

Adjustive Linear Regression and Its Application to the Inverse QSAR^{*}

Jianshen Zhu¹, Kazuya Haraguchi¹, Hiroshi Nagamochi¹ and Tatsuya Akutsu²

¹ Department of Applied Mathematics and Physics, Kyoto University, Kyoto 606-8501, Japan

² Bioinformatics Center, Institute for Chemical Research, Kyoto University, Uji 611-0011, Japan

Abstract

In this paper, we propose a new machine learning method, called adjustive linear regression, which can be regarded as an ANN on an architecture with an input layer and an output layer of a single node, wherein an error function is minimized by choosing not only weights of the arcs but also an activation function at each node in the two layers simultaneously. Under some conditions, such a minimization can be formulated as a linear program (LP) and a prediction function with adjustive linear regression is obtained as an optimal solution to the LP. We apply the new machine learning method to a framework of inferring a chemical compound with a desired property (i.e., inverse QSAR). From the results of our computational experiments, we observe that a prediction function constructed by adjustive linear regression for some chemical properties drastically outperforms that by Lasso linear regression.

Keywords: Machine Learning, Linear Regression, Integer Programming, Linear Program, Cheminformatics, Materials Informatics, QSAR/QSPR, Molecular Design.

1 Introduction

In this paper, we design a new learning method, called “adjustive linear regression” in order to construct a function that predicts a chemical property of a given chemical compound. We start with the background and the recent results on the research.

Background Analysis of chemical compounds is one of the important applications of intelligent computing. Indeed, various machine learning methods have been applied to the prediction of chemical activities from their structural data, where such a problem is often referred to as *quantitative structure activity relationship* (QSAR) [1, 2]. Recently, neural networks and deep-learning technologies have extensively been applied to QSAR [3].

In addition to QSAR, extensive studies have been done on inverse quantitative structure activity relationship (inverse QSAR), which seeks for chemical structures having desired chemical activities under some constraints. Since it is difficult to directly handle chemical structures in both QSAR and inverse QSAR, chemical compounds are usually represented as vectors of real or integer numbers, which are often called *descriptors* in chemoinformatics and correspond to *feature vectors* in machine learning. One major approach in inverse QSAR is to infer feature vectors from given chemical activities and constraints and then reconstruct chemical structures from these feature

^{*}Department of Applied Mathematics and Physics, Kyoto University, Technical Report, TR: 2021-002, September 3, 2021

vectors [4, 5, 6], where chemical structures are usually treated as undirected graphs. However, the reconstruction itself is a challenging task because the number of possible chemical graphs is huge. For example, chemical graphs with up to 30 atoms (vertices) \mathbf{C} , \mathbf{N} , \mathbf{O} , and \mathbf{S} may exceed 10^{60} [7]. Due to this difficulty, most existing methods for inverse QSAR do not guarantee optimal or exact solutions.

As a new approach, extensive studies have recently been done for inverse QSAR using *artificial neural networks* (ANNs), especially using graph convolutional networks [8]. For example, recurrent neural networks [10, 11], variational autoencoders [9], grammar variational autoencoders [12], generative adversarial networks [13], and invertible flow models [14, 15] have been applied. However, these methods do not yet guarantee optimal or exact solutions.

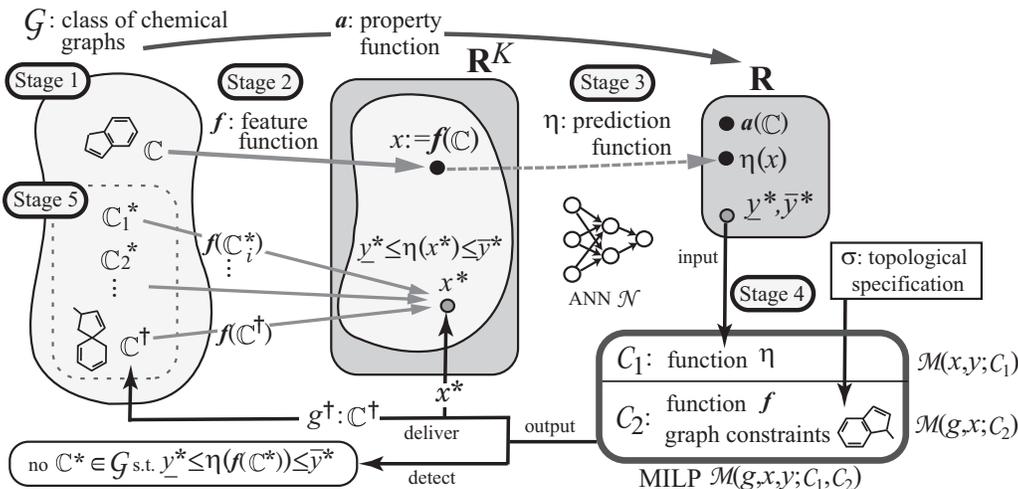


Figure 1: An illustration of a framework for inferring a set of chemical graphs \mathbf{C}^* .

Framework Akutsu and Nagamochi [16] proved that the computation process of a given ANN can be simulated with a mixed integer linear programming (MILP). Based on this, a novel framework for inferring chemical graphs has been developed and revised [17, 18], as illustrated in Figure 1. It constructs a prediction function in the first phase and infers a chemical graph in the second phase. The first phase of the framework consists of three stages. Stage 1 chooses a chemical property π and a class \mathcal{G} of graphs, where a property function a is defined so that $a(\mathbf{C})$ is the value of π for a compound $\mathbf{C} \in \mathcal{G}$, and collects a data set D_π of chemical graphs in \mathcal{G} such that $a(\mathbf{C})$ is available for every $\mathbf{C} \in D_\pi$. Stage 2 introduces a feature function $f : \mathcal{G} \rightarrow \mathbb{R}^K$ for a positive integer K . Stage 3 constructs a prediction function η with an ANN \mathcal{N} that, given a vector $x \in \mathbb{R}^K$, returns a value $y = \eta(x) \in \mathbb{R}$ so that $\eta(f(\mathbf{C}))$ serves as a predicted value to the real value $a(\mathbf{C})$ of π for each $\mathbf{C} \in D_\pi$. Given two reals y^* and \bar{y}^* as an interval for a target chemical value, the second phase infers chemical graphs \mathbf{C}^* with $y^* \leq \eta(f(\mathbf{C}^*)) \leq \bar{y}^*$ in the next two stages. After Stage 3, we have obtained a feature function f and a prediction function η . We can specify an additional constraint on the substructures of target chemical graphs, called a *topological specification* before we infer a target chemical graph. In Stage 4, the following two MILP formulations are prepared:

- MILP $\mathcal{M}(x, y; \mathcal{C}_1)$ with a set \mathcal{C}_1 of linear constraints on variables x and y (and some other auxiliary variables) simulates the process of computing $y := \eta(x)$ from a vector x ; and

- MILP $\mathcal{M}(g, x; \mathcal{C}_2)$ with a set \mathcal{C}_2 of linear constraints on variable x and a variable vector g that represents a chemical graph \mathbb{C} (and some other auxiliary variables) simulates the process of computing $x := f(\mathbb{C})$ from a chemical graph \mathbb{C} and chooses a chemical graph \mathbb{C} that satisfies the given topological specification σ .

Given an interval with boundaries $\underline{y}^*, \bar{y}^* \in \mathbb{R}$, Stage 4 solves the combined MILP $\mathcal{M}(g, x, y; \mathcal{C}_1, \mathcal{C}_2)$ to find a feature vector $x^* \in \mathbb{R}^K$ and a chemical graph \mathbb{C}^\dagger with the specification σ such that $f(\mathbb{C}^\dagger) = x^*$ and $\underline{y}^* \leq \eta(x^*) \leq \bar{y}^*$ (where if the MILP instance is infeasible then this suggests that there does not exist such a desired chemical graph). Stage 5 generates other chemical graphs \mathbb{C}^* such that $\underline{y}^* \leq \eta(f(\mathbb{C}^*)) \leq \bar{y}^*$ based on the output chemical graph \mathbb{C}^\dagger .

A modeling of chemical compounds together with an MILP formulation has been improved so that a chemical compound with any graph structure can be treated (see Shi et al. [18]). Not only ANNs but also other machine learning methods have been used to construct a prediction function η in Stage 3 recently. Tanaka et al. [19] (resp., Zhu et al. [20]) used a decision tree (resp., linear regression) to construct a prediction function η in Stage 3 in the framework and derived an MILP $\mathcal{M}(x, y; \mathcal{C}_1)$ that simulates the computation process of a decision tree (resp., linear regression).

The novelty of the framework is based on the fact that a prediction process by linear regression or ANNs can be modeled as an MILP, to which we can find a mathematically exact solution by relying on the state-of-the-art solvers from Operations Research (OR). A sophisticated method has been studied by Shi et al. [18] in order to formulate a sparse MILP instance even for a complicated requirements in a topological specification. Currently an MILP instance in Stage 4 for inferring a chemical compound with 50 non-hydrogen atoms contains around 10,000 variables and 10,000 linear constraints, and can be solved in a few seconds.

Contribution In this paper, we apply a mathematical programming in OR to Stage 3 in order to design a new machine learning method for QSAR. Let us compare linear regression and ANNs. The former uses a hyperplane to explain a given data set and the latter can represent a more complex subspace than a hyperplane. Importantly a best hyperplane that minimizes an error function can be found exactly in the former whereas a local optimum solution to an error function is constructed by an iterative procedure in the latter and different local optimum solutions often appear depending on how we have tuned many parameters in ANNs. Linear regression can be regarded as an ANN on an architecture with an input layer and an output layer of a single node with a linear activation function. We consider an ANN on the same architecture such that each node in the input and out layers is equipped with a set Φ of activation functions. Given a data set, we consider a problem of minimizing an error function on the data set by choosing a weight of each arc, a bias of the output node and a best activation function for each node simultaneously. With some restriction on the set Φ of activation functions and the definition of an error function, we show that such an minimization problem can be formulated as a linear program, which is much easier than an MILP to solve exactly. We call this new method "adjustive linear regression" and implemented it in Stages 3 and 4 in the framework. We used the same MILP $\mathcal{M}(g, x; \mathcal{C}_2)$ formulation proposed by Zhu et al. [20] and omit the details in this paper. We compared adjustive linear regression with Lasso linear regression in constructing prediction functions for several chemical properties. From the results of our computational experiments, we observe that a prediction function constructed by adjustive linear regression for some chemical properties drastically outperform that by Lasso

linear regression.

The paper is organized as follows. Section 2 reviews the idea of prediction functions based on linear regression and ANNs and designs “adjustive linear regression,” a new method for constructing a prediction function by solving a linear program to optimize a choice of weights/bias together with activation functions in an ANN with no hidden layers. Section 3.1 introduces some notions on graphs, a modeling of chemical compounds and a choice of descriptors. Section 4 reviews a method, called a *two-layered model* for representing the feature of a chemical graph in order to deal with an arbitrary graph in the framework. Section 5 reports the results on some computational experiments conducted for the framework of inferring chemical graphs by using our new method of adjustive linear regression. Section 6 makes some concluding remarks. Some technical details are given in Appendices: Appendix A for all descriptors in our feature function; Appendix B for a full description of a topological specification; and Appendix C for the detail of test instances used in our computational experiment for Stages 4 and 5.

2 Constructing Prediction Functions

Let \mathbb{R} , \mathbb{R}_+ , \mathbb{Z} and \mathbb{Z}_+ denote the sets of reals, non-negative reals, integers and non-negative integers, respectively. For two integers a and b , let $[a, b]$ denote the set of integers i with $a \leq i \leq b$. For a vector $x \in \mathbb{R}^p$, the j -th entry of x is denoted by $x(j)$, $j \in [1, p]$.

2.1 Linear Prediction Functions

For an integer $K \geq 1$, define a feature space \mathbb{R}^K . Let $\mathcal{X} = \{x_1, x_2, \dots, x_m\}$ be a set of feature vectors $x \in \mathbb{R}^K$ and let $a_i \in \mathbb{R}$ be a real assigned to a feature vector x_i . Let $A = \{a_i \mid i \in [1, m]\}$. A function $\eta : \mathbb{R}^K \rightarrow \mathbb{R}$ is called a *prediction function*. We wish to find a prediction function $\eta : \mathbb{R}^K \rightarrow \mathbb{R}$ based on a subset of $\{x_1, x_2, \dots, x_m\}$ so that $\eta(x_i)$ is closed to the value a_i for many indices $i \in [1, m]$.

For a prediction function $\eta : \mathbb{R}^K \rightarrow \mathbb{R}$, define an error function

$$\text{Err}(\eta; \mathcal{X}) \triangleq \sum_{i \in [1, m]} (a_i - \eta(x_i))^2,$$

and define the *coefficient of determination* $R^2(\eta, \mathcal{X})$ to be

$$R^2(\eta, \mathcal{X}) \triangleq 1 - \frac{\text{Err}(\eta; \mathcal{X})}{\sum_{i \in [1, m]} (a_i - \tilde{a})^2} \text{ for } \tilde{a} = \frac{1}{m} \sum_{i \in [1, m]} a_i.$$

Many methods have been proposed in order to find a prediction function η that minimizes the error function $\text{Err}(\eta_{w,b}; \mathcal{X})$ possibly without using all elements in \mathcal{X} .

For the feature space \mathbb{R}^K , a hyperplane is defined to be a pair (w, b) of a vector $w \in \mathbb{R}^K$ and a real $b \in \mathbb{R}$. A prediction function η is called *linear* if η is given by $\eta_{w,b}(x) = w \cdot x + b$, $x \in \mathbb{R}^K$ for a hyperplane (w, b) . The linear regression is to find a hyperplane (w, b) that minimizes $\text{Err}(\eta_{w,b}; \mathcal{X}) = \sum_{i \in [1, m]} (a_i - (w \cdot x_i + b))^2$.

In many cases, a feature vector f contains descriptors that do not play an essential role in constructing a good prediction function. When we solve the minimization problem, the entries $w(j)$ for some descriptors $j \in [1, K]$ in the resulting hyperplane (w, b) become zero, which means that these descriptors were not necessarily important for finding a prediction function $\eta_{w,b}$. It is proposed that solving the minimization with an additional penalty term to the error function often results in a more number of entries $w(j) = 0$, reducing a set of descriptors necessary for defining a prediction function $\eta_{w,b}$. For an error function with such a penalty term, a Ridge function $\frac{1}{2m}\text{Err}(\eta_{w,b}; \mathcal{X}) + \lambda[\sum_{j \in [1, K]} w(j)^2 + b^2]$ [21, 22] and a Lasso function $\frac{1}{2m}\text{Err}(\eta_{w,b}; \mathcal{X}) + \lambda[\sum_{j \in [1, K]} |w(j)| + |b|]$ [23] are known, where $\lambda \in \mathbb{R}$ is a given real number. As a hybridization of Ridge linear regression and Lasso linear regression, a linear regression that minimizes an error function defined to be $\frac{1}{2m}\text{Err}(\eta_{w,b}; \mathcal{X}) + \lambda_2[\sum_{j \in [1, K]} w(j)^2 + b^2] + \lambda_1[\sum_{j \in [1, K]} |w(j)| + |b|]$ is called elastic net linear regression [24], where $\lambda_1, \lambda_2 \in \mathbb{R}$ are given real numbers.

