## High-Throughput Predictions of Molecular Thermodynamics and Reactivity



Brett M. Savoie Charles Davidson Assistant Professor of Chemical Engineering 5/17/21, P2SAC Spring Meeting (Day 3)

### **Challenges of Contemporary Group Theories**

### **Benson Groups:**

• The idea is to decompose molecular properties ( $\Delta H_f$ , S°, C<sub>v</sub>) 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:**

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

• **Specificity:** the definition of a "group" has never been formalized and inconsistent granularity is applied. Limited information on group interactions.

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



Experimental: -5.15 +/- 0.34 kcal/mol

### **Challenges of Contemporary Group Theories**

### **Benson Groups:**

- The idea is to decompose molecular properties ( $\Delta H_f$ , S°, C<sub>v</sub>) 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:**

- **Provenance:** inconsistent thermodynamic data is available/used to determine group contributions.
- **Specificity:** the definition of a "group" has never been formalized and inconsistent granularity is applied. Limited information on group interactions.
- 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

## **Challenges of Contemporary Group Theories**



• **Specificity:** the definition of a "group" has never been formalized and inconsistent granularity is applied. Limited information on group interactions.

• 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.* In Press.

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." **2021** ChemRxiv. Preprint. https://doi.org/10.26434/chemrxiv.14527299.v1



Adjacency matrix for pedot monomer

## 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.* In Press.

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." **2021** ChemRxiv. Preprint. https://doi.org/10.26434/chemrxiv.14527299.v1

### TAFFI syntax for machinereadable atom types



## 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.* In Press.

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." **2021** ChemRxiv. Preprint. https://doi.org/10.26434/chemrxiv.14527299.v1

### TAFFI syntax for machinereadable atom types



## 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.* In Press.

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." **2021** ChemRxiv. Preprint. https://doi.org/10.26434/chemrxiv.14527299.v1

#### readable atom types **T**opology **A**utomated Force Field Interactions [8[6[6][1][1]][6[6][6]]] • 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 0 1 0 1 0 1 0 0 0 0 0 0 0 1 0 1 0 0 0 0 0 1 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 1 0 1 1 1 0 0 0 0 0 0 0 1 0 0 0 1 1 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 0 0 0 1 0 0 0 1 0 0 0 0 0 0 0 0 0 0 1 0 0 0 0 н 0 0 0 0 0 0 1 0 0 0 0 Adjacency matrix for pedot monomer

TAFFI syntax for machine-

## 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.* In Press.

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." **2021** ChemRxiv. Preprint. https://doi.org/10.26434/chemrxiv.14527299.v1



### 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. But 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.* In Press.

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." **2021** ChemRxiv. Preprint. https://doi.org/10.26434/chemrxiv.14527299.v1

### TAFFI syntax for machinereadable atom types



														4
С	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	
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	
0	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	
Η	0	0	0	0	0	0	1	0	0	0	0	0	0	
Н	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	
Η	0	0	0	0	0	0	0	1	0	0	0	0	0	
Н	0	0	0	0	0	0	0	1	0	0	0	0	0	
Adia	r ace	ene	cv	m	at	rix	fc	or	pe	dc	ot r	nc	no	omer

### **Graphical Decomposition of Model Compounds**



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

### **Graphical Decomposition of Model Compounds**



 $\Delta H_{f,G4} = -259.9 \text{ kJ/mol}$  $\Delta H_{f,TGIT} = -259.3 \text{ kJ/mol}$ no experimental data

Our implementation automatically plans calculations accounting for dependencies.

All calculations and components are databased for future use

## Benchmarking ΔH<sub>f,gas</sub> Predictions Against the PNK Dataset

Initial benchmarking set consists of ~1100 linear
C,H, and O containing compounds from PNK<sup>1</sup>

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2<sup>nd</sup> 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

572 small compounds from PNK

## Benchmarking ΔH<sub>f,gas</sub> Predictions Against the PNK Dataset

## Initial benchmarking set consists of ~1100 linear C,H, and O containing compounds from PNK<sup>1</sup>

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2<sup>nd</sup> 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<sub>f, τcir</sub> ( -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

150 medium compounds from PNK

## Benchmarking ΔH<sub>f,gas</sub> Predictions Against the PNK Dataset

## Initial benchmarking set consists of ~1100 linear C,H, and O containing compounds from PNK<sup>1</sup>

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2<sup>nd</sup> 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

~500 large compounds from PNK

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

## Initial benchmarking set consists of ~1100 linear C,H, and O containing compounds from PNK<sup>1</sup>

(1) J. B. Pedley, R. D. Naylor, S. P. Kirby "Thermochemical Data of Organic Compounds" 2<sup>nd</sup> 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.

~500 large compounds from PNK

### **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:



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

Zhao, Q.; Iovanac, N.; Savoie, B. M.; "Transferable Ring Corrections for Predicting Enthalpy of Formation of Cyclic Compounds" J. Chem. Info. Model. In Press.

