<u>Predicting ΔH_{f} and Degradation</u> Pathways from First Principles

Brett M. Savoie Davidson Assistant Professor of Chemical Engineering, Purdue University

P2SAC May Meeting, 5/10/22 C2H4

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.



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.

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

Challenges of Contemporary Group Theories



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

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

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

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.



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.



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.



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.

(2-bond specific) **T**opology **A**utomated [8[6[6][1][1]][6[6][6]]] • **Force Field Interactions** Depth 2 graph/structure equivalence 0 1 0 0 0 0 0 0 0 0 1 0 1 0 0 0 0 0 0 0 0 0 0 С 0 1 0 1 0 1 0 0 0 0 0 0 0 1 0 0 0 00 1 0 0 1 0 0 0 0 0 0 0 0 0 0 Adjacency 1 0 0 0 1 0 0 0 0 0 0 matrix for 0 1 0 1 1 1 0 0 0 C 0 0 PEDOT 0 0 1 0 0 0 1 1 1 0 0 0 1 0 0 0 0 0 0 P2SAC 0 monomer 0 0 0 0 1 0 0 0 0 0 0 **Publications** 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 Η 0

0 0 0 0 0 0 1 0 0 0 0

H

TCIT is a <u>component</u> theory

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.



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.





How will we select molecules for parameterizing TCIT components?

identify components*



How will we select molecules for parameterizing TCIT components?













Savoie Research Group

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



Benchmarking ΔH_{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.

200 • TCIT (kJ/mol) -200ΔH_{f,τciτ} (-600 MSE: -0.18 kJ/mol MAE: 2.30 kJ/mol -1000-600 -1000-200200 $\Delta H_{f,G4}$ (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



Benchmarking ΔH_{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 ΔH_{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:

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

 $\mathrm{RC} = H_f(ring) - H_f(\bullet) -$

Benchmarking Ring-Correction Performance

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

(b) The neural-network based ringcorrection exhibits excellent reproduction of the G4 predictions (MSE: ~3kJ/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, <u>these compounds are</u> within BGIT's training data. ~120 ring-containing compounds from PNK (excluding training)

BGIT cannot make predictions for ~2% of PNK compounds

How Many Components are Possible?

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?

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

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 2000 H,C,N, and O containing molecules for Components per 1k Molecules distinct components and the model 1500 compounds necessary to predict ΔH_{f} 1000 The derivative plot shows that TCIT New 500 initially generates ~2 new components per molecule, but by the end of the 0 sampling ~100 molecules need to be 200k 600k 800k 1m 400k 0 sampled to find a new component. 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 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

PubChem is a repository of chemical properties that contains many millions of organic By the end of the summer, we will have 60k model compounds and 2k distinct rings calculated at the G4 level. This will be the largest repository of G4 calculations in the world. 3000 st abarra that It is foreseeable that we could complete all B, F, Cl, Br, S, and P containing structures over the next year.

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

 $A \rightarrow ?$: 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

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 quantum chemistry and ML make it worthwhile to revisit this problem.

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.

Savoie Research Group I

Zhao, Q.; Savoie, B. M. Nature Computational Science 2021, 479-490. (P2SAC Supported)

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.

Zhao, Q.; Savoie, B. M. *Nature Computational Science* 2021, 479-490. (P2SAC Support)

3-Hydroperoxypropanal - Reaction Network

We used YARP to recursively elucidate the 3hydroperoxypropanal 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).

Zhao, Q.; Savoie, B. M. Nature Computational Science 2021, 479-490. (P2SAC Support)

3-Hydroperoxypropanal - Reaction Network

We used YARP to recursively elucidate the 3hydroperoxypropanal 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).

Zhao, Q.; Savoie, B. M. *Nature Computational Science* 2021, 479-490. (P2SAC Support)

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.

Savoie Research Group

Network Exploration 2: Glucose Pyrolysis

Figure 1. Proposed pathways in literature from gluclose 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).

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)

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)

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)

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)

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)

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)

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)

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)

Outlook

The throughput enabled by YARP creates many new opportunities:

- (i) Broader reaction discovery \rightarrow Lowe Dataset as a discovery testbed
- (ii) Generating positive and negative exemplary reaction datasets
- (iii) Exploring deeper networks (e.g., materials degradation, catalysis)

In the process safety space, it seems like predicting kinetics in addition to thermochemistry could be useful.

Qiyuan, Z.; Savoie, B. M. "Simultaneously Improving Reaction Coverage and Computational Cost in Automated Reaction Prediction Tasks." Nature Computational Science **2021**, 1, 479-490.

Outlook and Acknowledgements

Students: Qiyuan Zhao, Tyler Pasut

Project Accomplishments:

- Implemented a fully-consistent 2-bond (i.e., component) increment theory based on G4 data
- Automated model compound generation and fitting algorithms.
- Built a database infrastructure for reusing calculations and parameter fitting.
- Developed a ring-correction for TCIT to improve performance on conjugated and non-benzene structures.
- Extended TCIT to condensed phases and new thermodynamic properties and radicals.

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