Zhu et al. [20] used Lasso linear regression to construct a prediction function η in Stage 3 in the framework.

2.2 ANNs for Linear Prediction Functions

It is not difficult to see that a linear prediction function η with a hyperplane (w, b) can be represented by an ANN \mathcal{N} with an input layer $L_{\text{in}} = \{u_1, u_2, \dots, u_K\}$ of K input nodes and an output layer $L_{\text{out}} = \{v\}$ of a single output node v such that the weight of an arc (u_j, v) from an input node u_j to the output node v is given by $w(j), j \in [1, K]$; the bias at node v is given by b ; and the activation function at node v is linear. See Figure 2(a) for an illustration of an ANN \mathcal{N} that represents a linear prediction function η with a hyperplane (w, b) . Given a vector $x \in \mathbb{R}^K$, the ANN \mathcal{N} outputs $y := \sum_{j \in [1, K]} w(j)x(j) + b$.

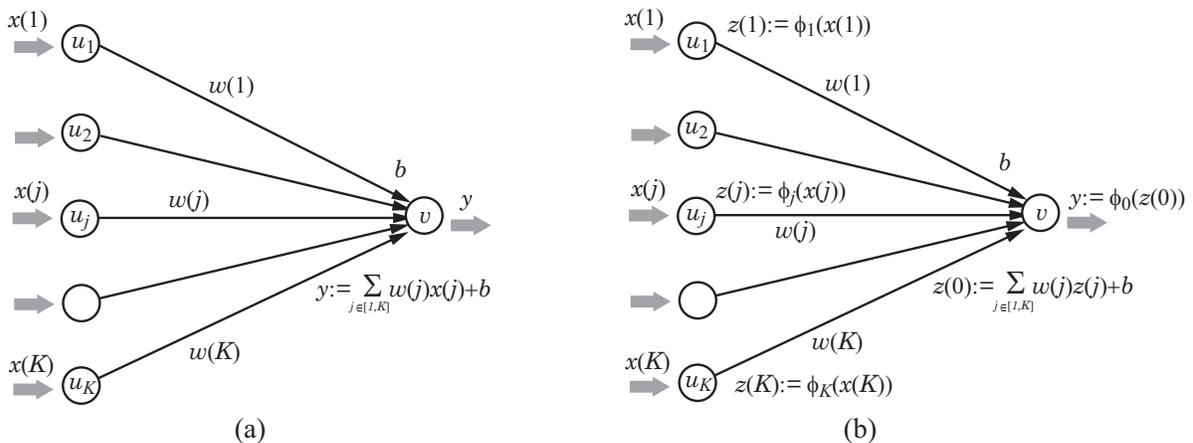


Figure 2: An illustration of the process in ANNs with no hidden layers: (a) An ANN \mathcal{N} that represents a linear prediction function η with a hyperplane (w, b) ; (b) an ANN \mathcal{N}_ϕ with activation functions $\phi_j, j \in [0, K]$ at all nodes.

We consider an ANN \mathcal{N}_ϕ with the same architecture with the ANN \mathcal{N} and introduce activation functions ϕ_j at nodes $u_j, j \in [1, K]$ and an activation function ϕ_0 at node v . Given a vector $x \in \mathbb{R}^K$,

the ANN \mathcal{N}_ϕ outputs $y := \phi_0(z(0))$ for $z(0) := \sum_{j \in [1, K]} w(j)z(j) + b$ and $z(j) := \phi_j(x(j))$, $j \in [1, K]$.

In a standard method of a prediction function $\eta_{\mathcal{N}_\phi}$ with the above ANN \mathcal{N}_ϕ , we specify each activation function ϕ_j and determine weights w and a bias b by executing an iterative procedure that tries to minimize an error function between the real values a_i and the predicted values $\eta_{\mathcal{N}_\phi}(x_i)$.

2.3 Adjustive Linear Regression

In this paper, we design a new method of constructing a prediction function with the above ANN \mathcal{N}_ϕ so that (i) not only weights w and a bias b but also prediction functions ϕ_j are chosen so as to minimize an error function and (ii) the minimization problem is formulated as a linear programming problem.

We introduce a class Φ_j of functions for a choice of each activation function ϕ_j , $j \in [0, K]$. When we choose a function $\phi_j \in \Phi_j$ for each $j \in [0, K]$ and a hyperplane (w, b) , we define a prediction function $\eta_{\Psi, w, b}$ such that

$$\eta_{\Psi, w, b}(x) \triangleq \phi_0\left(\sum_{j \in [1, K]} w(j)(\phi_j(x(j)))\right) - b$$

for the set $\Psi = \{\phi_j \mid j \in [0, K]\}$ of the functions.

In this paper, we use a function $\xi(t) = ct + c't^2 + c''(1 - (t-1)^2)$, $0 \leq t \leq 1$ for a function ϕ_j , $j \in [1, K]$ or the inverse ϕ_0^{-1} of a function ϕ_0 , where c, c' and c'' are nonnegative constant constants with $c + c' + c'' = 1$ which will be determined for each $j \in [0, K]$ by our method. Note that, for a domain $0 \leq t \leq 1$, $\xi(t)$ is a monotone increasing function and admits an inverse function $\xi^{-1}(t)$, where $\xi^{-1}(t) = t/(c + 2c'')$ when $c' = c''$ and $\xi^{-1}(t) = (-c - 2c'' + \sqrt{(c + 2c'')^2 + 4(c' - c'')t})/(2c' - 2c'')$ when $c' \neq c''$.

We introduce a class Φ_j of functions in the following way.

1. Normalize the set $\{x_i(j) \mid x_i \in \mathcal{X}\}$, $j \in [1, K]$ and the set $\{a_i(j) \mid x_i \in \mathcal{X}\}$ so that the minimum and maximum in the set become 0 and 1.
2. For each index $j \in [0, K]$, define a class Φ_j of functions to be

$$\Phi_j \triangleq \{c_0(j)t + c_1(j)t^2 + c_2(j)(1 - (t-1)^2), 0 \leq t \leq 1 \mid c_q(j) \geq 0, q \in [0, 2], \sum_{q \in [0, 2]} c_q(j) = 1\}, j \in [1, K].$$

Define

$$\tilde{\Phi}_0 \triangleq \{c_0(0)t + c_1(0)t^2 + c_2(0)(1 - (t-1)^2), 0 \leq t \leq 1 \mid c_q(0) \geq 0, q \in [0, 2], \sum_{q \in [0, 2]} c_q(0) = 1\},$$

$$\Phi_0 \triangleq \{\xi^{-1}(t), 0 \leq t \leq 1 \mid \xi(t) \in \tilde{\Phi}_0\}.$$

To use linear programming, we measure an error of a prediction function η over a data set \mathcal{X} by the sum of the absolute errors:

$$\text{SAE}(\eta; \mathcal{X}) \triangleq \sum_{x_i \in \mathcal{X}} |a_i - \eta(x_i)|.$$

Now our aim is to find a prediction function $\eta_{\Psi,w,b}$ that minimizes the sum of the absolute errors $\text{SAE}(\eta_{\Psi,w,b}; \mathcal{X})$ or equivalently

$$\sum_{i \in [1,m]} |\phi_0^{-1}(a_i) - (\sum_{j \in [1,K]} w(j)(\phi_j(x_i(j)))) - b| \quad (1)$$

over all functions $\phi_0 \in \tilde{\Phi}_0, \phi_j \in \Phi_j, j \in [1, K]$ and hyperplanes (w, b) .

To formulate this minimization problem as a linear programming problem, we predetermine the sign of $w(j)$ for each descriptor j in a hyperplane (w, b) that we will choose. Compute the correlation coefficient $\sigma(X_j, A)$ between $X_j = \{x_i(j) \mid i \in [1, m]\}$ and $A = \{a_i \mid i \in [1, m]\}$ and partition the set of descriptors into two sets $I^+ := \{j \in [1, K] \mid \sigma(X_j, A) \geq 0\}$ and $I^- := \{j \in [1, K] \mid \sigma(X_j, A) < 0\}$. We impose an additional constraint that $w(j) \geq 0, j \in I^+$ and $w(j) \leq 0, j \in I^-$. Then the objective function (1) is described as follows, where we rewrite each term $w(j), j \in I^+$ (resp., $-w(j), j \in I^-$) as $w'(j)$:

$$\begin{aligned} & \sum_{i \in [1,m]} \left| c_0(0)a_i + c_1(0)a_i^2 + c_2(0)(1 - (a_i - 1)^2) \right. \\ & \quad - \sum_{j \in I^+} [w'(j)(c_0(j)x_i(j) + c_1(j)x_i(j)^2 + c_2(j)(1 - (x_i(j) - 1)^2))] \\ & \quad \left. + \sum_{j \in I^-} [w'(j)(c_0(j)x_i(j) + c_1(j)x_i(j)^2 + c_2(j)(1 - (x_i(j) - 1)^2))] - b \right|. \end{aligned}$$

We minimize (1) over all nonnegative reals $c_q(j), q \in [0, 2], j \in [1, K]$, nonnegative reals $w(j), j \in [1, K]$ and a real $b \in \mathbb{R}$ such that $\sum_{q \in [0,2]} c_q(j) = 1, j \in [1, K]$.

Before we derive a linear programming formulation to the above minimization problem, we include a penalty term for the weights $w(j), j \in [1, K]$ analogously with the Lasso linear regression. We consider the following problem which we call *adjustive linear regression*, where $w'(j)c_q(j), q \in [0, 2]$ is rewritten as $w_q(j)$.

Adjustive Linear Regression(\mathcal{X}, λ)

$$\begin{aligned} \text{Minimize: } & \frac{1}{2m} \sum_{i \in [1,m]} \left| c_0(0)a_i + c_1(0)a_i^2 + c_2(0)(1 - (a_i - 1)^2) \right. \\ & \quad - \sum_{j \in I^+} [w_0(j)x_i(j) + w_1(j)x_i(j)^2 + w_2(j)(1 - (x_i(j) - 1)^2)] \\ & \quad \left. + \sum_{j \in I^-} [w_0(j)x_i(j) + w_1(j)x_i(j)^2 + w_2(j)(1 - (x_i(j) - 1)^2)] - b \right| \quad (2) \\ & + \lambda \sum_{j \in [1,K]} w_0(j) + \lambda |b| \end{aligned}$$

subject to

$$c_0(0) + c_1(0) + c_2(0) = 1.$$

When the set $S = \{a_i \mid i \in [1, m]\}$ (resp., $S = \{x_i(j) \mid i \in [1, m]\}$ for an index $j \in [1, m]$) is binary (i.e., $S = \{0, 1\}$), we always set $c_1(0) = c_2(0) = 0$ (resp., $c_1(j) = c_2(j) = 0$).

We observe that adjustive linear regression is an extension of the Lasso linear regression except that the error function is the sum of absolute errors in the former and the sum of square errors in the latter.

We solve the above minimization problem (2) to construct a prediction function $\eta_{\Psi, w, b}$. Let $c_q^*(0), q \in [0, 2]$, $w_q^*(j), q \in [0, 2], j \in [1, K]$ and b^* denote the values of variables $c_q(0), q \in [0, 2]$, $w_q(j), q \in [0, 2], j \in [1, K]$ and b in an optimal solution, respectively. Let K' denote the number of descriptors $j \in [1, K]$ with $w_0^*(j) > 0$ and $I_{K'}$ denote the set of $j \in [1, K]$ with $w_0^*(j) > 0$. Then we set

$$\begin{aligned} w^*(j) &:= 0 \text{ for } j \in [1, K] \text{ with } w_0^*(j) = 0, \\ w^*(j) &:= w_0^*(j)/(w_0^*(j) + w_1^*(j) + w_2^*(j)) \text{ for } j \in I^+ \cap I_{K'}, \\ w^*(j) &:= -w_0^*(j)/(w_0^*(j) + w_1^*(j) + w_2^*(j)) \text{ for } j \in I^- \cap I_{K'}, \\ c_q^*(j) &:= w_q^*(j)/w^*(j), q \in [0, 2] \text{ for } j \in I_{K'} \text{ and} \\ w^* &:= (w_0^*(1), w_0^*(2), \dots, w_0^*(K)) \in \mathbb{R}^K. \end{aligned}$$

For a set Ψ^* of selected functions $\phi_j(t) = c_0^*(j)t + c_1^*(j)t^2 + c_2^*(j)(1 - (t - 1)^2), j \in I_{K'}$ with and $\phi_0(t)$ with $\phi_0^{-1}(t) = c_0^*(0)t + c_1^*(0)t^2 + c_2^*(0)(1 - (t - 1)^2)$ and a hyperplane (w^*, b^*) , we construct a prediction function η_{Ψ^*, w^*, b^*} .

We propose the following scheme of executing adjustive linear regression for constructing a prediction function and evaluating the performance.

1. Given a data set $\mathcal{X} = \{x_i \in \mathbb{R}^K \mid i \in [1, m]\}$ of normalized feature vectors and a set $A = \{a_i \in \mathbb{R} \mid i \in [1, m]\}$ of normalized observed values, we choose a real $\lambda > 0$ possibly from a set of candidates for $\lambda > 0$ so that the performance of a prediction function η_{Ψ^*, w^*, b^*} obtained from an optimal solution (Ψ^*, w^*, b^*) to the adjustive linear regression (\mathcal{X}, λ) attains a criterion, where we may use cross-validation and the test coefficient of determination to know the performance. Note that the reals $c_q^*(j), q \in [1, 2], j \in I_{K'}$ indicate the non-linearity of a function ϕ_j and the resulting prediction function η_{Ψ^*, w^*, b^*} may easily cause overfitting when the bounds $c_q^*(j)$ are unnecessarily large. To bound reals $c_q^*(j), q \in [1, 2]$, we also penalize weights $w_q(j), q \in [1, 2], j \in I_{K'}$ with the same real λ .
2. With the real λ determined in 1, we evaluate the performance of a prediction function obtained with adjustive linear regression based on cross-validation. We divide the entire set \mathcal{X} into five subsets $\mathcal{X}^{(k)}, k \in [1, 5]$. For each $k \in [1, 5]$, we use the set $\mathcal{X} \setminus \mathcal{X}^{(k)}$ as a training data to construct a prediction function $\eta_{\Psi, w, b}$ with adjustive linear regression $(\mathcal{X} \setminus \mathcal{X}^{(k)}, \lambda)$ and compute the coefficient of determination $R^2(\eta_{\Psi, w, b}; \mathcal{X}^{(k)})$.

