

# Generic Group Contribution Method

*Authors:* Christoph Flamm<sup>5,8</sup>, Marc Hellmuth<sup>1,7</sup>, Daniel Merkle<sup>2</sup>, Nikolai Nøjgaard<sup>1,2</sup>, Peter F. Stadler<sup>3,4,5,7</sup>

---

February 11, 2019

<sup>1</sup> Dpt. of Mathematics and Computer Science, University of Greifswald

<sup>2</sup> Department of Mathematics and Computer Science, University of Southern Denmark, Denmark

<sup>3</sup> Bioinformatics Group, Department of Computer Science; and Interdisciplinary Center of Bioinformatics, University of Leipzig

<sup>4</sup> Max-Planck-Institute for Mathematics in the Sciences

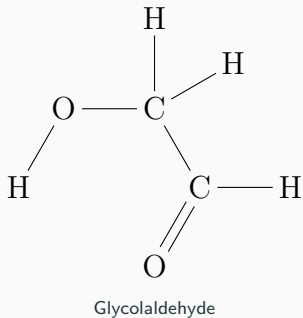
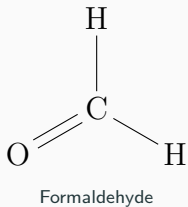
<sup>5</sup> Inst. f. Theoretical Chemistry, University of Vienna

<sup>6</sup> Santa Fe Institute, 1399 Hyde Park Rd., Santa Fe

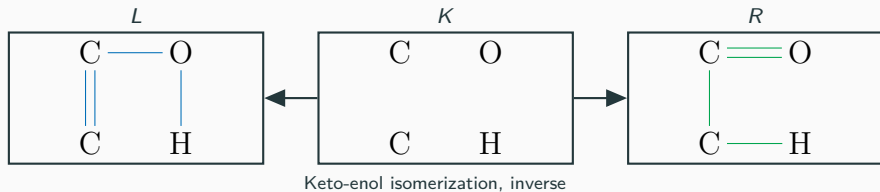
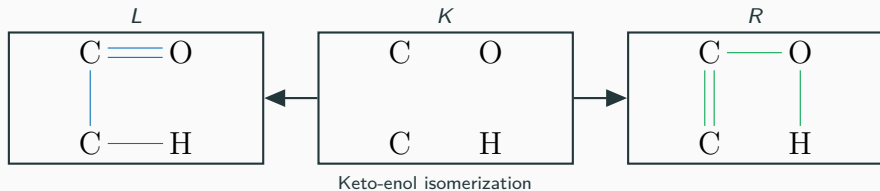
<sup>7</sup> Saarland University, Center for Bioinformatics

<sup>8</sup> Center for Anatomy and Cell Biology, Medical University of Vienna

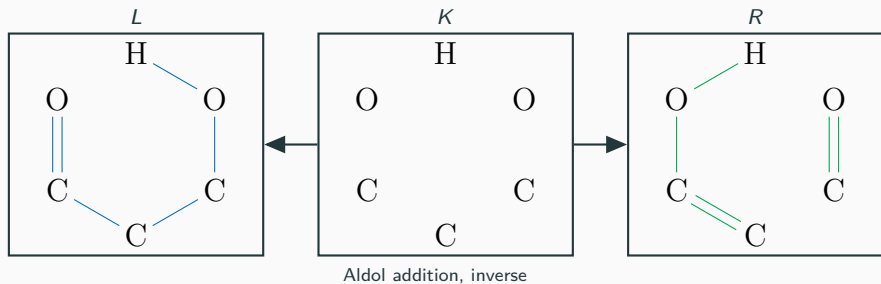
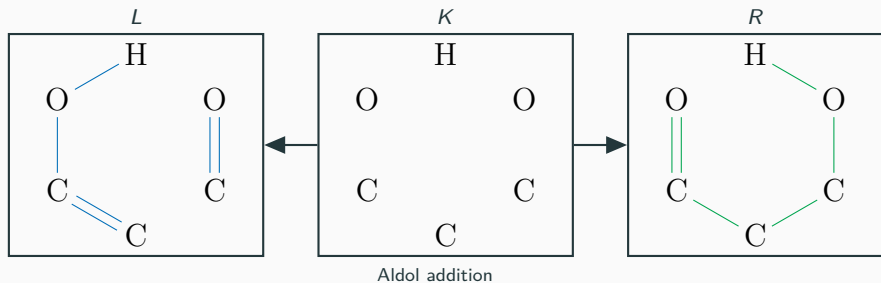
## The Beginning: A look at MØD



