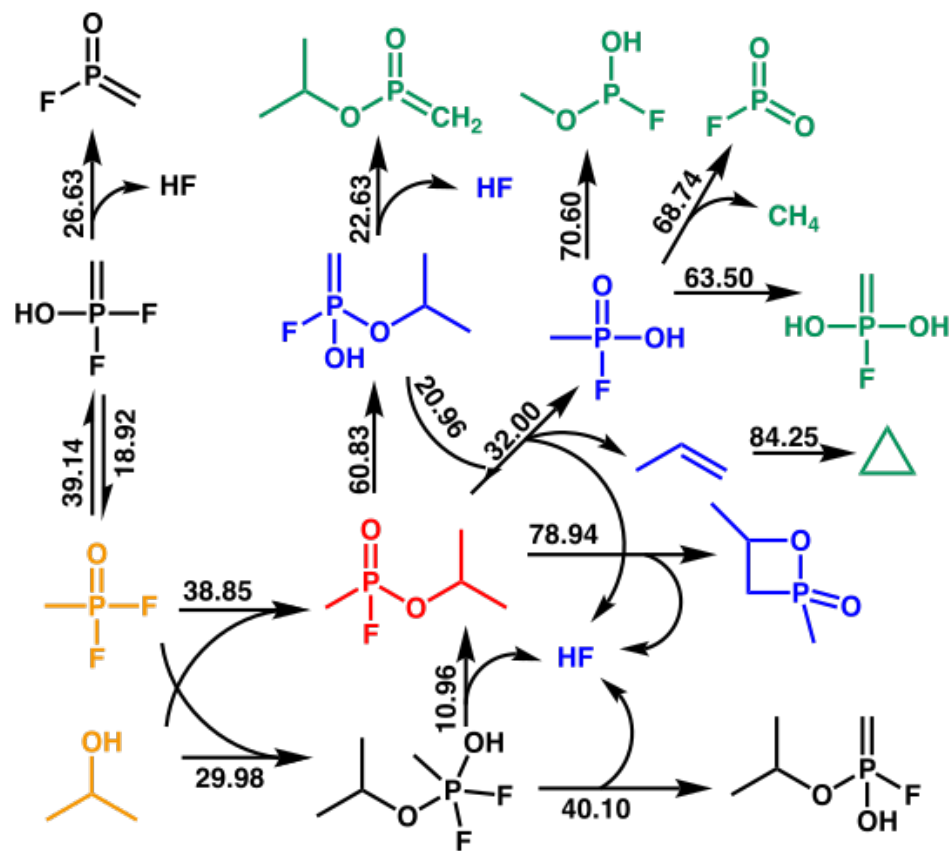| CWA type | Chemical agents | Method of exposure | Clinical symptoms |
|---|---|---|---|
| Nerve agents | G-agents (sarin, cyclosarin, tabun, soman) | Inhalation | SLUDGE, miotic pupils, bradycardia, bronchospasm, bronchorrea, muscle spasms/fasciculations, weakness, flaccid paralysis, tachycardia, seizures, respiratory failure |
| | V-agents (VE, VG, VM, VR, VX) | | |
| Blistering agents | Nitrogen mustard & sulfur mustard (mustard gas) | Inhalation | Acute: Skin, eye and lung damage (pulmonary edema and pulmonary hemmorhage), erythematous rash, skin blistering |
| | | | Chronic: Lung damage (chronic obstructive pulmonary disease, asthma, bronchiolitis obliterans), neutropenia, pancytopenia |
| Asphyxiants | Carbon monoxide, chlorine, phosgene, hydrogen sulfide gases | Inhalation | Upper airway distress, skin and eye irritation, fatal pulmonary edema and acute respiratory distress syndrome |
| Blood agents | Cyanide | Skin absorption, inhalation and ingestion | Severe distress, tachycardia, cyanosis, hypotension, severe metabolic acidosis, seizures, cardiac arrest |
| Hydrofluoric acid | — | Skin absorption, inhalation and ingestion | Severe pain in exposed area, gastrointestinal distress, vomiting, cardiac arrhythmias, hypocalcemia, hyperkalemia |

Sarin

Mustard Gas (HB)

VX

Mechanism of Action

Lowest barrier bimolecular reaction

**Students:** Qiyuan Zhao, Tyler Pasut, Michael Woulfe

**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→?→B, A→B+?, and A→? 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)
• Spencer Goldrich(PMP)

**P²SAC**
Purdue Process Safety & Assurance Center