2.4 An LP formulation for Adjustive Linear Regression

We formulate a linear programming problem to the adjustive linear regression (\mathcal{X}, λ) .

LP (\mathcal{X}, λ) :

constants:

- A set $\mathcal{X} = \{x_i \in \mathbb{R}^K \mid i \in [1, m]\}$ of feature vectors and a set $A = \{a_i \in \mathbb{R} \mid i \in [1, m]\}$ of observed values. Assume that each of the sets $X_j = \{x_i(j) \mid i \in [1, m]\}, j \in [1, K]$ and A is standardized;
- A positive real $\lambda \in \mathbb{R}$: a coefficient for the penalty term;

variables:

- Nonnegative reals $c_q(0) \in \mathbb{R}, q \in [0, 2]$;

- Nonnegative vectors $w_q \in \mathbb{R}^K$, $q \in [0, 2]$ and a real $b \in \mathbb{R}$;
- Nonnegative real $\bar{b} \in \mathbb{R}$;
- Nonnegative reals $\Delta_i \geq 0$, $i \in [1, m]$;

constraints:

$$c_0(0) + c_1(0) + c_2(0) = 1, \quad (3)$$

$$\begin{aligned} & \Delta_i \geq c_0(0)a_i + c_1(0)a_i^2 + c_2(0)(1 - (a_i - 1)^2) \\ & - \sum_{j \in I^+} [w_0(j)x_i(j) + w_1(j)x_i(j)^2 + w_2(j)(1 - (x_i(j) - 1)^2)] \\ & + \sum_{j \in I^-} [w_0(j)x_i(j) + w_1(j)x_i(j)^2 + w_2(j)(1 - (x_i(j) - 1)^2)] - b \geq -\Delta_i, \quad i \in [1, m], \end{aligned} \quad (4)$$

$$\bar{b} \geq b \geq -\bar{b}, \quad (5)$$

objective function:

$$\text{Minimize } \frac{1}{2m} \sum_{i \in [1, m]} \Delta_i + \lambda \sum_{q \in [0, 2], j \in [1, K]} w_q(j) + \lambda \bar{b}.$$

We see that the numbers of variables and constraints in the linear program $\text{LP}(\mathcal{X}, \lambda)$ are both $O(m + K)$.

Let $w_q^*(j)$, $q \in [0, 2]$, $j \in [1, K]$ and b^* denote the values of variables $w_q(j)$, $q \in [0, 2]$, $j \in [1, K]$ and b in an optimal solution to linear program $\text{LP}(\mathcal{X}, \lambda)$, respectively. Let K' denote the number of descriptors $j \in [1, K]$ with $w_0^*(j) > 0$ and $I_{K'}$ denote the set of $j \in [1, K]$ with $w_0^*(j) > 0$. Then we obtain an optimal solution to the adjustive linear regression (2) by setting $w^*(j) := w_0^*(j)/(w_0^*(j) + w_1^*(j) + w_2^*(j))$, $j \in I^+ \cap I_{K'}$, $w^*(j) := -w_0^*(j)/(w_0^*(j) + w_1^*(j) + w_2^*(j))$, $j \in I^- \cap I_{K'}$, and $c_q^*(j) := w_q^*(j)/w^*(j)$, $q \in [1, 2]$, $j \in I_{K'}$.

3 Modeling of Chemical Compounds

This section introduces some notions and terminologies on graphs and reviews the modeling of chemical compounds due to Zhu et al. [20].

Graph Given a graph G , let $V(G)$ and $E(G)$ denote the sets of vertices and edges, respectively. For a subset $V' \subseteq V(G)$ (resp., $E' \subseteq E(G)$) of a graph G , let $G - V'$ (resp., $G - E'$) denote the graph obtained from G by removing the vertices in V' (resp., the edges in E'), where we remove all edges incident to a vertex in V' in $G - V'$. An edge subset $E' \subseteq E(G)$ in a connected graph G is called *separating* (resp., *non-separating*) if $G - E'$ remains connected (resp., $G - E'$ becomes disconnected). The *rank* $r(G)$ of a graph G is defined to be the minimum $|F|$ of an edge subset

$F \subseteq E(G)$ such that $G - F$ contains no cycle, where $r(G) = |E(G)| - |V(G)| + 1$. Observe that $r(G - E') = r(G) - |E'|$ holds for any non-separating edge subset $E' \subseteq E(G)$. An edge $e = u_1u_2 \in E(G)$ in a connected graph G is called a *bridge* if $\{e\}$ is separating, i.e., $G - e$ consists of two connected graphs G_i containing vertex u_i , $i = 1, 2$. For a connected cyclic graph G , an edge e is called a *core-edge* if it is in a cycle of G or is a bridge $e = u_1u_2$ such that each of the connected graphs G_i , $i = 1, 2$ of $G - e$ contains a cycle. A vertex incident to a core-edge is called a *core-vertex* of G . A path with two end-vertices u and v is called a u, v -*path*.

A vertex designated in a graph G is called a *root*. In this paper, we designate at most two vertices as roots, and denote by $\text{Rt}(G)$ the set of roots of G . We call a graph G *rooted* (resp., *bi-rooted*) if $|\text{Rt}(G)| = 1$ (resp., $|\text{Rt}(G)| = 2$), where we call G *unrooted* if $\text{Rt}(G) = \emptyset$.

For a graph G possibly with roots a *leaf-vertex* is defined to be a non-root vertex $v \in V(G) \setminus \text{Rt}(G)$ with degree 1, call the edge uv incident to a leaf vertex v a *leaf-edge*, and denote by $V_{\text{leaf}}(G)$ and $E_{\text{leaf}}(G)$ the sets of leaf-vertices and leaf-edges in G , respectively. For a graph or a rooted graph G , we define graphs $G_i, i \in \mathbb{Z}_+$ obtained from G by removing the set of leaf-vertices i times so that

$$G_0 := G; \quad G_{i+1} := G_i - V_{\text{leaf}}(G_i),$$

where we call a vertex $v \in V_{\text{leaf}}(G_k)$ a *leaf k -branch* and we say that a vertex $v \in V_{\text{leaf}}(G_k)$ has *height* $\text{ht}(v) = k$ in G . The *height* $\text{ht}(T)$ of a rooted tree T is defined to be the maximum of $\text{ht}(v)$ of a vertex $v \in V(T)$. For an integer $k \geq 0$, we call a rooted tree T *k -lean* if T has at most one leaf k -branch. For an unrooted cyclic graph G , we regard that the set of non-core-edges in G induces a collection \mathcal{T} of trees each of which is rooted at a core-vertex, where we call G *k -lean* if each of the rooted trees in \mathcal{T} is k -lean.

3.1 Chemical Graphs

To represent a chemical compound, we introduce a set of chemical elements such as H (hydrogen), C (carbon), O (oxygen), N (nitrogen) and so on. To distinguish a chemical element \mathbf{a} with multiple valences such as S (sulfur), we denote a chemical element \mathbf{a} with a valence i by $\mathbf{a}_{(i)}$, where we do not use such a suffix (i) for a chemical element \mathbf{a} with a unique valence. Let Λ be a set of chemical elements $\mathbf{a}_{(i)}$. For example, $\Lambda = \{\mathbf{H}, \mathbf{C}, \mathbf{O}, \mathbf{N}, \mathbf{P}, \mathbf{S}_{(2)}, \mathbf{S}_{(4)}, \mathbf{S}_{(6)}\}$. Let $\text{val} : \Lambda \rightarrow [1, 6]$ be a valence function. For example, $\text{val}(\mathbf{H}) = 1$, $\text{val}(\mathbf{C}) = 4$, $\text{val}(\mathbf{O}) = 2$, $\text{val}(\mathbf{P}) = 5$, $\text{val}(\mathbf{S}_{(2)}) = 2$, $\text{val}(\mathbf{S}_{(4)}) = 4$ and $\text{val}(\mathbf{S}_{(6)}) = 6$. For each chemical element $\mathbf{a} \in \Lambda$, let $\text{mass}(\mathbf{a})$ denote the mass of \mathbf{a} .

A chemical compound is represented by a *chemical graph* defined to be a tuple $\mathbb{C} = (H, \alpha, \beta)$ of a simple, connected undirected graph H and functions $\alpha : V(H) \rightarrow \Lambda$ and $\beta : E(H) \rightarrow [1, 3]$. The set of atoms and the set of bonds in the compound are represented by the vertex set $V(H)$ and the edge set $E(H)$, respectively. The chemical element assigned to a vertex $v \in V(H)$ is represented by $\alpha(v)$ and the bond-multiplicity between two adjacent vertices $u, v \in V(H)$ is represented by $\beta(e)$ of the edge $e = uv \in E(H)$. We say that two tuples $(H_i, \alpha_i, \beta_i), i = 1, 2$ are *isomorphic* if they admit an isomorphism ϕ , i.e., a bijection $\phi : V(H_1) \rightarrow V(H_2)$ such that $uv \in E(H_1), \alpha_1(u) = \mathbf{a}, \alpha_1(v) = \mathbf{b}, \beta_1(uv) = m \leftrightarrow \phi(u)\phi(v) \in E(H_2), \alpha_2(\phi(u)) = \mathbf{a}, \alpha_2(\phi(v)) = \mathbf{b}, \beta_2(\phi(u)\phi(v)) = m$. When H_i is rooted at a vertex $r_i, i = 1, 2$, $(H_i, \alpha_i, \beta_i), i = 1, 2$ are *rooted-isomorphic* (r-isomorphic) if they admit an isomorphism ϕ such that $\phi(r_1) = r_2$.

For a notational convenience, we use a function $\beta_{\mathbb{C}} : V(H) \rightarrow [0, 12]$ for a chemical graph $\mathbb{C} = (H, \alpha, \beta)$ such that $\beta_{\mathbb{C}}(u)$ means the sum of bond-multiplicities of edges incident to a vertex u ; i.e.,

$$\beta_{\mathbb{C}}(u) \triangleq \sum_{uv \in E(H)} \beta(uv) \text{ for each vertex } u \in V(H).$$

For each vertex $u \in V(H)$, define the *electron-degree* $\text{eledeg}_{\mathbb{C}}(u)$ to be

$$\text{eledeg}_{\mathbb{C}}(u) \triangleq \beta_{\mathbb{C}}(u) - \text{val}(\alpha(u)).$$

For each vertex $u \in V(H)$, let $\text{deg}_{\mathbb{C}}(u)$ denote the number of vertices adjacent to the vertex u in \mathbb{C} .

For a chemical graph $\mathbb{C} = (H, \alpha, \beta)$, let $V_{\mathbf{a}}(\mathbb{C})$, $\mathbf{a} \in \Lambda$ denote the set vertices $v \in V(H)$ such that $\alpha(v) = \mathbf{a}$ in \mathbb{C} and define the *hydrogen-suppressed chemical graph* $\langle \mathbb{C} \rangle$ to be the graph obtained from H by removing all the vertices $v \in V_{\mathbf{H}}(\mathbb{C})$.

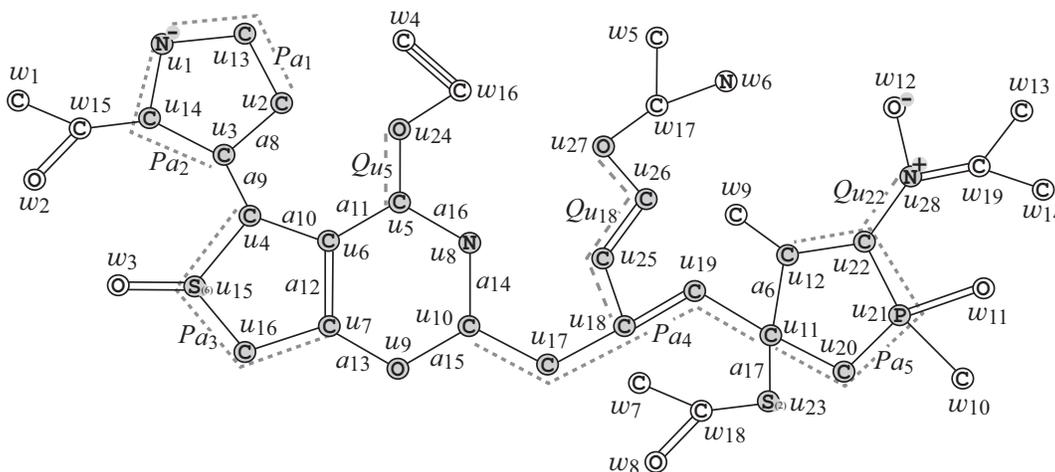


Figure 3: An illustration of a hydrogen-suppressed chemical graph $\langle \mathbb{C} \rangle$ obtained from a chemical graph \mathbb{C} with $r(\mathbb{C}) = 4$ by removing all the hydrogens, where for $\rho = 2$, $V^{\text{ex}}(\mathbb{C}) = \{w_i \mid i \in [1, 19]\}$ and $V^{\text{int}}(\mathbb{C}) = \{u_i \mid i \in [1, 28]\}$.

4 Two-layered Model

This section reviews the two-layered model introduced by Zhu et al. [20].

Let $\mathbb{C} = (H, \alpha, \beta)$ be a chemical graph and $\rho \geq 1$ be an integer, which we call a *branch-parameter*.

A *two-layered model* of \mathbb{C} is a partition of the hydrogen-suppressed chemical graph $\langle \mathbb{C} \rangle$ into an “interior” and an “exterior” in the following way. We call a vertex $v \in V(\langle \mathbb{C} \rangle)$ (resp., an edge $e \in E(\langle \mathbb{C} \rangle)$) of \mathbb{C} an *exterior-vertex* (resp., *exterior-edge*) if $\text{ht}(v) < \rho$ (resp., e is incident to an exterior-vertex) and denote the sets of exterior-vertices and exterior-edges by $V^{\text{ex}}(\mathbb{C})$ and $E^{\text{ex}}(\mathbb{C})$, respectively and denote $V^{\text{int}}(\mathbb{C}) = V(\langle \mathbb{C} \rangle) \setminus V^{\text{ex}}(\mathbb{C})$ and $E^{\text{int}}(\mathbb{C}) = E(\langle \mathbb{C} \rangle) \setminus E^{\text{ex}}(\mathbb{C})$, respectively. We call a vertex in $V^{\text{int}}(\mathbb{C})$ (resp., an edge in $E^{\text{int}}(\mathbb{C})$) an *interior-vertex* (resp., *interior-edge*).