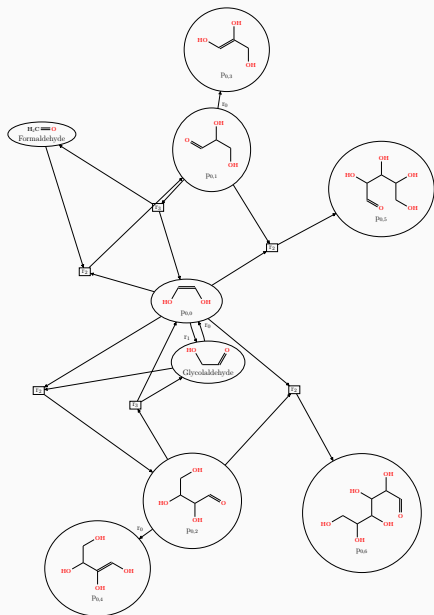
# The Beginning: A look at MØD



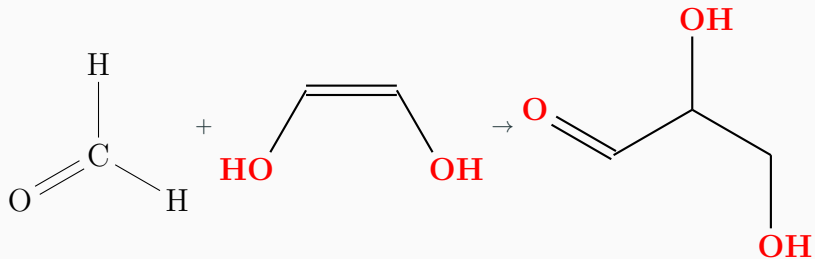
# The Beginning: A look at MØD



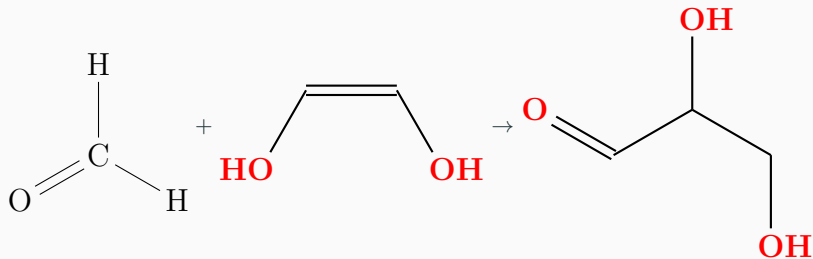
# The Beginning: A look at MØD



## Reactions: Fact or Fiction?

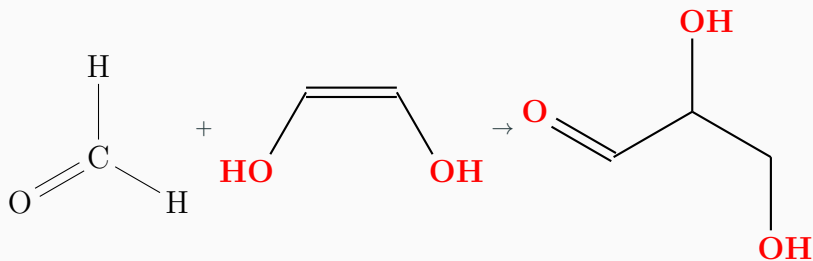


## Reactions: Fact or Fiction?



Real or theoretical reaction? Let's look at its energy change!

## Reactions: Fact or Fiction?



Real or theoretical reaction? Let's look at its energy change!

Gibbs Free Energy:  $G = H - TS$ , where  $H$  is enthalpy,  $T$  temperature and  $S$  entropy.

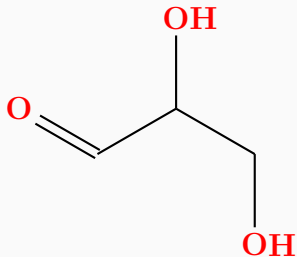
Gibbs Free Energy Change:  $\Delta G = G_{products} - G_{educts}$



## So How Do We Compute the Energy

- The Gibbs Free Energy of a molecule can be measured in the lab.
- But our chemical universe can (in theory) be infinite.
- Hence, we want to create a predictive model on a sampled population.

# State of the Art: Group Contribution Method

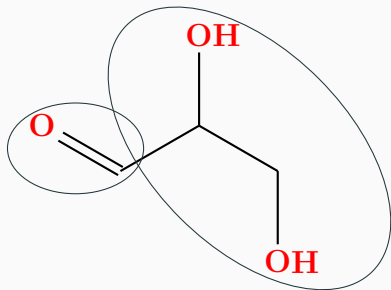


- We can decompose a molecule into functional groups that linearly relates to  $G$ .

Problems with the Group Contribution Method in a Generic Framework:

- What are the functional groups?
- How to tile a graph?
- Introducing new functional group changes the entire input.

## State of the Art: Group Contribution Method

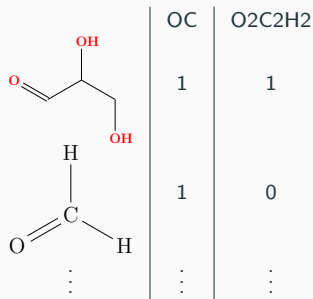
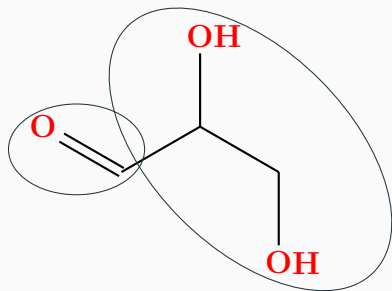


- We can decompose a molecule into functional groups that linearly relates to  $G$ .

Problems with the Group Contribution Method in a Generic Framework:

- What are the functional groups?
- How to tile a graph?
- Introducing new functional group changes the entire input.

# State of the Art: Group Contribution Method

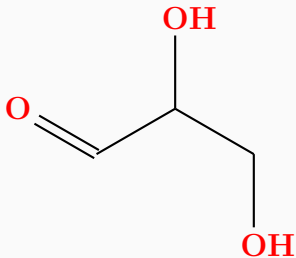


- We can decompose a molecule into functional groups that linearly relates to  $G$ .

Problems with the Group Contribution Method in a Generic Framework:

- What are the functional groups?
- How to tile a graph?
- Introducing new functional group changes the entire input.