### **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, these compounds are within BGIT's training data. ~120 ring-containing compounds from PNK (excluding training)



BGIT cannot make predictions\_for ~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)

### 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<sup>[1]</sup>:

$$\Delta_{\rm f} H^{\circ}_{\rm (s)} = \Delta_{\rm f} H^{\circ}_{\rm (g)} - \Delta_{\rm sub} H^{\circ}$$
$$\Delta_{\rm f} H^{\circ}_{(\ell)} = \Delta_{\rm f} H^{\circ}_{\rm (g)} - \Delta_{\rm vap} H^{\circ}$$

We have implemented group contribution models for heat of vaporization<sup>[2]</sup> and sublimation<sup>[3]</sup>, 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**<sub>v</sub>): 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<sub>rot</sub>) and molecular symmetry:

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

Murray, J.S., Brinck, T. and Politzer, P., **1996**. *Chemical physics*, 204, 289-299.
Pankow, J.F. and Asher, W.E., **2008**. *Atmospheric Chemistry and Physics*.
Bagheri, M.; Gandomi, A. H.; Golbraikh, A. **2012**, *Thermochim. Acta*, 543, 96–106

- <.> Indicates conformational averaging
- R: ideal gas constant
- $\sigma$ : symmetry number
- $\alpha$ ,  $\beta$ : regressed constants

### Benchmarking Condensed Phase ΔH<sub>f</sub> 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  $\Delta H_{vap}$  and  $\Delta H_{sub}$  models.

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

### **Benchmarking TCIT S° and C<sub>v</sub> 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<sup>[1]</sup>

(d) TCIT heat capacity comparison for 904 large molecules from NIST<sup>[1]</sup>

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

TCIT now supports S° and C $_{\rm v}$  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<sub>v</sub> Predictions**



### What's Next: Ions, Radicals, and other Properties

Ionic and radical species present acute challenges for experimental data collection and have limited representability within BGIT. We propose to continue leveraging the intrinsic advantages of TCIT to address these gaps.

We are extending TCIT to cover ionic and radical species by propagating formal charge information through model compound generation and component definitions:



Because of the way that we have developed TCIT, the extension is relatively straightforward, it just requires the additional prediction of hydrogen-bond increments for conversion from the radical to the closed-shell analogue.

### Preliminary TCIT Radical Benchmarks (ΔH<sub>f</sub>)



• In this case, there simply isn't much experimental data available, so most of our validations are occurring at the G4 level (some experimental comparisons are presented later).

•  $\Delta H_f$  errors are consistent with previous benchmark for closed-shell species.

### Preliminary TCIT Radical Benchmarks (BDE)



• The bond-dissociation energy (BDE) is another test that is relevant to radicals.

• Errors are consistent with error propagation (Note: BDE is calculated as the enthalpy difference between the neutral and two radicals generated by the scission)

### Preliminary TCIT Radical Benchmarks (C<sub>v</sub>)



• Similar accuracy compared with closed-shell species.

### Preliminary TCIT Radical Benchmarks (S°)



• Similar accuracy compared with closed-shell species.

### Preliminary TCIT Radical Benchmarks (ΔH<sub>f</sub> exp)

### ~40 experimental values have been cobbled together from NIST and ATCT testing structures

- Preliminary validation is looking promising, with similar overall performance for TCIT on radical and closed-shell species.
- This is a big victory for TCIT, since the limited experimental data for radical species makes CHETAH predictions impossible in many cases.



### **Trying to Solve 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 reaction enumeration.
- Transition state algorithms for A→B problems are mature. Let the TS algorithm identify physical reactions
- Recent developments in semiempirical quantum chemistry can be leveraged here.



Qiyuan, Z.; Savoie, B. M. "More and Faster: Simultaneously Improving Reaction Coverage and Computational Cost in Automated Reaction Prediction Tasks." (2020) *ChemRxiv*. https://doi.org/10.26434/chemrxiv.13076087.v1

### **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.



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

### **YARP: Elementary Reaction Step(s)**

**Elementary Reaction Step:** a reaction with a single saddle point between products and reactants (we can't always know this in advance).

For closed-shell neutral organic molecules, "break 2 bonds form 2 bonds" is the simplest ERS that yields non-trivial closed-shell neutral products.







**Note:** b3f3 and b4f4 products are decomposable into b2f2 ERS(s).

A more general suite of ERS(s) are certainly possible. These will certainly be required for radicals, organics with expanded octets, and metals.

### **Testing YARP on a Network 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. *JACS* **2018**, 140 (3), 1035–1048.









### **Testing YARP on a Network Problem**



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

Grambow, C. A, Suleimanov, Y. V. JACS 2018, 140 (3), 1035-1048.

### **3-Hydroperoxypropanal - Reaction Network**

We used YARP to recursively elucidate all b2f2 products out to depth two, in order to achieve parity with Grambow et al.



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

### 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.

## 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)

### YARP needs to be extended:

- (i) Broader suite of ERS(s) for organometallics and expanded octets
- (ii) Validation of our treatment of ionic, radical, and resonance species

(iii) Kinetic modeling for prioritizing branches of deep networks

Qiyuan, Z.; Savoie, B. M. "More and Faster: Simultaneously Improving Reaction Coverage and Computational Cost in Automated Reaction Prediction Tasks." (2020) *ChemRxiv*. https://doi.org/10.26434/chemrxiv.13076087.v1

### **Outlook and Acknowledgements**

# **Qiyuan Zhao** performed all of the work. **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.



- P2SAC for funding.
- Ray Mentzer (Purdue)



- Madison Sprecher (Purdue UG)
- Caitlin Justice (Purdue UG)