The set $E^{\text{ex}}(\mathbb{C})$ of exterior-edges forms a collection of connected graphs each of which is regarded as a rooted tree T rooted at the vertex $v \in V(T)$ with the maximum $\text{ht}(v)$. Let $\mathcal{T}^{\text{ex}}(\langle \mathbb{C} \rangle)$ denote the set of these chemical rooted trees in $\langle \mathbb{C} \rangle$. The *interior* \mathbb{C}^{int} of \mathbb{C} is defined to be the subgraph $(V^{\text{int}}(\mathbb{C}), E^{\text{int}}(\mathbb{C}))$ of $\langle \mathbb{C} \rangle$.

Figure 3 illustrates an example of a hydrogen-suppressed chemical graph $\langle \mathbb{C} \rangle$. For a branch-parameter $\rho = 2$, the interior of the chemical graph $\langle \mathbb{C} \rangle$ in Figure 3 is obtained by removing the set of vertices with degree 1 $\rho = 2$ times; i.e., first remove the set $V_1 = \{w_1, w_2, \dots, w_{14}\}$ of vertices of degree 1 in $\langle \mathbb{C} \rangle$ and then remove the set $V_2 = \{w_{15}, w_{16}, \dots, w_{19}\}$ of vertices of degree 1 in $\langle \mathbb{C} \rangle - V_1$, where the removed vertices become the exterior-vertices of $\langle \mathbb{C} \rangle$.

For each interior-vertex $u \in V^{\text{int}}(\mathbb{C})$, let $T_u \in \mathcal{T}^{\text{ex}}(\langle \mathbb{C} \rangle)$ denote the chemical tree rooted at u (where possibly T_u consists of vertex u) and define the ρ -fringe-tree $\mathbb{C}[u]$ to be the chemical rooted tree obtained from T_u by putting back the hydrogens originally attached with T_u in \mathbb{C} . Let $\mathcal{T}(\mathbb{C})$ denote the set of ρ -fringe-trees $\mathbb{C}[u], u \in V^{\text{int}}(\mathbb{C})$. Figure 4 illustrates the set $\mathcal{T}(\mathbb{C}) = \{\mathbb{C}[u_i] \mid i \in [1, 28]\}$ of the 2-fringe-trees of the example \mathbb{C} with $\langle \mathbb{C} \rangle$ in Figure 3.

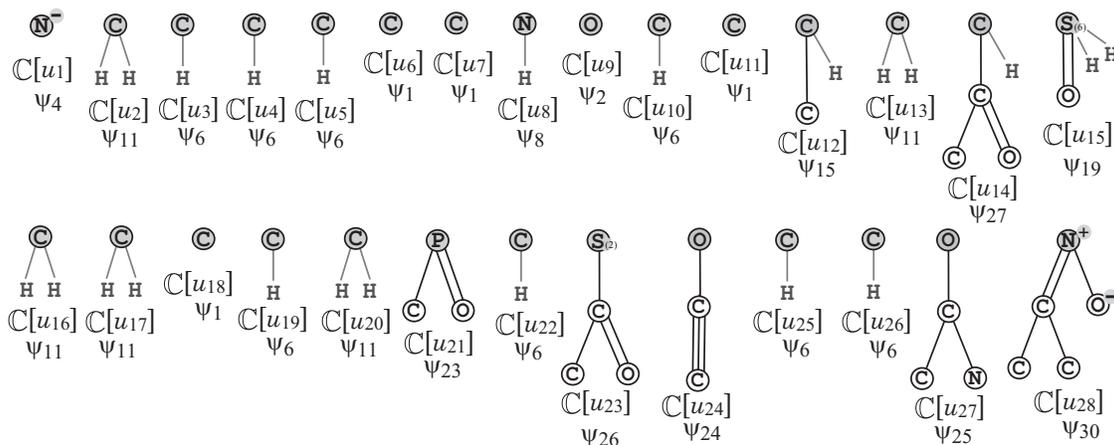


Figure 4: The set $\mathcal{T}(\mathbb{C})$ of 2-fringe-trees $\mathbb{C}[u_i], i \in [1, 28]$ of the example \mathbb{C} with $\langle \mathbb{C} \rangle$ in Figure 3, where the root of each tree is depicted with a gray circle and the hydrogens attached to non-root vertices are omitted in the figure.

Feature Function The feature of an interior-edge $e = uv \in E^{\text{int}}(\mathbb{C})$ such that $\alpha(u) = \mathbf{a}$, $\deg_{\langle \mathbb{C} \rangle}(u) = d$, $\alpha(v) = \mathbf{b}$, $\deg_{\langle \mathbb{C} \rangle}(v) = d'$ and $\beta(e) = m$ is represented by a tuple (ad, bd', m) , which is called the *edge-configuration* of the edge e , where we call the tuple $(\mathbf{a}, \mathbf{b}, m)$ the *adjacency-configuration* of the edge e .

For an integer K , a feature vector $f(\mathbb{C})$ of a chemical graph \mathbb{C} is defined by a *feature function* f that consists of K descriptors. We call \mathbb{R}^K the *feature space*.

Tanaka et al. [19] defined a feature vector $f(\mathbb{C}) \in \mathbb{R}^K$ to be a combination of the frequency of edge-configurations of the interior-edges and the frequency of chemical rooted trees among the set of chemical rooted trees $\mathbb{C}[u]$ over all interior-vertices u .

Topological Specification A topological specification is described as a set of the following rules proposed by Shi et al. [18] and modified by Tanaka et al. [19]:

- (i) a *seed graph* G_C as an abstract form of a target chemical graph \mathbb{C} ;
- (ii) a set \mathcal{F} of chemical rooted trees as candidates for a tree $\mathbb{C}[u]$ rooted at each interior-vertex u in \mathbb{C} ; and
- (iii) lower and upper bounds on the number of components in a target chemical graph such as chemical elements, double/triple bonds and the interior-vertices in \mathbb{C} .

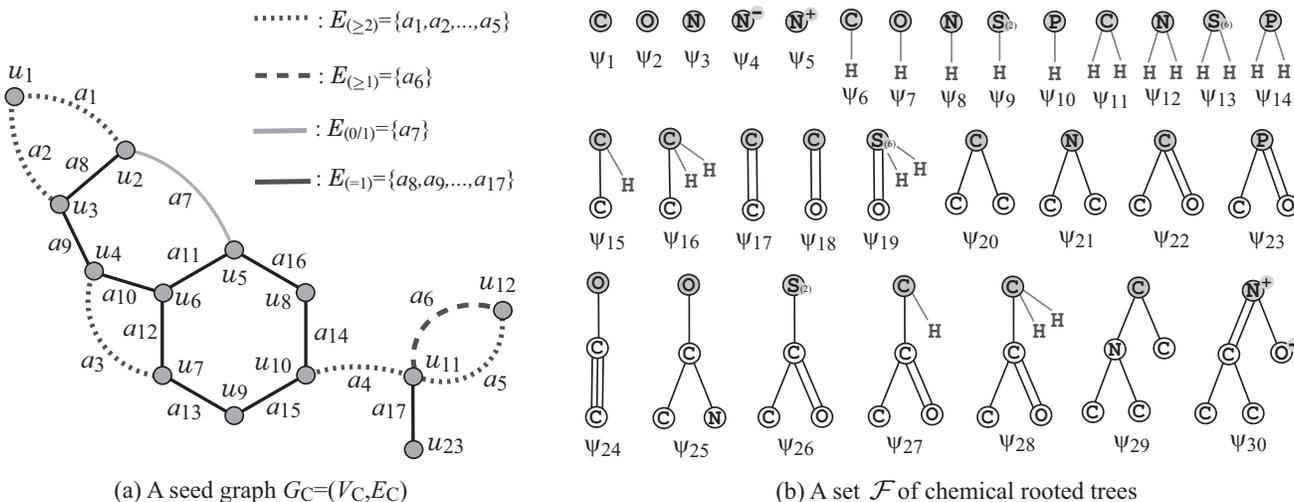


Figure 5: (a) An illustration of a seed graph G_C with $r(G_C) = 5$, where the vertices in V_C are depicted with gray circles, the edges in $E_{\geq 2}$ are depicted with dotted lines, the edges in $E_{\geq 1}$ are depicted with dashed lines, the edges in $E_{(0/1)}$ are depicted with gray bold lines and the edges in $E_{(=1)}$ are depicted with black solid lines; (b) A set $\mathcal{F} = \{\psi_1, \psi_2, \dots, \psi_{30}\} \subseteq \mathcal{F}(D_\pi)$ of 30 chemical rooted trees $\psi_i, i \in [1, 30]$, where the root of each tree is depicted with a gray circle, where the hydrogens attached to non-root vertices are omitted in the figure.

Figure 5(a) and (b) illustrate examples of a seed graph G_C and a set \mathcal{F} of chemical rooted trees, respectively. Given a seed graph G_C , the interior of a target chemical graph \mathbb{C} is constructed from G_C by replacing some edges $a = uv$ with paths P_a between the end-vertices u and v and by attaching new paths Q_v to some vertices v . For example, a chemical graph \mathbb{C} with $\langle \mathbb{C} \rangle$ in Figure 3 is constructed from the seed graph G_C in Figure 5(a) as follows.

- First replace five edges $a_1 = u_1u_2, a_2 = u_1u_3, a_3 = u_4u_7, a_4 = u_{10}u_{11}$ and $a_5 = u_{11}u_{12}$ in G_C with new paths $P_{a_1} = (u_1, u_{13}, u_2), P_{a_2} = (u_1, u_{14}, u_3), P_{a_3} = (u_4, u_{15}, u_{16}, u_7), P_{a_4} = (u_{10}, u_{17}, u_{18}, u_{19}, u_{11})$ and $P_{a_5} = (u_{11}, u_{20}, u_{21}, u_{22}, u_{12})$, respectively to obtain a subgraph G_1 of $\langle \mathbb{C} \rangle$.
- Next attach to this graph G_1 three new paths $Q_{u_5} = (u_5, u_{24}), Q_{u_{18}} = (u_{18}, u_{25}, u_{26}, u_{27})$ and $Q_{u_{22}} = (u_{22}, u_{28})$ to obtain the interior of $\langle \mathbb{C} \rangle$ in Figure 3.
- Finally attach to the interior 28 trees selected from the set \mathcal{F} and assign chemical elements and bond-multiplicities in the interior to obtain a chemical graph \mathbb{C} with $\langle \mathbb{C} \rangle$ in Figure 3. In Figure 4, $\psi_1 \in \mathcal{F}$ is selected for $\mathbb{C}[u_i], i \in \{6, 7, 11\}$. Similarly ψ_2 for $\mathbb{C}[u_9], \psi_4$ for $\mathbb{C}[u_1], \psi_6$ for $\mathbb{C}[u_i], i \in \{3, 4, 5, 10, 19, 22, 25, 26\}, \psi_8$ for $\mathbb{C}[u_8], \psi_{11}$ for $\mathbb{C}[u_i], i \in \{2, 13, 16, 17, 20\}, \psi_{15}$ for $\mathbb{C}[u_{12}], \psi_{19}$ for $\mathbb{C}[u_{15}], \psi_{23}$ for $\mathbb{C}[u_{21}], \psi_{24}$ for $\mathbb{C}[u_{24}], \psi_{25}$ for $\mathbb{C}[u_{27}], \psi_{26}$ for $\mathbb{C}[u_{23}], \psi_{27}$ for $\mathbb{C}[u_{14}]$

and ψ_{30} for $\mathbb{C}[u_{28}]$.

Our definition of a topological specification is analogous with the one by Tanaka et al. [19] except for a necessary modification due to the introduction of multiple valences of chemical elements, cations and anions (see Appendix B for a full description of topological specification).

5 Results

We implemented our method of Stages 1 to 5 for inferring chemical graphs under a given topological specification and conducted experiments to evaluate the computational efficiency. We executed the experiments on a PC with Processor: Core i7-9700 (3.0GHz; 4.7 GHz at the maximum) and Memory: 16 GB RAM DDR4. We used scikit-learn version 0.23.2 with Python 3.8.5 for executing linear regression with Lasso function or constructing an ANN. To solve an LP in Stage 3 or an MILP in Stage 4, we used CPLEX version 12.10.

Results on Phase 1. We implemented Stages 1, 2 and 3 in Phase 1 as follows.

We have conducted experiments of adjustive linear regression and for 37 chemical properties of monomers (resp., ten chemical properties of polymers) using the feature function [20] (resp., [25]) and we found that the test coefficient of determination R^2 of ALR exceeds 0.6 for the following 28 properties of monomers:

isotropic polarizability (ALPHA); boiling point (BP); critical pressure (CP); critical temperature (CT); heat capacity at 298.15K (CV); dissociation constants (DC); electron density on the most positive atom (EDPA); flash point (FP); energy difference between the highest and lowest unoccupied molecular orbitals (GAP); heat of atomization (HA); heat of combustion (HC); heat of formation (HF); energy of highest occupied molecular orbital (HOMO); heat of vaporization (HV); isobaric heat capacities in liquid phase (IHCL); isobaric heat capacities in solid phase (IHCS); K_vI-ats retention index (KVI), octanol/water partition coefficient (KOW); lipophilicity (LP); energy of lowest unoccupied molecular orbital (LUMO); melting point (MP); optical rotation (OPTR); refractive index (RF); solubility (SL); surface tension (SFT); internal energy at 0K (U0); viscosity (VIS); and vapor density (VD) and that the test coefficient of determination R^2 of ALR exceeds 0.8 for the following eight properties of polymers:

experimental amorphous density (AMD); characteristic ratio (CHAR); dielectric constant(DEC); heat capacity liquid (HCL); heat capacity solid (HCS); mol volume (MLV); refractive index (RFID); and glass transition (TG), where we include the result of property permittivity (PRM) for a comparison with Lasso linear regression and ANN.

We used data sets are provided by HSDB from PubChem [27] for CP, CT, DC, FP, HC, HV, KOW, OPTR, RF and VD. M. Jalali-Heravi and M. Fatemi [28] for EDPA and KVI, Roy and Saha [29] for BP, HA, HF and MP, Ramakrishnan et al. [30] for ALPHA, CV, LUMO and U0, Goussard et al. [31] for SFT, Goussard et al. [32] for VIS, R. Naef [33] for IHCL and IHCS, Xiao [34] for LP and Delaney [35] for SL. Properties ALPHA, CV, HOMO, LUMO and U0 share a common original data set D^* with more than 130,000 compounds, and we used a set D_π of 1,000 compounds randomly selected from D^* as a common data set of these four properties π in this experiment.