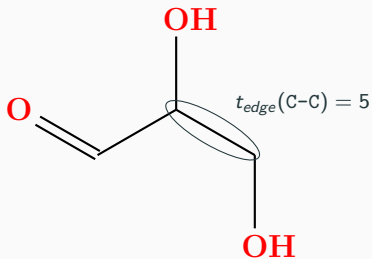
## A Closer Look at Molecular Energies



- The energy of a molecule can be approximated as the sum of its bond energies.

$$t_{obs}(G) = \sum_{e \in E(G)} t_{edge}(e)$$

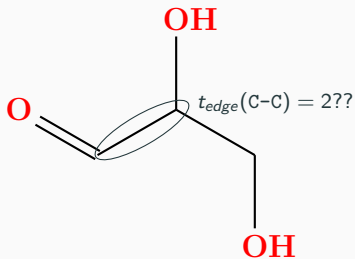
## A Closer Look at Molecular Energies



- The energy of a molecule can be approximated as the sum of its bond energies.

$$t_{obs}(G) = \sum_{e \in E(G)} t_{edge}(e)$$

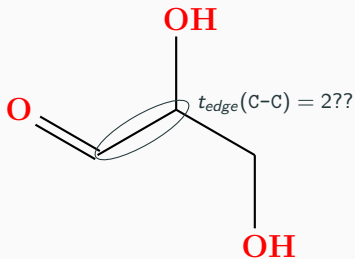
## A Closer Look at Molecular Energies



- The energy of a molecule can be approximated as the sum of its bond energies.

$$t_{obs}(G) = \sum_{e \in E(G)} t_{edge}(e)$$

## A Closer Look at Molecular Energies



- The energy of a molecule can be approximated as the sum of its bond energies.

$$t_{obs}(G) = \sum_{e \in E(G)} t_{edge}(e)$$

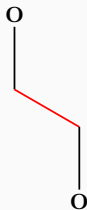
- The bond energy is determined by its surrounding **context**.



# Defining Contexts

## Definition (Context)

A **context** is a pair  $C = (G, e)$ , where  $G$  is a graph and  $e$  is an edge in  $G$ . The size of  $C$  is defined as the number of edges in  $G$  and we call  $e$  the origin edge.



# Defining Contexts

## Definition (Context)

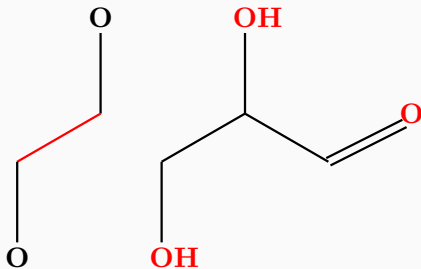
A **context** is a pair  $C = (G, e)$ , where  $G$  is a graph and  $e$  is an edge in  $G$ . The size of  $C$  is defined as the number of edges in  $G$  and we call  $e$  the origin edge.

## Definition (Frequency)

Given a graph  $G$  and a context  $C = (H, e')$  we say that  $C$  is a *context around*  $e \in E(G)$ , if there is a subgraph isomorphism  $\varphi$  from  $H$  to  $G$  that satisfy  $\varphi(e') = e$ . The **frequency**  $f(C, G, e)$  of  $C$  around some edge  $e \in E(G)$  is the number of subgraph isomorphisms  $\varphi_1, \varphi_2, \dots$  from  $C$  to  $G$  that satisfy  $\varphi_i(e') = e$ .

The *frequency* of  $C$  in  $G$  is defined as:

$$f(C, G) = \sum_{e \in E(G)} f(C, G, e)$$



# Defining Contexts

## Definition (Context)

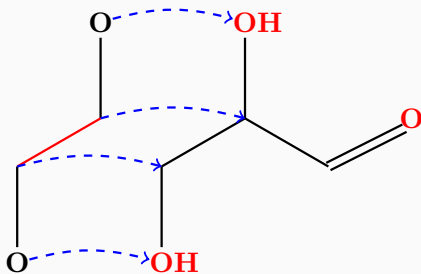
A **context** is a pair  $C = (G, e)$ , where  $G$  is a graph and  $e$  is an edge in  $G$ . The size of  $C$  is defined as the number of edges in  $G$  and we call  $e$  the origin edge.

## Definition (Frequency)

Given a graph  $G$  and a context  $C = (H, e')$  we say that  $C$  is a *context around*  $e \in E(G)$ , if there is a subgraph isomorphism  $\varphi$  from  $H$  to  $G$  that satisfy  $\varphi(e') = e$ . The **frequency**  $f(C, G, e)$  of  $C$  around some edge  $e \in E(G)$  is the number of subgraph isomorphisms  $\varphi_1, \varphi_2, \dots$  from  $C$  to  $G$  that satisfy  $\varphi_i(e') = e$ .

The *frequency* of  $C$  in  $G$  is defined as:

$$f(C, G) = \sum_{e \in E(G)} f(C, G, e)$$



# Defining Contexts

## Definition (Context)

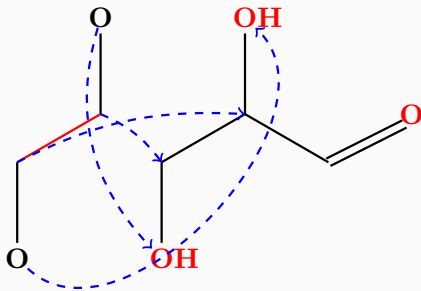
A **context** is a pair  $C = (G, e)$ , where  $G$  is a graph and  $e$  is an edge in  $G$ . The size of  $C$  is defined as the number of edges in  $G$  and we call  $e$  the origin edge.

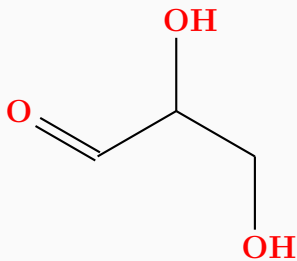
## Definition (Frequency)

Given a graph  $G$  and a context  $C = (H, e')$  we say that  $C$  is a *context around*  $e \in E(G)$ , if there is a subgraph isomorphism  $\varphi$  from  $H$  to  $G$  that satisfy  $\varphi(e') = e$ . The **frequency**  $f(C, G, e)$  of  $C$  around some edge  $e \in E(G)$  is the number of subgraph isomorphisms  $\varphi_1, \varphi_2, \dots$  from  $C$  to  $G$  that satisfy  $\varphi_i(e') = e$ .