We used data sets of polymers provided by Bicerano [36], where we did not include any polymer

whose chemical formula could not be found by its name in the book. For property CHAR (resp., RFID), we remove the following polymer as an outlier from the original data set: ethyleneTerephthalate, oxy(2-methyl-6-phenyl-1_4-phenylene) and N-vinylCarbazole (resp., 2-decyl-1_4-butadiene).

Stage 1. We set a graph class \mathcal{G} to be the set of all chemical graphs with any graph structure, and set a branch-parameter ρ to be 2.

For each of the properties, we first select a set Λ of chemical elements and then collect a data set D_π on chemical graphs over the set Λ of chemical elements. To construct the data set D_π , we eliminated chemical compounds that do not satisfy one of the following: the graph is connected, the number of carbon atoms is at least four, and the number of non-hydrogen neighbors of each atom is at most 4.

Table 1 shows the size and range of data sets that we prepared for each chemical property in Stage 1, where we denote the following:

- Λ : the set of elements used in the data set D_π ; Λ is one of the following 12 sets: $\Lambda_1 = \{\text{H, C, O}\}$; $\Lambda_2 = \{\text{H, C, O, N}\}$; $\Lambda_3 = \{\text{H, C, O, S}\}$; $\Lambda_4 = \{\text{H, C, O, Si}_{(4)}\}$; $\Lambda_5 = \{\text{H, C, O, N, S}_{(2), \text{F}}\}$; $\Lambda_6 = \{\text{H, C}_{(2), \text{C}_{(3), \text{C}_{(4), \text{O, N}_{(2), \text{N}_{(3)}}\}$; $\Lambda_7 = \{\text{H, C, O, N, Cl, Pb}\}$; $\Lambda_8 = \{\text{H, C, O, N, S}_{(2), \text{S}_{(6), \text{Cl}}\}$; $\Lambda_9 = \{\text{H, C, O, N, S}_{(2), \text{S}_{(4), \text{S}_{(6), \text{Cl}}\}$; $\Lambda_{10} = \{\text{H, C}_{(2), \text{C}_{(3), \text{C}_{(4), \text{C}_{(5), \text{O, N}_{(1), \text{N}_{(2), \text{N}_{(3), \text{F}}\}$; $\Lambda_{11} = \{\text{H, C}_{(2), \text{C}_{(3), \text{C}_{(4), \text{O, N}_{(2), \text{N}_{(3), \text{S}_{(2), \text{S}_{(4), \text{S}_{(6), \text{Cl}}\}$; $\Lambda_{12} = \{\text{H, C, O, N, P}_{(2), \text{P}_{(5), \text{Cl}}\}$; $\Lambda_{13} = \{\text{H, C, O}_{(1), \text{O}_{(2), \text{N}}\}$; $\Lambda_{14} = \{\text{H, C, O, N, Cl}\}$; $\Lambda_{15} = \{\text{H, C, O, N, Cl, S}_{(2)}\}$; and $\Lambda_{16} = \{\text{H, C, O}_{(1), \text{O}_{(2), \text{N, Cl, Si}_{(4), \text{F}}\}$, where $\mathbf{a}_{(i)}$ for a chemical element \mathbf{a} and an integer $i \geq 1$ means that a chemical element \mathbf{a} with valence i .
- $|D_\pi|$: the size of data set D_π over Λ for the property π ;
- \underline{n}, \bar{n} : the minimum and maximum values of the number $n(\mathbb{C})$ of non-hydrogen atoms in compounds \mathbb{C} in D_π ;
- \underline{a}, \bar{a} : the minimum and maximum values of $a(\mathbb{C})$ for π over compounds \mathbb{C} in D_π ;
- $|\Gamma|$: the number of different edge-configurations of interior-edges over the compounds in D_π ;
- $|\mathcal{F}|$: the number of non-isomorphic chemical rooted trees in the set of all 2-fringe-trees in the compounds in D_π ; and
- K : the number of descriptors in a feature vector $f(\mathbb{C})$.

Stage 2. We used the feature function defined in our chemical model without suppressing hydrogen (see Appendix A for the detail). We standardize the range of each descriptor and the range $\{t \in \mathbb{R} \mid \underline{a} \leq t \leq \bar{a}\}$ of property values $a(\mathbb{C}), \mathbb{C} \in D_\pi$.

Stage 3. For each chemical property π , we select a penalty value λ_π for a constant λ in $\text{ALR}(\mathcal{X}, \lambda)$ by conducting linear regression as a preliminary experiment.

We conducted an experiment in Stage 3 to evaluate the performance of the prediction function based on cross-validation. For a property π , an execution of a *cross-validation* consists of five trials of constructing a prediction function as follows. First partition the data set D_π into five subsets $D^{(k)}, k \in [1, 5]$ randomly. For each $k \in [1, 5]$, the k -th trial constructs a prediction function $\eta^{(k)}$ by conducting a linear regression with the penalty term λ_π using the set $D_\pi \setminus D^{(k)}$ as a training data set. For each property, we executed ten cross-validations and we show the median of test $\text{R}^2(\eta^{(k)}, D^{(k)}), k \in [1, 5]$ over all ten cross-validations. Recall that a subset of descriptors is selected

in linear regression with ALR and let K' denote the average number of descriptors selected by ALR over all 50 trials.

Table 1: Results in Phase 1 for monomers.

π	Λ	$ D_\pi $	n, \bar{n}	\underline{a}, \bar{a}	$ \Gamma $	$ \mathcal{F} $	K	λ_π	K'	L-time	ALR	LLR	ANN
ALPHA	Λ_{10}	977	4, 9	50.9, 99.6	59	190	297	6.0e-4	88.4	3.00	0.953	0.961	0.888
BP	Λ_2	370	4, 67	-11.7, 470.0	22	130	184	5.3e-6	97.3	1.42	0.816	0.599	0.765
BP	Λ_8	444	4, 67	-11.7, 470.0	26	163	230	2.8e-6	112.6	2.02	0.832	0.663	0.720
CP	Λ_2	125	4, 63	$4.7 \times 10^{-6}, 5.52$	8	75	112	5.9e-3	17.4	0.15	0.650	0.445	0.694
CP	Λ_7	131	4, 63	$4.7 \times 10^{-6}, 5.52$	8	79	119	8.3e-3	12.3	0.12	0.690	0.555	0.727
CT	Λ_2	125	4, 63	56.1, 3607.5	8	76	113	3.1e-3	37.5	0.24	0.900	0.037	0.357
CT	Λ_7	132	4, 63	56.1, 3607.5	8	81	121	2.6e-3	43.0	0.28	0.895	0.048	0.356
CV	Λ_{10}	977	4, 9	19.2, 44.0	59	190	297	2.1e-4	143.9	4.57	0.966	0.970	0.911
DC	Λ_8	161	5, 44	0.5, 17.11	25	69	130	2.5e-3	45.3	0.35	0.602	0.574	0.622
EDPA	Λ_1	52	11, 16	0.80, 3.76	9	33	64	2.6e-3	19.0	0.06	0.999	0.999	0.992
FP	Λ_2	368	4, 67	-82.99, 300.0	20	131	183	6.0e-4	86.6	1.31	0.719	0.589	0.746
FP	Λ_8	424	4, 67	-82.99, 300.0	25	161	229	2.1e-4	109.4	1.92	0.684	0.571	0.745
GAP	Λ_{10}	977	4, 9	0.15, 0.41	59	190	297	1.6e-4	145.2	4.77	0.755	0.784	0.795
HA	Λ_3	115	4, 11	1100.6, 3009.6	8	83	115	1.1e-3	56.5	0.29	0.998	0.997	0.926
HC	Λ_2	255	4, 63	49.6, 35099.6	17	106	154	1.0e-3	69.9	0.74	0.986	0.946	0.848
HC	Λ_8	282	4, 63	49.6, 35099.6	21	118	177	1.6e-3	69.8	0.84	0.986	0.951	0.903
HF	Λ_1	82	4, 16	30.2, 94.8	5	50	74	8.4e-4	5.8	0.05	0.982	0.987	0.928
HOMO	Λ_{10}	977	4, 9	-0.11, 0.10	59	190	297	1.0e-4	158.7	4.95	0.689	0.841	0.689
HV	Λ_2	95	4, 16	19.12, 5193.1	12	63	105	1.6e-3	40.7	0.19	0.626	-13.7	-8.44
IHCL	Λ_2	770	4, 78	106.3, 1956.1	23	200	256	3.6e-5	126.3	3.24	0.987	0.986	0.974
IHCL	Λ_8	865	4, 78	106.3, 1956.1	29	246	316	6.0e-4	64.2	1.98	0.989	0.985	0.971
IHCS	Λ_6	581	5, 70	67.4, 1220.9	33	124	192	1.1e-5	86.5	1.72	0.971	0.985	0.971
IHCS	Λ_{11}	668	5, 70	67.4, 1220.9	40	140	228	1.0e-6	96.1	2.21	0.974	0.982	0.968
KVI	Λ_1	52	11, 16	1422.0, 1919.0	9	33	64	5.9e-3	18.8	0.05	0.838	0.677	0.727
KOW	Λ_2	684	4, 58	-7.5, 15.6	25	166	223	3.1e-4	119.2	3.13	0.964	0.953	0.952
KOW	Λ_9	899	4, 69	-7.5, 15.6	37	219	303	5.5e-4	141.4	4.95	0.952	0.927	0.937
LP	Λ_2	615	6, 60	-3.62, 6.84	32	116	186	3.1e-4	85.6	1.81	0.844	0.856	0.867
LP	Λ_9	936	6, 74	-3.62, 6.84	44	136	231	1.0e-4	108.4	3.37	0.807	0.840	0.859
LUMO	Λ_{10}	977	4, 9	-0.11, 0.10	59	190	297	6.0e-4	81.1	2.75	0.833	0.841	0.860
MP	Λ_2	467	4, 122	-185.33, 300.0	23	142	197	8.0e-4	95.7	1.78	0.831	0.810	0.799
MP	Λ_9	577	4, 122	-185.33, 300.0	32	176	255	1.1e-4	136.6	2.99	0.807	0.810	0.820
OPTR	Λ_2	147	5, 44	-117.0, 165.0	21	55	107	1.0e-3	30.2	0.24	0.876	0.825	0.919
OPTR	Λ_5	157	5, 69	-117.0, 165.0	25	62	123	1.1e-3	32.3	0.27	0.870	0.825	0.878
RF	Λ_2	166	4, 26	1.3326, 1.613	14	98	142	3.5e-3	25.4	0.24	0.685	0.619	0.521
SL	Λ_2	673	4, 55	-9.332, 1.11	27	154	217	1.0e-3	45.2	1.21	0.784	0.808	0.848
SL	Λ_9	915	4, 55	-11.6, 1.11	42	207	300	6.0e-4	73.4	2.33	0.828	0.808	0.861
SFT	Λ_4	247	5, 33	12.3, 45.1	11	91	128	1.0e-3	63.1	0.67	0.847	0.927	0.859
U0	Λ_{10}	977	4, 9	-570.6, -272.84	59	190	297	1.1e-3	69.7	2.40	0.995	0.999	0.890
VIS	Λ_4	282	5, 36	-0.64, 1.63	12	88	126	2.1e-3	23.2	0.37	0.911	0.893	0.929
VD	Λ_2	474	4, 30	0.7, 20.6	21	160	214	1.0e-4	119.1	2.24	0.985	0.927	0.912
VD	Λ_{12}	551	4, 30	0.7, 20.6	24	191	256	6.0e-4	101.1	2.28	0.980	0.942	0.889

Tables 1 and 2 show the results on Stages 2 and 3 for the properties on monomers and polymers, respectively, where we denote the following:

- λ_π : the penalty value in the Lasso function selected for a property π , where $a \ e \ b$ means $a \times 10^b$;
- K' : the average of the number of descriptors selected in the linear regression over all 50 trials in ten cross-validations;
- L-time: the average time (sec.) to construct a prediction function with ALR by solving an LP with $O(|D_\pi| + K)$ variables and constraints over all 50 trials in ten cross-validations;
- ALR: the median of test R^2 over all 50 trials in ten cross-validations for prediction functions constructed with ALR;
- LLR: the median of test R^2 over all 50 trials in ten cross-validations for prediction functions constructed with Lasso linear regression; and
- ANN: the median of test R^2 over all 50 trials in ten cross-validations for prediction functions constructed with ANN (see [26] for the details of the implementation).

Table 2: Results in Phase 1 for polymers.

π	Λ	$ D_\pi $	\underline{n}, \bar{n}	\underline{a}, \bar{a}	$ \Gamma $	$ \mathcal{F} $	K	λ_π	K'	L-time	ALR	LLR	ANN
AMD	Λ_2	86	4, 45	0.838, 1.34	28	25	83	1.5E-3	22.2	0.09	0.933	0.914	0.885
AMD	Λ_{15}	93	4, 45	0.838, 1.45	31	30	94	2.5E-3	24.4	0.10	0.917	0.918	0.823
CHAR	Λ_1	24	4, 18	5.5, 13.2	15	15	56	5.0E-5	15.0	0.02	0.904	0.650	0.616
CHAR	Λ_2	27	4, 18	5.5, 13.2	22	17	67	5.9E-3	17.8	0.03	0.835	0.431	0.641
DEC	Λ_{15}	37	4, 22	2.13, 3.4	22	19	72	3.1E-3	18.2	0.04	0.918	0.761	0.641
HcL	Λ_2	52	4, 25	105.7, 677.8	22	17	67	2.1E-3	19.3	0.06	0.996	0.990	0.969
HcL	Λ_8	55	4, 32	105.7, 678.1	27	20	81	2.6E-3	17.3	0.05	0.992	0.987	0.970
HcS	Λ_2	54	4, 45	84.5, 720.5	26	20	75	1.0E-3	23.5	0.07	0.963	0.968	0.893
HcS	Λ_8	59	4, 45	84.5, 720.5	32	24	92	4.1E-4	30.7	0.09	0.983	0.961	0.880
MLV	Λ_2	86	4, 45	60.7, 466.6	28	25	83	1.0E-3	22.2	0.10	0.998	0.996	0.931
MLV	Λ_{15}	93	4, 45	60.7, 466.6	31	30	94	5.6E-4	27.5	0.09	0.997	0.994	0.894
PRM	Λ_2	112	4, 45	2.23, 4.91	25	15	69	2.6E-4	14.5	0.09	0.505	0.801	0.801
PRM	Λ_{14}	131	4, 45	2.23, 4.91	25	17	73	5.9E-3	13.9	0.09	0.489	0.784	0.735
RFID	Λ_{13}	91	4, 29	1.4507, 1.683	26	35	96	3.1E-4	32.9	0.15	0.953	0.852	0.871
RFID	Λ_{16}	124	4, 29	1.339, 1.683	32	50	124	2.1E-3	37.5	0.21	0.956	0.832	0.891
TG	Λ_2	204	4, 58	171, 673	32	36	101	2.5E-3	26.9	0.23	0.923	0.902	0.883
TG	Λ_8	232	4, 58	171, 673	36	43	118	1.1E-3	38.6	0.54	0.927	0.894	0.881

From Tables 1 and 2, we see that ALR performs well for most of the properties in our experiments, The performance by ALR is inferior to that by LLR or ANN for some properties such as GAP, HOMO, LUMO, OPTR, SL, SFT and PRM, whereas ALR outperforms LLR and ANN for properties BP, CT, HV, KVI, VD, CHAR, RFID and TG. It should be noted that ALR drastically improves the result for properties CT and HV.