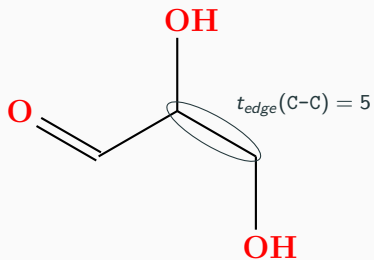
The *frequency* of  $C$  in  $G$  is defined as:

$$f(C, G) = \sum_{e \in E(G)} f(C, G, e)$$

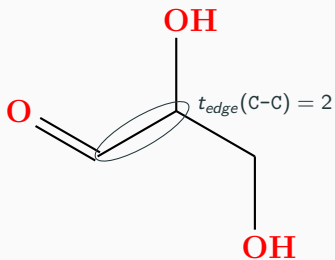




$$\mathcal{K}_1 = \{O=C, C-C, O-C, O-H\}$$



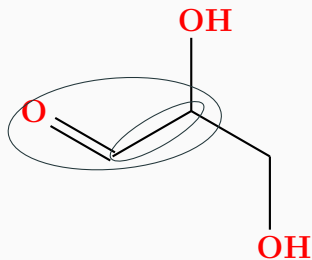
$$\mathcal{K}_1 = \{O=C, C-C, O-C, O-H\}$$



$$\mathcal{K}_1 = \{O=C, C-C, O-C, O-H\}$$

$$t_C(C-C) = \text{avg. energy} = 3.5$$

$$t_{edge}(e) \approx f(C, G, e) \cdot t_C(C-C) = 3.5$$

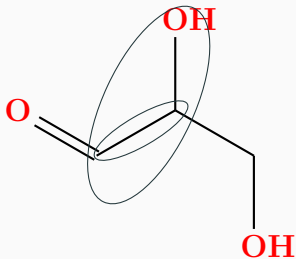


$$\mathcal{K}_1 = \{O=C, C-C, O-C, O-H\}$$

$$t_C(C-C) = \text{avg. energy} = 3.5$$

$$t_{edge}(e) \approx f(C-C, G, e) \cdot t_C(C-C) + f(G, O=C-C) \cdot t_C(O=C-C)$$



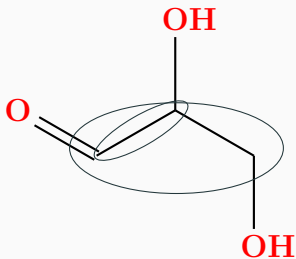


$$\mathcal{K}_1 = \{\text{O}=\text{C}, \text{C}-\text{C}, \text{O}-\text{C}, \text{O}-\text{H}\}$$

$$t_{\mathcal{C}}(\text{C}-\text{C}) = \text{avg. energy} = 3.5$$

$$t_{\text{edge}}(e) \approx f(\text{C}-\text{C}, G, e) \cdot t_{\mathcal{C}}(\text{C}-\text{C}) + f(G, \text{O}=\text{C}-\text{C}) \cdot t_{\mathcal{C}}(\text{O}=\text{C}-\text{C}) + f(G, \text{O}-\text{C}-\text{C}) \cdot t_{\mathcal{C}}(\text{O}-\text{C}-\text{C})$$

$$\mathcal{K}_1 \subseteq \mathcal{K}_2 \subseteq \dots \subseteq \mathcal{K}_k$$



$$\mathcal{K}_1 = \{\text{O}=\text{C}, \text{C}-\text{C}, \text{O}-\text{C}, \text{O}-\text{H}\}$$

$$t_{\mathcal{C}}(\text{C}-\text{C}) = \text{avg. energy} = 3.5$$

$$t_{\text{edge}}(e) \approx f(\text{C}-\text{C}, G, e) \cdot t_{\mathcal{C}}(\text{C}-\text{C}) + f(G, \text{O}=\text{C}-\text{C}) \cdot t_{\mathcal{C}}(\text{O}=\text{C}-\text{C}) + f(G, \text{O}-\text{C}-\text{C}) \cdot t_{\mathcal{C}}(\text{O}-\text{C}-\text{C})$$

$$\mathcal{K}_1 \subseteq \mathcal{K}_2 \subseteq \dots \subseteq \mathcal{K}_k$$

$$t_{\text{edge}}(e) \approx \sum_{C \in \mathcal{K}_i} f(C, G, e) \cdot t_{\mathcal{C}}(C)$$

$$t_{edge}(e) \approx \sum_{C \in \mathcal{K}_i} f(C, G, e) \cdot t_C(C)$$

$$t_{edge}(e) \approx \sum_{C \in \mathcal{K}_i} f(C, G, e) \cdot t_C(C)$$

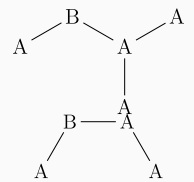
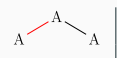
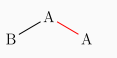
$$\Rightarrow t_{obs}(S) - \sum_{C \in \mathcal{K}_{k-1}} f(C, S) \cdot t_C(C) = \sum_{C \in \mathcal{C}_k^S} f(C, S) \cdot t_C(C) + \epsilon$$

# Frequency Matrix and Learning the Significant Contexts

$$t_{edge}(e) \approx \sum_{C \in \mathcal{K}_i} f(C, G, e) \cdot t_C(C)$$

$$\Rightarrow t_{obs}(S) - \sum_{C \in \mathcal{K}_{k-1}} f(C, S) \cdot t_C(C) = \sum_{C \in \mathcal{K}_k^S} f(C, S) \cdot t_C(C) + \epsilon$$