Results on Phase 2. To execute Stages 4 and 5 in Phase 2, we used a set of seven instances I_a , $I_b^i, i \in [1, 4]$, I_c and I_d based on the seed graphs prepared by Zhu et al. [20]. We here present their seed graphs G_C (see Appendix B for the details of I_a and Appendix C for the details of $I_b^i, i \in [1, 4]$, I_c and I_d). The seed graph G_C of I_a is given by the graph in Figure 5(a). The seed graph G_C^1 of

I_b^1 (resp., $G_C^i, i = 2, 3, 4$ of $I_b^i, i = 2, 3, 4$) is illustrated in Figure 6.

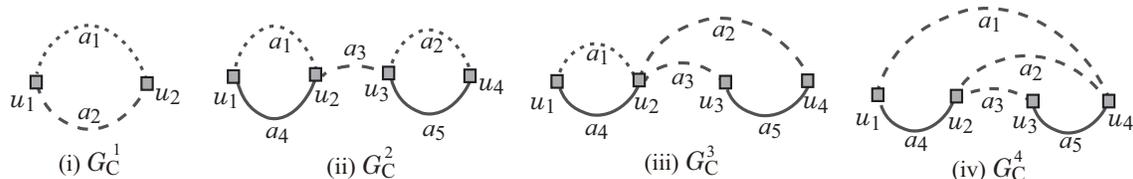


Figure 6: Seed graphs: (i) G_C^1 for I_b^1 and I_d ; (ii) G_C^2 for I_b^2 ; (iii) G_C^3 for I_b^3 ; (iv) G_C^4 for I_b^4 .

Instance I_c has been introduced in order to infer a chemical graph \mathbb{C}^\dagger such that the core of \mathbb{C}^\dagger is equal to the core of chemical graph \mathbb{C}_A : CID 24822711 in Figure 7(a) and the frequency of each edge-configuration in the non-core of \mathbb{C}^\dagger is equal to that of chemical graph \mathbb{C}_B : CID 59170444 in Figure 7(b). This means that the seed graph G_C of I_c is the core of \mathbb{C}_A which is indicated by a shaded area in Figure 7(a).

Instance I_d has been introduced in order to infer a chemical monocyclic graph \mathbb{C}^\dagger such that the frequency vector of edge-configurations in \mathbb{C}^\dagger is a vector obtained by merging those of chemical graphs \mathbb{C}_A : CID 10076784 and \mathbb{C}_B : CID 44340250 in Figure 7(c) and (d), respectively. The seed graph G_C of I_d is given by G_C^1 in Figure 6(i).

Stage 4. We executed Stage 4 for two properties $\pi \in \{\text{CT}, \text{HV}\}$.

For the MILP formulation $\mathcal{M}(x, y; \mathcal{C}_1)$ in the framework, we use the prediction function $\eta_{w,b}$ constructed in Stage 3. Tables 3 and 4 show the computational results of the experiment in Stage 4 for the two properties, where we denote the following:

- $\underline{y}^*, \bar{y}^*$: lower and upper bounds $\underline{y}^*, \bar{y}^* \in \mathbb{R}$ on the value $a(\mathbb{C})$ of a chemical graph \mathbb{C} to be inferred;
- $\#v$ (resp., $\#c$): the number of variables (resp., constraints) in the MILP in Stage 4;
- I-time: the time (sec.) to solve the MILP in Stage 4;
- n : the number $n(\mathbb{C}^\dagger)$ of non-hydrogen atoms in the chemical graph \mathbb{C}^\dagger inferred in Stage 4;
- n^{int} : the number $n^{\text{int}}(\mathbb{C}^\dagger)$ of interior-vertices in the chemical graph \mathbb{C}^\dagger inferred in Stage 4; and
- η : the predicted property value $\eta(f(\mathbb{C}^\dagger))$ of the chemical graph \mathbb{C}^\dagger inferred in Stage 4.

From Tables 3 and 4, we observe that an instance with around 50 non-hydrogen atoms is solved in at most around 1000 seconds. Instances with the size in Stage 4 were solved at most around 100 seconds in the experiments due to Azam et al. [26]. The computation time for Stage 4 increased possibly because we included a new set of constraints for representing activation functions ϕ_j in the MILP $\mathcal{M}(x, y; \mathcal{C}_1)$ of the framework (see Appendix D.11 for the details). Note that for the properties $\pi \in \{\text{CT}, \text{HV}\}$, no prediction functions constructed for the framework [20, 26] performed well.

Stage 5. We executed Stage 5 to generate a more number of target chemical graphs \mathbb{C}^* , where we call a chemical graph \mathbb{C}^* a *chemical isomer* of a target chemical graph \mathbb{C}^\dagger of a topological specification σ if $f(\mathbb{C}^*) = f(\mathbb{C}^\dagger)$ and \mathbb{C}^* also satisfies the same topological specification σ . We computed chemical isomers \mathbb{C}^* of each target chemical graph \mathbb{C}^\dagger inferred in Stage 4. We execute an algorithm for generating chemical isomers of \mathbb{C}^\dagger up to 100 when the number of all chemical

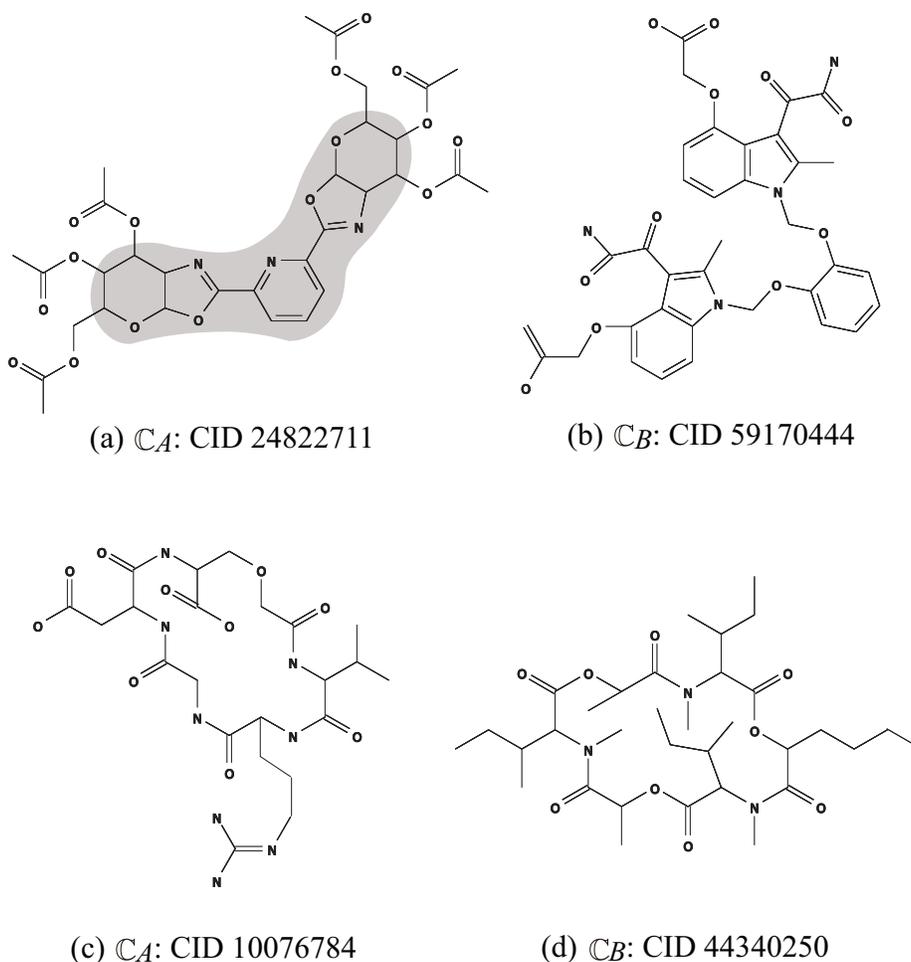


Figure 7: Chemical compounds: (a) CID 24822711; (b) CID 59170444; (c) CID 10076784; (d) CID 44340250, where hydrogens are omitted.

Table 3: Results of Stages 4 and 5 for Cr.

inst.	$\underline{y}^*, \bar{y}^*$	#v	#c	I-time	n	n^{int}	η	D-time	C-LB	#C
I_a	1180, 1220	16190	16895	8.4	42	25	1197.0	0.0619	1	1
I_b^1	1080, 1120	86470	85993	119.2	35	14	1118.0	0.266	84	84
I_b^2	1980, 2020	113055	114138	365.5	45	30	1986.9	2.62	4.4×10^6	100
I_b^3	1930, 1970	112745	114086	555.4	45	25	1950.9	4.57	1.1×10^6	100
I_b^4	1630, 1670	112455	114054	661.2	47	29	1663.8	0.171	5.5×10^5	100
I_c	1980, 2020	14794	15226	7.2	50	34	2011.0	0.0161	1	1
I_d	1430, 1470	12790	14104	11.9	45	23	1447.5	0.164	5184	100

isomers exceeds 100. Such an algorithm can be obtained from the dynamic programming proposed by Tanaka et al. [19] with a slight modification. The algorithm first decomposes \mathbb{C}^\dagger into a set of

Table 4: Results of Stages 4 and 5 for HV.

inst.	$\underline{y}^*, \bar{y}^*$	#v	#c	I-time	n	n^{int}	η	D-time	C-LB	#C
I_a	145, 150	16190	16879	24.9	37	23	147.3	0.0632	2	2
I_b^1	190, 195	14179	13334	146.6	35	12	190.2	0.121	30	30
I_b^2	290, 295	17986	18547	188.8	46	25	294.2	0.154	604	100
I_b^3	165, 170	17683	18495	1167.2	45	25	166.8	36.8	7.5×10^6	100
I_b^4	250, 255	17400	18463	313.7	50	25	251.8	0.166	2208	100
I_c	285, 290	14838	15254	102.5	50	32	286.8	0.016	1	1
I_d	245, 250	12828	14126	351.9	40	23	249.2	5.53	3.9×10^5	100

acyclic chemical graphs, next replace each acyclic chemical graph T with another acyclic chemical graph T' that admits the same feature vector as that of T and finally assemble the resulting acyclic chemical graphs into a chemical isomer \mathbb{C}^* of \mathbb{C}^\dagger . The algorithm can compute a lower bound on the total number of all chemical isomers \mathbb{C}^\dagger without generating all of them.

Tables 3 and 4 show the computational results of the experiment in Stage 5 for properties CT and HV, where we denote the following:

- D-time: the running time (sec.) to execute the dynamic programming algorithm in Stage 5 to compute a lower bound on the number of all chemical isomers \mathbb{C}^* of \mathbb{C}^\dagger and generate all (or up to 100) chemical isomers \mathbb{C}^* ;
- C-LB: a lower bound on the number of all chemical isomers \mathbb{C}^* of \mathbb{C}^\dagger ; and
- #C: the number of all (or up to 100) chemical isomers \mathbb{C}^* of \mathbb{C}^\dagger generated in Stage 5.

From Tables 3 and 4, we observe that the running time for generating up to 100 target chemical graphs in Stage 5 is less than around 5 second for many cases. For some chemical graph \mathbb{C}^\dagger , no chemical isomer was found by our algorithm. For such a case, we may use a method of generating other chemical graphs \mathbb{C}^\dagger in Stage 4 by solving the MILP again with additional linear constraints in a systematical way (see [26] for the details).

6 Concluding Remarks

In this paper, we proposed a new machine learning method, adjustive linear regression (ALR), which has the following feature: (i) ALR is an extension of the Lasso linear regression except for the definition of error functions; (ii) ALR is a special case of an ANN except that a choice of activation functions is also optimized differently from the standard ANNs and the definition of error functions; and (iii) ALR can be executed exactly by solving the equivalent linear program with $O(m+K)$ variables and constraints for a set of m data with K descriptors. Even though ALR is a special case of an ANN with non-linear activation functions, we still can read the relationship between cause and effect from a prediction function due to the simple structure of ALR.

When the size of the original data set is large, we can choose a smaller subset for constructing a prediction function. For example, we chose 1,000 compounds from the data set of 130,000

compounds provided by Ramakrishnan et al. [30] for conducting experiments in Stage 3.

In this paper, we used a quadratic function for a set Ψ of activation functions ϕ . We can use many different functions such as sigmoid function and ramp functions, where the non-linearity of a function does not affect to derive a linear program for ALR.

References

- [1] Lo, Y-C., Rensi, S.E., Torng, W., Altman, R.B.: Machine learning in chemoinformatics and drug discovery. *Drug Discovery Today* **23**, 1538–1546 (2018)
- [2] Tetko, I.V., Engkvist, O.: From Big Data to Artificial Intelligence: chemoinformatics meets new challenges. *J. Cheminformatics* **12**, 74 (2020)
- [3] Ghasemi, F., Mehridehnavi, A., Pérez-Garrido, A., Pérez-Sánchez, H.: Neural network and deep-learning algorithms used in QSAR studies: merits and drawbacks. *Drug Discovery Today* **23**, 1784–1790 (2018)
- [4] Miyao, T., Kaneko, H., Funatsu, K.: Inverse QSPR/QSAR analysis for chemical structure generation (from y to x). *J. Chem. Inf. Model.* **56**, 286–299 (2016)
- [5] Ikebata, H., Hongo, K., Isomura, T., Maezono, R., Yoshida, R.: Bayesian molecular design with a chemical language model. *J. Comput. Aided Mol. Des.* **31**, 379–391 (2017)
- [6] Rupakheti, C., Virshup, A., Yang, W., Beratan, D.N.: Strategy to discover diverse optimal molecules in the small molecule universe. *J. Chem. Inf. Model.* **55**, 529–537 (2015)
- [7] Bohacek, R.S., McMartin, C., Guida, W.C.: The art and practice of structure-based drug design: A molecular modeling perspective. *Med. Res. Rev.* **16**, 3–50 (1996)
- [8] Kipf, T. N., Welling, M.: Semi-supervised classification with graph convolutional networks, arXiv:1609.02907 (2016)
- [9] Gómez-Bombarelli, R., Wei, J.N., Duvenaud, D., Hernández-Lobato, J.M., Sánchez-Lengeling, B., Sheberla, D., Aguilera-Iparraguirre, J., Hirzel, T.D., Adams, R.P., Aspuru-Guzik, A.: Automatic chemical design using a data-driven continuous representation of molecules. *ACS Cent. Sci.* **4**, 268–276 (2018)
- [10] Segler, M.H.S., Kogej, T., Tyrchan, C., Waller, M.P.: Generating focused molecule libraries for drug discovery with recurrent neural networks. *ACS Cent. Sci.* **4**, 120–131 (2017)
- [11] Yang, X., Zhang, J., Yoshizoe, K., Terayama, K., Tsuda, K.: ChemTS: an efficient python library for de novo molecular generation. *STAM* **18**, 972–976 (2017)
- [12] Kusner, M.J., Paige, B., Hernández-Lobato, J.M.: Grammar variational autoencoder. *Proc. of the 34th International Conference on Machine Learning-Volume 70*, 1945–1954 (2017)

- [13] De Cao, N., Kipf, T.: MolGAN: An implicit generative model for small molecular graphs. arXiv:1805.11973 (2018)
- [14] Madhawa, K., Ishiguro, K., Nakago, K., Abe, M.: GraphNVP: an invertible flow model for generating molecular graphs. arXiv:1905.11600 (2019)
- [15] Shi, C., Xu, M., Zhu, Z., Zhang, W., Zhang, M., Tang, J.: GraphAF: a flow-based autoregressive model for molecular graph generation. arXiv:2001.09382 (2020)
- [16] Akutsu, T., Nagamochi, H.: A mixed integer linear programming formulation to artificial neural networks. Proc. of the 2nd Int. Conf. on Information Science and Systems, 215–220 (2019)
- [17] Azam, N. A., Chiewvanichakorn, R., Zhang, F., Shurbevski, A., Nagamochi, H., Akutsu, T.: A method for the inverse QSAR/QSPR based on artificial neural networks and mixed integer linear programming. Proc. of the 13th International Joint Conference on Biomedical Engineering Systems and Technologies – Volume 3: BIOINFORMATICS, 101–108 (2020)
- [18] Shi, Y., Zhu, J., Azam, N. A., Haraguchi, K., Zhao, L., Nagamochi, H., Akutsu, T.: An inverse QSAR method based on a two-layered model and integer programming. International Journal of Molecular Sciences. **22**, 2847 (2021)
- [19] Tanaka, K., Zhu, J., Azam, N. A., Haraguchi, K., Zhao, L., Nagamochi, H., Akutsu, T.: An inverse QSAR method based on decision tree and integer programming, The 17th International Conference on Intelligent Computing, August 12-15, 2021, in Shenzhen, China, In: Huang D.S., Jo K.H., Li J., Gribova V., Hussain A. (eds) Intelligent Computing Theories and Application, ICIC 2021, Lecture Notes in Computer Science, vol. 12837. Springer, Cham.
- [20] Zhu, J., Azam, N. A., Haraguchi, K., Zhao, L., Nagamochi, H., Akutsu, T.: A method for molecular design based on linear regression and integer programming. arXiv: 2107.02381 (2021) <http://arxiv.org/abs/2107.02381>
- [21] Hoerl, A., Kennard, R.: Ridge Regression: Biased Estimation for Nonorthogonal Problems. Technometrics, **12**(1), 55–67 (1970)
- [22] Hoerl, A., Kennard, R.: Ridge Regression: Applications to Nonorthogonal Problems. Technometrics, **12**(1), 69–82 (1970)
- [23] Tibshirani, R.: Regression shrinkage and selection via the lasso. J. R. Statist. Soc. B **58**, 267–288 (1996)
- [24] Zou, H., and Hastie, T.: Regularization and variable selection via the elastic net. J. Royal Statistical Society: Series B (Statistical Methodology), **67**(2), 301-320 (2005)
- [25] Ido, R., Cao, S., Zhu, J., Azam, N. A., Haraguchi, K., Zhao, L., Nagamochi, H., Akutsu, T.: A method for inferring polymers based on linear regression and integer programming. Department of Applied Mathematics and Physics, Kyoto University, Technical Report, TR: 2021-001 <http://www.amp.i.kyoto-u.ac.jp/tecrep/> (2021)

- [26] Azam, N. A., Zhu, J., Haraguchi, K., Zhao, L., Nagamochi, H., Akutsu, T.: Molecular design based on artificial neural networks, integer programming and grid neighbor search. arXiv: 2108.10266 (2021)
- [27] Annotations from HSDB (on pubchem): <https://pubchem.ncbi.nlm.nih.gov/>. Accessed: 2021-8-26.
- [28] Jalali-Heravi, M., Fatemi, M. : H.Artificial neural network modeling of Kovats retention indices for noncyclic and monocyclic terpenes (2001) [https://doi.org/10.1016/S0021-9673\(00\)01274-7/](https://doi.org/10.1016/S0021-9673(00)01274-7/)
- [29] Roy, K., Saha, A.: Comparative QSPR studies with molecular connectivity, molecular negentropy and TAU indices (2003) <https://doi.org/10.1007/s00894-003-0135-z/>
- [30] Ramakrishnan, R., Dral, P., Rupp, M., Anatole von Lilienfeld, O.: Quantum chemistry structures and properties of 134 kilo molecules (2014) <https://doi.org/10.6084/m9.figshare.c.978904.v5>. Accessed: 2021-8-26.
- [31] Goussard, V., Duprat, F., Gerbaud, V., Ploix, J.-J., Dreyfus, G., Nardello-Rataj, V., Aubry, J.-M.: Predicting the surface tension of liquids: comparison of four modeling approaches and application to cosmetic oils, *J. Chem. Inf. Model.*, 57, 12, 29862995 (2017) <https://pubs.acs.org/doi/full/10.1021/acs.jC1m.7b00512>
- [32] Goussard, V., Francois Duprat F., Ploix, J.-L., Dreyfus, G., Nardello-Rataj, V., Aubry, J.-M.: A new machine-learning tool for fast estimation of liquid viscosity. application to cosmetic oils, *J. Chem. Inf. Model.*, 60, 4, 20122023 (2020) <https://pubs.acs.org/doi/10.1021/acs.jcim.0c00083>
- [33] Naef, R.: Calculation of the isobaric heat capacities of the liquid and solid phase of organic compounds at and around 298.15 K based on their “ true ” molecular volume. *Molecules*, 24 (8) (2019), <https://www.mdpi.com/1420-3049/24/8/1626/>
- [34] Xiao, N.: Lipophilicity Dataset - logD7.4 of 1,130 Compounds (2017) <https://doi.org/10.6084/m9.figshare.5596750>. Accessed: 2021-8-26.
- [35] Delaney, J.S.: ESOL: Estimating Aqueous Solubility Directly from Molecular Structure (2019) https://figshare.com/articles/dataset/ESOL_Estimating_Aqueous_Solubility_Directly_from_Molecular_Structure/7944677/1. Accessed: 2021-8-26.
- [36] Bicerano, J.: Prediction of Polymer Properties. 3rd Edition, Revised and Expanded. CRC Press (2002)

Appendix

A A Full Description of Descriptors

Associated with the two functions α and β in a chemical graph $\mathbb{C} = (H, \alpha, \beta)$, we introduce functions $ac : V(E) \rightarrow (\Lambda \setminus \{\mathbf{H}\}) \times (\Lambda \setminus \{\mathbf{H}\}) \times [1, 3]$, $cs : V(E) \rightarrow (\Lambda \setminus \{\mathbf{H}\}) \times [1, 6]$ and $ec : V(E) \rightarrow ((\Lambda \setminus \{\mathbf{H}\}) \times [1, 6]) \times ((\Lambda \setminus \{\mathbf{H}\}) \times [1, 6]) \times [1, 3]$ in the following.

To represent a feature of the exterior of \mathbb{C} , a chemical rooted tree in $\mathcal{T}(\mathbb{C})$ is called a *fringe-configuration* of \mathbb{C} .

We also represent leaf-edges in the exterior of \mathbb{C} . For a leaf-edge $uv \in E(\langle \mathbb{C} \rangle)$ with $\deg_{\langle \mathbb{C} \rangle}(u) = 1$, we define the *adjacency-configuration* of e to be an ordered tuple $(\alpha(u), \alpha(v), \beta(uv))$. Define

$$\Gamma_{ac}^{\text{lf}} \triangleq \{(\mathbf{a}, \mathbf{b}, m) \mid \mathbf{a}, \mathbf{b} \in \Lambda, m \in [1, \min\{\text{val}(\mathbf{a}), \text{val}(\mathbf{b})\}]\}$$

as a set of possible adjacency-configurations for leaf-edges.

To represent a feature of an interior-vertex $v \in V^{\text{int}}(\mathbb{C})$ such that $\alpha(v) = \mathbf{a}$ and $\deg_{\langle \mathbb{C} \rangle}(v) = d$ (i.e., the number of non-hydrogen atoms adjacent to v is d) in a chemical graph $\mathbb{C} = (H, \alpha, \beta)$, we use a pair $(\mathbf{a}, d) \in (\Lambda \setminus \{\mathbf{H}\}) \times [1, 4]$, which we call the *chemical symbol* $cs(v)$ of the vertex v . We treat (\mathbf{a}, d) as a single symbol ad , and define Λ_{dg} to be the set of all chemical symbols $\mu = ad \in (\Lambda \setminus \{\mathbf{H}\}) \times [1, 4]$.

We define a method for featuring interior-edges as follows. Let $e = uv \in E^{\text{int}}(\mathbb{C})$ be an interior-edge $e = uv \in E^{\text{int}}(\mathbb{C})$ such that $\alpha(u) = \mathbf{a}$, $\alpha(v) = \mathbf{b}$ and $\beta(e) = m$ in a chemical graph $\mathbb{C} = (H, \alpha, \beta)$. To feature this edge e , we use a tuple $(\mathbf{a}, \mathbf{b}, m) \in (\Lambda \setminus \{\mathbf{H}\}) \times (\Lambda \setminus \{\mathbf{H}\}) \times [1, 3]$, which we call the *adjacency-configuration* $ac(e)$ of the edge e . We introduce a total order $<$ over the elements in Λ to distinguish between $(\mathbf{a}, \mathbf{b}, m)$ and $(\mathbf{b}, \mathbf{a}, m)$ ($\mathbf{a} \neq \mathbf{b}$) notationally. For a tuple $\nu = (\mathbf{a}, \mathbf{b}, m)$, let $\bar{\nu}$ denote the tuple $(\mathbf{b}, \mathbf{a}, m)$.

Let $e = uv \in E^{\text{int}}(\mathbb{C})$ be an interior-edge $e = uv \in E^{\text{int}}(\mathbb{C})$ such that $cs(u) = \mu$, $cs(v) = \mu'$ and $\beta(e) = m$ in a chemical graph $\mathbb{C} = (H, \alpha, \beta)$. To feature this edge e , we use a tuple $(\mu, \mu', m) \in \Lambda_{\text{dg}} \times \Lambda_{\text{dg}} \times [1, 3]$, which we call the *edge-configuration* $ec(e)$ of the edge e . We introduce a total order $<$ over the elements in Λ_{dg} to distinguish between (μ, μ', m) and (μ', μ, m) ($\mu \neq \mu'$) notationally. For a tuple $\gamma = (\mu, \mu', m)$, let $\bar{\gamma}$ denote the tuple (μ', μ, m) .

Let π be a chemical property for which we will construct a prediction function η from a feature vector $f(\mathbb{C})$ of a chemical graph \mathbb{C} to a predicted value $y \in \mathbb{R}$ for the chemical property of \mathbb{C} .

We first choose a set Λ of chemical elements and then collect a data set D_π of chemical compounds C whose chemical elements belong to Λ , where we regard D_π as a set of chemical graphs \mathbb{C} that represent the chemical compounds C in D_π . To define the interior/exterior of chemical graphs $\mathbb{C} \in D_\pi$, we next choose a branch-parameter ρ , where we recommend $\rho = 2$.

Let $\Lambda^{\text{int}}(D_\pi) \subseteq \Lambda$ (resp., $\Lambda^{\text{ex}}(D_\pi) \subseteq \Lambda$) denote the set of chemical elements used in the set $V^{\text{int}}(\mathbb{C})$ of interior-vertices (resp., the set $V^{\text{ex}}(\mathbb{C})$ of exterior-vertices) of \mathbb{C} over all chemical graphs $\mathbb{C} \in D_\pi$, and $\Gamma^{\text{int}}(D_\pi)$ denote the set of edge-configurations used in the set $E^{\text{int}}(\mathbb{C})$ of interior-edges in \mathbb{C} over all chemical graphs $\mathbb{C} \in D_\pi$. Let $\mathcal{F}(D_\pi)$ denote the set of chemical rooted trees ψ r-isomorphic to a chemical rooted tree in $\mathcal{T}(\mathbb{C})$ over all chemical graphs $\mathbb{C} \in D_\pi$, where possibly a chemical rooted tree $\psi \in \mathcal{F}(D_\pi)$ consists of a single chemical element $\mathbf{a} \in \Lambda \setminus \{\mathbf{H}\}$.

We define an integer encoding of a finite set A of elements to be a bijection $\sigma : A \rightarrow [1, |A|]$, where we denote by $[A]$ the set $[1, |A|]$ of integers. Introduce an integer coding of each of the sets $\Lambda^{\text{int}}(D_\pi)$, $\Lambda^{\text{ex}}(D_\pi)$, $\Gamma^{\text{int}}(D_\pi)$ and $\mathcal{F}(D_\pi)$. Let $[\mathbf{a}]^{\text{int}}$ (resp., $[\mathbf{a}]^{\text{ex}}$) denote the coded integer of an element $\mathbf{a} \in \Lambda^{\text{int}}(D_\pi)$ (resp., $\mathbf{a} \in \Lambda^{\text{ex}}(D_\pi)$), $[\gamma]$ denote the coded integer of an element γ in $\Gamma^{\text{int}}(D_\pi)$ and $[\psi]$ denote an element ψ in $\mathcal{F}(D_\pi)$.

We assume that a chemical graph \mathbb{C} treated in this paper satisfies $\deg_{\langle \mathbb{C} \rangle}(v) \leq 4$ in the hydrogen-suppressed graph $\langle \mathbb{C} \rangle$.

In our model, we use an integer $\text{mass}^*(\mathbf{a}) = \lfloor 10 \cdot \text{mass}(\mathbf{a}) \rfloor$, for each $\mathbf{a} \in \Lambda$.