$X =$

			...
2	2	...	
0	1	...	
⋮	⋮	⋮	

# Frequency Matrix and Learning the Significant Contexts

$$t_{edge}(e) \approx \sum_{C \in \mathcal{K}_i} f(C, G, e) \cdot t_C(C)$$

$$\Rightarrow t_{obs}(S) - \sum_{C \in \mathcal{K}_{k-1}} f(C, S) \cdot t_C(C) = \sum_{C \in \mathcal{C}_k^S} f(C, S) \cdot t_C(C) + \epsilon$$

$$X = \begin{array}{c|cc|c} \begin{array}{c} \begin{array}{c} A \quad B \quad A \quad A \\ \diagdown \quad \diagup \quad \diagdown \quad \diagup \\ \quad \quad B \quad \quad A \\ \quad \quad \diagdown \quad \diagup \\ \quad \quad \quad \quad A \\ \quad \quad \quad \quad \diagdown \quad \diagup \\ \quad \quad \quad \quad \quad \quad A \quad A \end{array} \\ \vdots \end{array} & \begin{array}{c} \begin{array}{c} A \quad A \quad A \\ \diagdown \quad \diagup \quad \diagdown \end{array} \\ 2 \\ \vdots \end{array} & \begin{array}{c} \begin{array}{c} B \quad A \quad A \\ \diagdown \quad \diagup \quad \diagdown \end{array} \\ 2 \\ 1 \\ \vdots \end{array} & \begin{array}{c} \dots \\ \dots \\ \dots \\ \ddots \end{array} \end{array}$$

$$\text{LASSO: } \min \left( \sum_{i=1}^{|\mathcal{C}_1^S| + |S|} \left( y_i - \sum_{j=1}^{|\mathcal{C}_k^S|} X_{ij} t_j \right)^2 + \lambda \sum_{j=1}^{|\mathcal{C}_k^S|} |t_j| \right).$$

## Some considerations about contexts

- Structural information that determines bond energies must be stored in their frequencies in  $\mathcal{S}$ .

## Some considerations about contexts

- Structural information that determines bond energies must be stored in their frequencies in  $\mathcal{S}$ .
- The number of non-isomorphic subgraphs in a graph can grow exponentially with its size.



## Some considerations about contexts

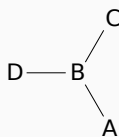
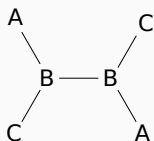
- Structural information that determines bond energies must be stored in their frequencies in  $\mathcal{S}$ .
- The number of non-isomorphic subgraphs in a graph can grow exponentially with its size.
- Contexts that only occur in "few" samples are unreliable.

## Some considerations about contexts

- Structural information that determines bond energies must be stored in their frequencies in  $\mathcal{S}$ .
- The number of non-isomorphic subgraphs in a graph can grow exponentially with its size.
- Contexts that only occur in "few" samples are unreliable.
- The frequencies of two contexts  $C_1 = (G_1, e_1)$  and  $C_2 = (G_2, e_2)$  where  $G_1 \simeq G_2$  are collinear in  $\mathcal{S}$ .

## Definition

The support  $\text{sup}(C)$  of a context  $C \in \mathcal{C}^S$  is the number of graphs in  $S$  which  $C$  can be embedded into. Given a positive integer  $\tau$  we say that  $C$  is supported if  $\text{sup}(C) \geq \tau$ .

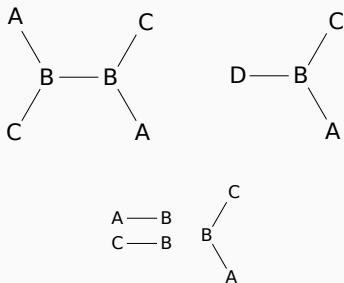


## Definition

The support  $\text{sup}(C)$  of a context  $C \in \mathcal{C}^S$  is the number of graphs in  $S$  which  $C$  can be embedded into. Given a positive integer  $\tau$  we say that  $C$  is supported if  $\text{sup}(C) \geq \tau$ .

## Definition

Let  $\mathcal{G}$  be a set of graphs and  $k$  and  $\tau$  two integers such that  $k > 0$  and  $\tau > 0$ . Then  $\text{FSM}(\mathcal{G}, k, \tau)$  is the set of all subgraphs in  $\mathcal{G}$  that contains  $k$  edges and are subgraph isomorphic to at least  $\tau$  graphs in  $\mathcal{G}$ .

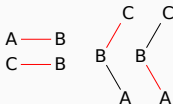
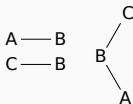
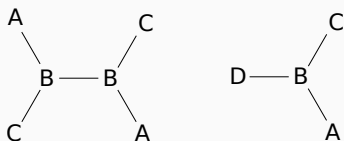


## Definition

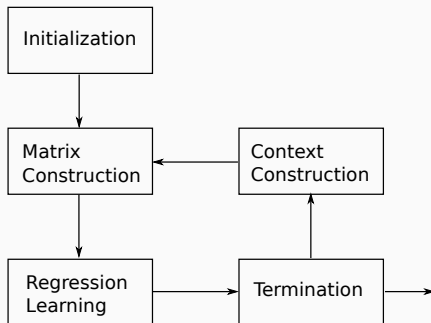
The support  $\text{sup}(C)$  of a context  $C \in \mathcal{C}^S$  is the number of graphs in  $S$  which  $C$  can be embedded into. Given a positive integer  $\tau$  we say that  $C$  is supported if  $\text{sup}(C) \geq \tau$ .

## Definition

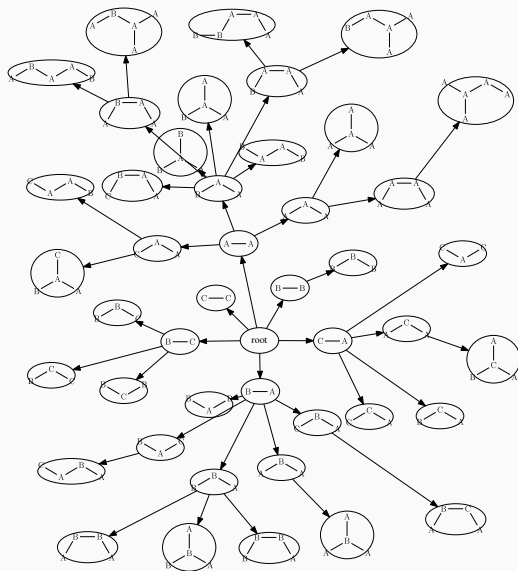
Let  $\mathcal{G}$  be a set of graphs and  $k$  and  $\tau$  two integers such that  $k > 0$  and  $\tau > 0$ . Then  $\text{FSM}(\mathcal{G}, k, \tau)$  is the set of all subgraphs in  $\mathcal{G}$  that contains  $k$  edges and are subgraph isomorphic to at least  $\tau$  graphs in  $\mathcal{G}$ .



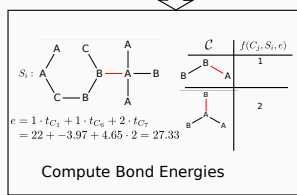
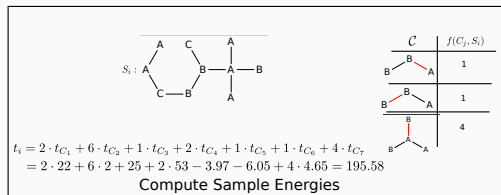
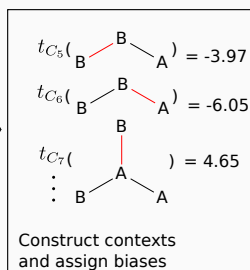
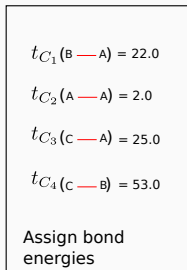
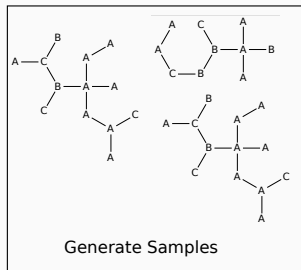
# Algorithm



# Predicting new graphs

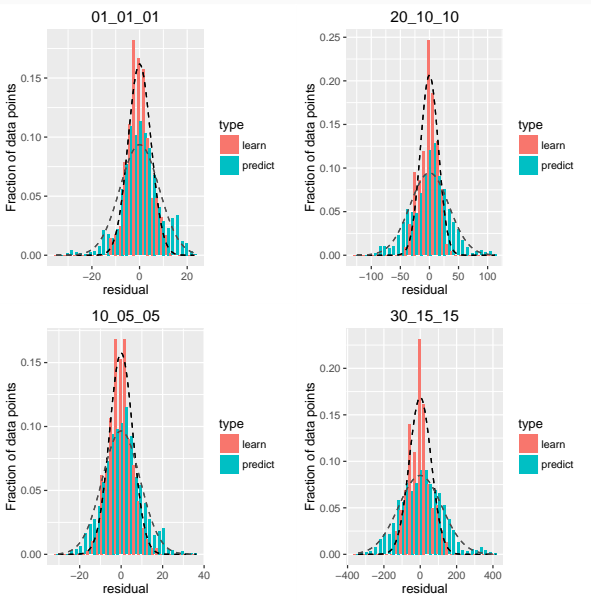


# How does it perform: Construction of synthetic dataset.

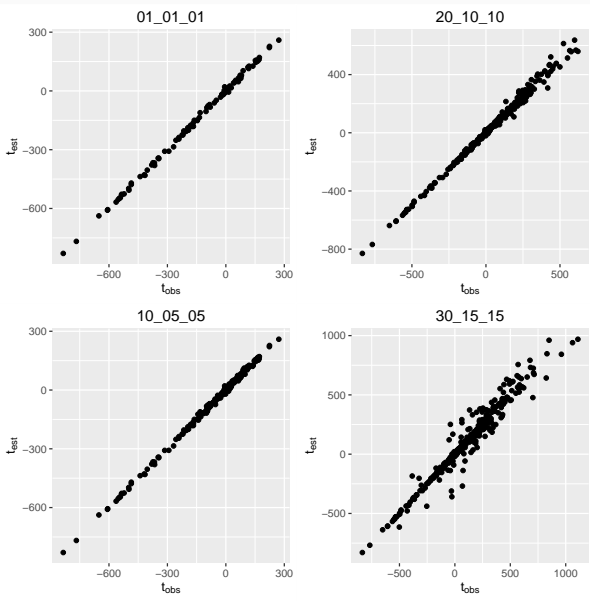




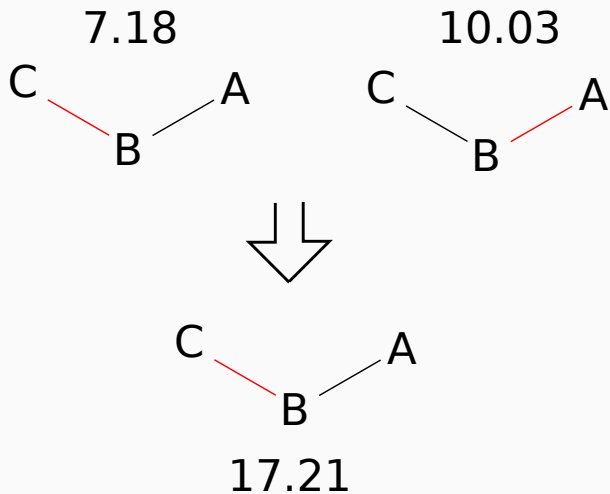
# Results: Synthetic dataset



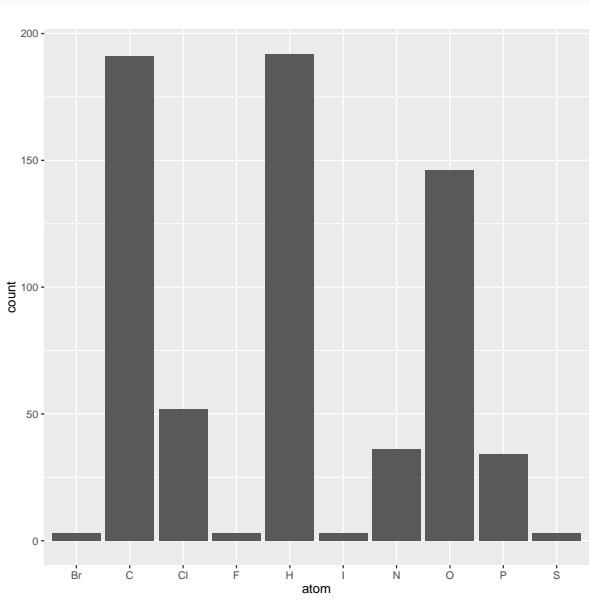
# Results: Synthetic dataset



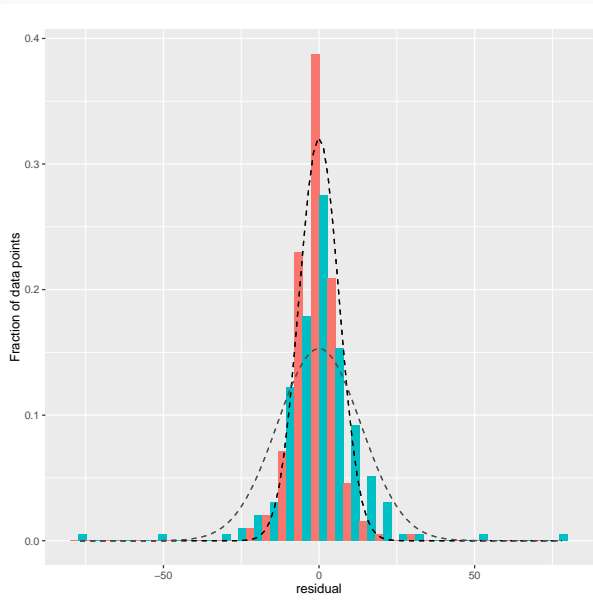
## Results: Synthetic dataset



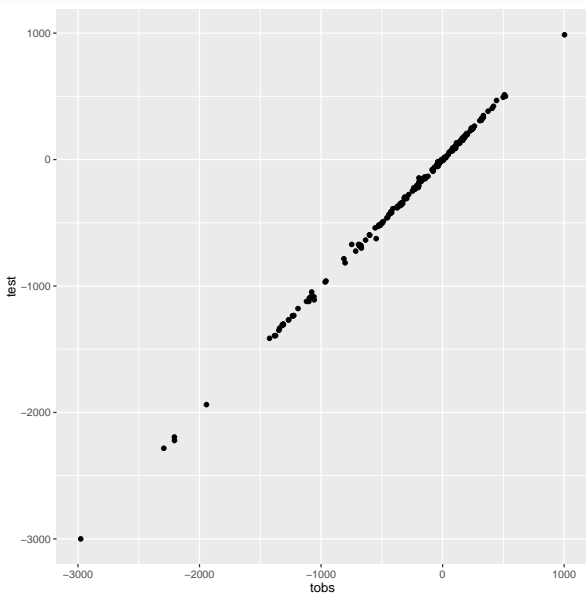
## Results: Gibbs Free Energy in metabolic networks



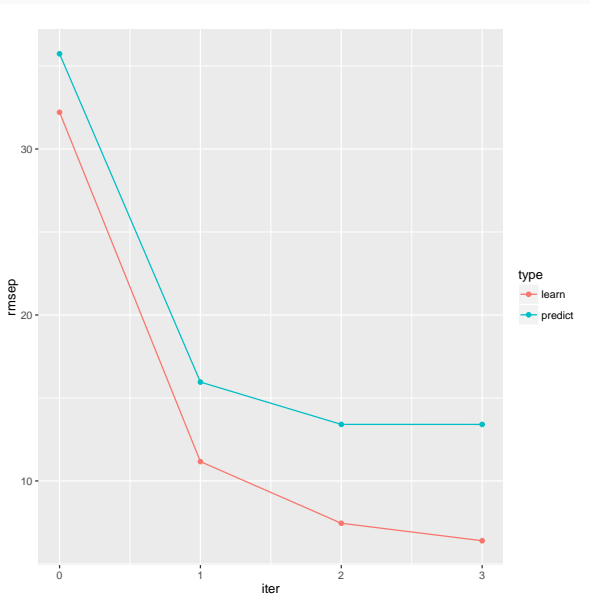
# Results: Gibbs Free Energy in metabolic networks



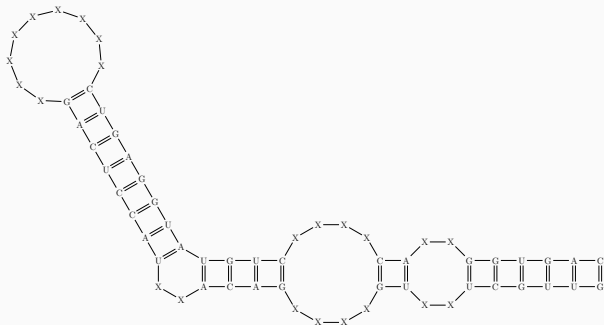
## Results: Gibbs Free Energy in metabolic networks



## Results: Gibbs Free Energy in metabolic networks

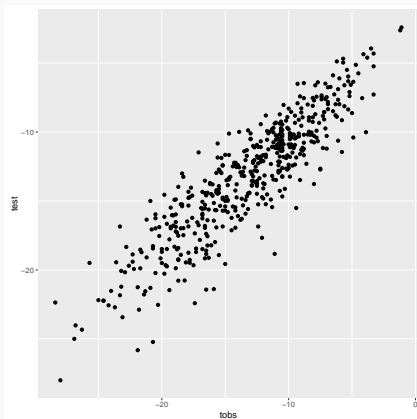
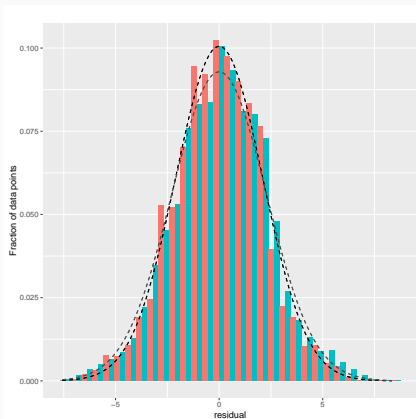


# Results: Minimum Free Energy of RNA secondary structures

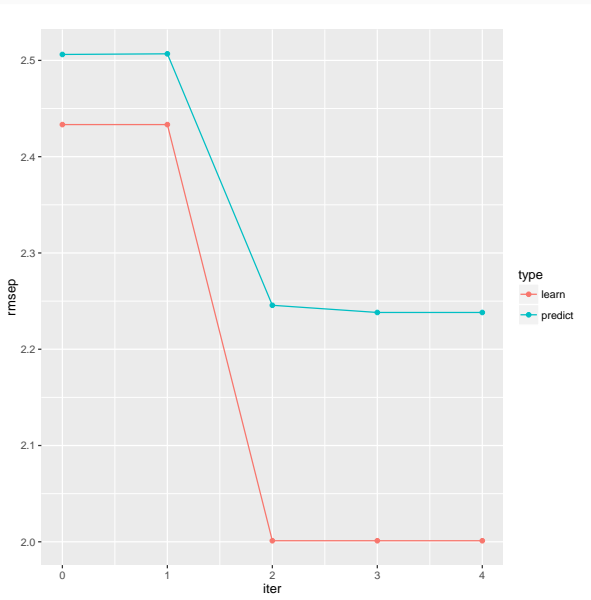




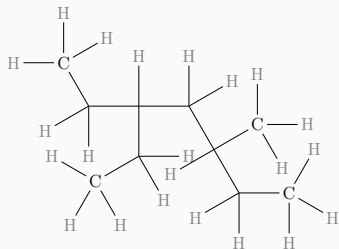
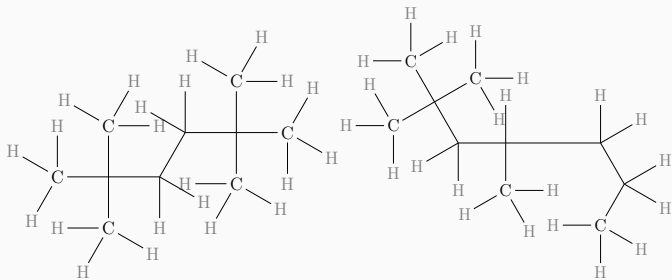
# Results: Minimum Free Energy of RNA secondary structures



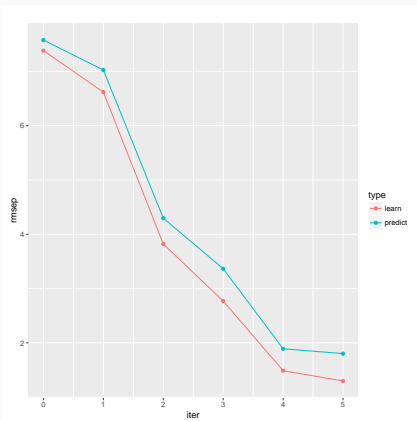
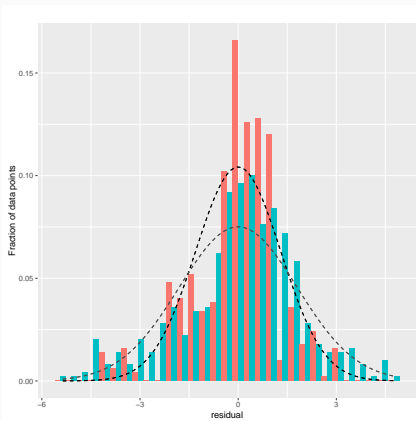
# Results: Minimum Free Energy of RNA secondary structures



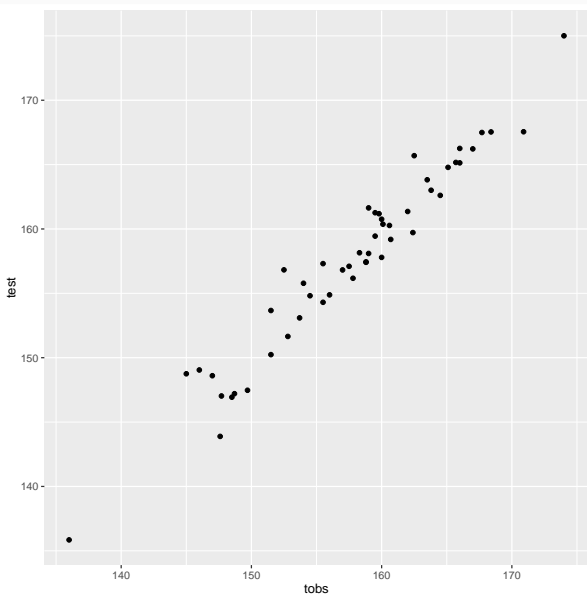
## Results: Boiling point acyclic molecules



# Results: Boiling point acyclic molecules



# Results: Boiling point acyclic molecules



## Conclusion

- Constructed a generic group contribution method based on the approximation of molecular energies.
- Can be used on a wide range of thermo dynamic properties.