We define the *feature vector* $f(\mathbb{C})$ of a chemical graph $\mathbb{C} = (H, \alpha, \beta) \in D_\pi$ to be a vector that consists of the following non-negative integer descriptors $\text{dcp}_i(\mathbb{C})$, $i \in [1, K]$, where $K = 14 + |\Lambda^{\text{int}}(D_\pi)| + |\Lambda^{\text{ex}}(D_\pi)| + |\Gamma^{\text{int}}(D_\pi)| + |\mathcal{F}(D_\pi)| + |\Gamma_{\text{ac}}^{\text{lf}}|$.

1. $\text{dcp}_1(\mathbb{C})$: the number $|V(H)| - |V_{\text{H}}|$ of non-hydrogen atoms in \mathbb{C} .
2. $\text{dcp}_2(\mathbb{C})$: the rank $r(\mathbb{C})$ of \mathbb{C} .
3. $\text{dcp}_3(\mathbb{C})$: the number $|V^{\text{int}}(\mathbb{C})|$ of interior-vertices in \mathbb{C} .
4. $\text{dcp}_4(\mathbb{C})$: the average $\overline{\text{ms}}(\mathbb{C})$ of mass^* over all atoms in \mathbb{C} ;
i.e., $\overline{\text{ms}}(\mathbb{C}) \triangleq \frac{1}{|V(H)|} \sum_{v \in V(H)} \text{mass}^*(\alpha(v))$.
5. $\text{dcp}_i(\mathbb{C})$, $i = 4 + d, d \in [1, 4]$: the number $\text{dg}_d^{\overline{\text{H}}}(\mathbb{C})$ of non-hydrogen vertices $v \in V(H) \setminus V_{\text{H}}$ of degree $\deg_{\langle \mathbb{C} \rangle}(v) = d$ in the hydrogen-suppressed chemical graph $\langle \mathbb{C} \rangle$.
6. $\text{dcp}_i(\mathbb{C})$, $i = 8 + d, d \in [1, 4]$: the number $\text{dg}_d^{\text{int}}(\mathbb{C})$ of interior-vertices of interior-degree $\deg_{\mathbb{C}^{\text{int}}}(v) = d$ in the interior $\mathbb{C}^{\text{int}} = (V^{\text{int}}(\mathbb{C}), E^{\text{int}}(\mathbb{C}))$ of \mathbb{C} .
7. $\text{dcp}_i(\mathbb{C})$, $i = 12 + m, m \in [2, 3]$: the number $\text{bd}_m^{\text{int}}(\mathbb{C})$ of interior-edges with bond multiplicity m in \mathbb{C} ; i.e., $\text{bd}_m^{\text{int}}(\mathbb{C}) \triangleq \{e \in E^{\text{int}}(\mathbb{C}) \mid \beta(e) = m\}$.
8. $\text{dcp}_i(\mathbb{C})$, $i = 14 + [\mathbf{a}]^{\text{int}}$, $\mathbf{a} \in \Lambda^{\text{int}}(D_\pi)$: the frequency $\text{na}_{\mathbf{a}}^{\text{int}}(\mathbb{C}) = |V_{\mathbf{a}}(\mathbb{C}) \cap V^{\text{int}}(\mathbb{C})|$ of chemical element \mathbf{a} in the set $V^{\text{int}}(\mathbb{C})$ of interior-vertices in \mathbb{C} .
9. $\text{dcp}_i(\mathbb{C})$, $i = 14 + |\Lambda^{\text{int}}(D_\pi)| + [\mathbf{a}]^{\text{ex}}$, $\mathbf{a} \in \Lambda^{\text{ex}}(D_\pi)$: the frequency $\text{na}_{\mathbf{a}}^{\text{ex}}(\mathbb{C}) = |V_{\mathbf{a}}(\mathbb{C}) \cap V^{\text{ex}}(\mathbb{C})|$ of chemical element \mathbf{a} in the set $V^{\text{ex}}(\mathbb{C})$ of exterior-vertices in \mathbb{C} .
10. $\text{dcp}_i(\mathbb{C})$, $i = 14 + |\Lambda^{\text{int}}(D_\pi)| + |\Lambda^{\text{ex}}(D_\pi)| + [\gamma]$, $\gamma \in \Gamma^{\text{int}}(D_\pi)$: the frequency $\text{ec}_\gamma(G)$ of edge-configuration γ in the set $E^{\text{int}}(\mathbb{C})$ of interior-edges in \mathbb{C} .
11. $\text{dcp}_i(\mathbb{C})$, $i = 14 + |\Lambda^{\text{int}}(D_\pi)| + |\Lambda^{\text{ex}}(D_\pi)| + |\Gamma^{\text{int}}(D_\pi)| + [\psi]$, $\psi \in \mathcal{F}(D_\pi)$: the frequency $\text{fc}_\psi(\mathbb{C})$ of fringe-configuration ψ in the set of ρ -fringe-trees in \mathbb{C} .
12. $\text{dcp}_i(\mathbb{C})$, $i = 14 + |\Lambda^{\text{int}}(D_\pi)| + |\Lambda^{\text{ex}}(D_\pi)| + |\Gamma^{\text{int}}(D_\pi)| + |\mathcal{F}(D_\pi)| + [\nu]$, $\nu \in \Gamma_{\text{ac}}^{\text{lf}}$: the frequency $\text{ac}_\nu^{\text{lf}}(\mathbb{C})$ of adjacency-configuration ν in the set of leaf-edges in $\langle \mathbb{C} \rangle$.

B Specifying Target Chemical Graphs

Given a prediction function η and a target value $y^* \in \mathbb{R}$, we call a chemical graph \mathbb{C}^* such that $\eta(x^*) = y^*$ for the feature vector $x^* = f(\mathbb{C}^*)$ a *target chemical graph*. This section presents a set of rules for specifying topological substructure of a target chemical graph in a flexible way in Stage 4.

We first describe how to reduce a chemical graph $\mathbb{C} = (H, \alpha, \beta)$ into an abstract form based on which our specification rules will be defined. To illustrate the reduction process, we use the chemical graph $\mathbb{C} = (H, \alpha, \beta)$ such that $\langle \mathbb{C} \rangle$ is given in Figure 3.

R1 Removal of all ρ -fringe-trees: The interior $H^{\text{int}} = (V^{\text{int}}(\mathbb{C}), E^{\text{int}}(\mathbb{C}))$ of \mathbb{C} is obtained by removing the non-root vertices of each ρ -fringe-trees $\mathbb{C}[u] \in \mathcal{T}(\mathbb{C}), u \in V^{\text{int}}(\mathbb{C})$. Figure 8 illustrates the interior H^{int} of chemical graph \mathbb{C} with $\rho = 2$ in Figure 3.

R2 Removal of some leaf paths: We call a u, v -path Q in H^{int} a *leaf path* if vertex v is a leaf-vertex of H^{int} and the degree of each internal vertex of Q in H^{int} is 2, where we regard that Q is rooted at vertex u . A connected subgraph S of the interior H^{int} of \mathbb{C} is called a *cyclical-base* if S is obtained from H by removing the vertices in $V(Q_u) \setminus \{u\}, u \in X$ for a subset X of interior-vertices and a set $\{Q_u \mid u \in X\}$ of leaf u, v -paths Q_u such that no two paths Q_u and $Q_{u'}$ share a vertex. Figure 9(a) illustrates a cyclical-base $S = H^{\text{int}} - \bigcup_{u \in X} (V(Q_u) \setminus \{u\})$ of the interior H^{int} for a set $\{Q_{u_5} = (u_5, u_{24}), Q_{u_{18}} = (u_{18}, u_{25}, u_{26}, u_{27}), Q_{u_{22}} = (u_{22}, u_{28})\}$ of leaf paths in Figure 8.

R3 Contraction of some pure paths: A path in S is called *pure* if each internal vertex of the path is of degree 2. Choose a set \mathcal{P} of several pure paths in S so that no two paths share vertices except for their end-vertices. A graph S' is called a *contraction* of a graph S (with respect to \mathcal{P}) if S' is obtained from S by replacing each pure u, v -path with a single edge $a = uv$, where S' may contain multiple edges between the same pair of adjacent vertices. Figure 9(b) illustrates a contraction S' obtained from the chemical graph S by contracting each uv -path $P_a \in \mathcal{P}$ into a new edge $a = uv$, where $a_1 = u_1u_2, a_2 = u_1u_3, a_3 = u_4u_7, a_4 = u_{10}u_{11}$ and $a_5 = u_{11}u_{12}$ and $\mathcal{P} = \{P_{a_1} = (u_1, u_{13}, u_2), P_{a_2} = (u_1, u_{14}, u_3), P_{a_3} = (u_4, u_{15}, u_{16}, u_7), P_{a_4} = (u_{10}, u_{17}, u_{18}, u_{19}, u_{11}), P_{a_5} = (u_{11}, u_{20}, u_{21}, u_{22}, u_{12})\}$ of pure paths in Figure 9(a).

We will define a set of rules so that a chemical graph can be obtained from a graph (called a seed graph in the next section) by applying processes R3 to R1 in a reverse way. We specify topological substructures of a target chemical graph with a tuple $(G_C, \sigma_{\text{int}}, \sigma_{\text{ce}})$ called a *target specification* defined under the set of the following rules.

Seed Graph

A *seed graph* $G_C = (V_C, E_C)$ is defined to be a graph (possibly with multiple edges) such that the edge set E_C consists of four sets $E_{(\geq 2)}, E_{(\geq 1)}, E_{(0/1)}$ and $E_{(=1)}$, where each of them can be empty. A seed graph plays a role of the most abstract form S' in R3. Figure 5(a) illustrates an example of a seed graph G_C with $r(G_C) = 5$, where $V_C = \{u_1, u_2, \dots, u_{12}, u_{23}\}$, $E_{(\geq 2)} = \{a_1, a_2, \dots, a_5\}$, $E_{(\geq 1)} = \{a_6\}$, $E_{(0/1)} = \{a_7\}$ and $E_{(=1)} = \{a_8, a_9, \dots, a_{16}\}$.

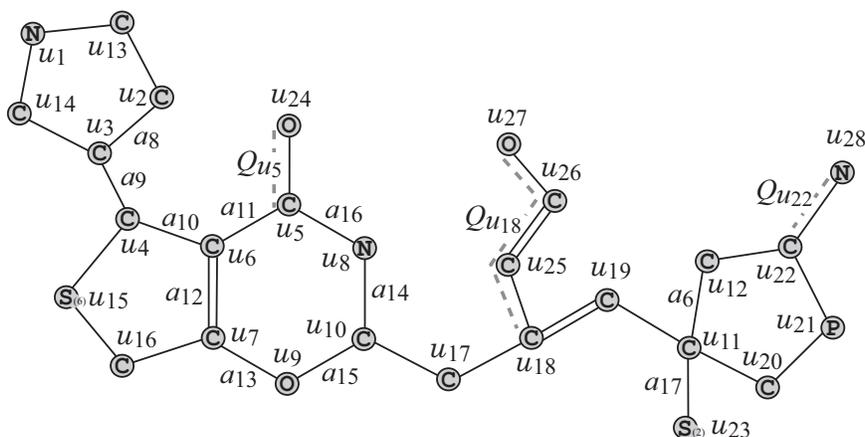
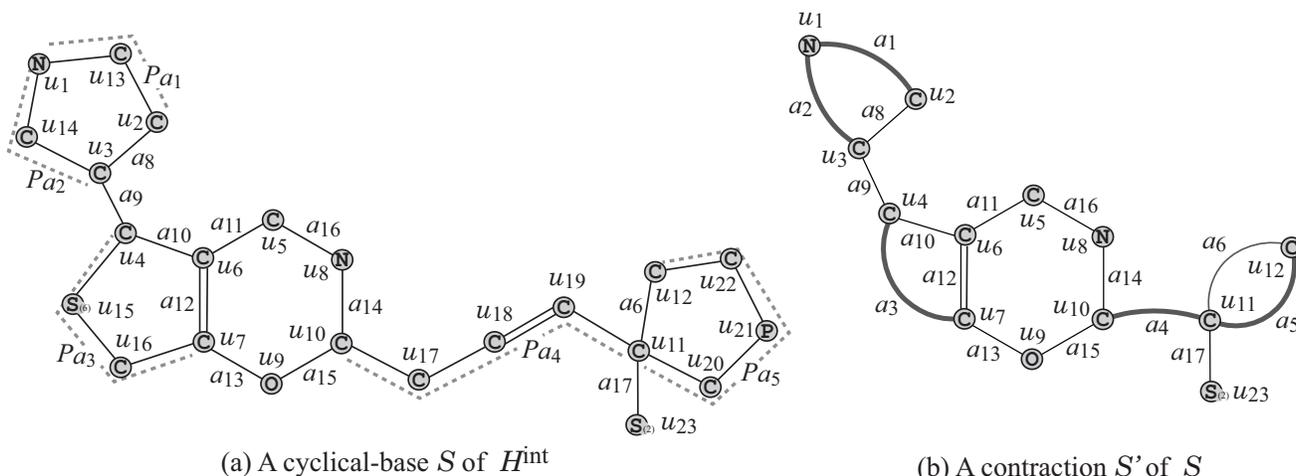


Figure 8: The interior H^{int} of chemical graph \mathbb{C} with $\langle \mathbb{C} \rangle$ in Figure 3 for $\rho = 2$.



(a) A cyclical-base S of H^{int}

(b) A contraction S' of S

Figure 9: (a) A cyclical-base $S = H^{\text{int}} - \bigcup_{u \in \{u_5, u_{18}, u_{22}\}} (V(Q_u) \setminus \{u\})$ of the interior H^{int} in Figure 8; (b) A contraction S' of S for a pure path set $\mathcal{P} = \{P_{a_1}, P_{a_2}, \dots, P_{a_5}\}$ in (a), where a new edge obtained by contracting a pure path is depicted with a thick line.

A *subdivision* S of $G_{\mathbb{C}}$ is a graph constructed from a seed graph $G_{\mathbb{C}}$ according to the following rules:

- Each edge $e = uv \in E_{(\geq 2)}$ is replaced with a u, v -path P_e of length at least 2;
- Each edge $e = uv \in E_{(\geq 1)}$ is replaced with a u, v -path P_e of length at least 1 (equivalently e is directly used or replaced with a u, v -path P_e of length at least 2);
- Each edge $e \in E_{(0/1)}$ is either used or discarded, where $E_{(0/1)}$ is required to be chosen as a non-separating edge subset of $E(G_{\mathbb{C}})$ since otherwise the connectivity of a final chemical graph \mathbb{C} is not guaranteed; $r(\mathbb{C}) = r(G_{\mathbb{C}}) - |E'|$ holds for a subset $E' \subseteq E_{(0/1)}$ of edges discarded in a final chemical graph \mathbb{C} ; and

- Each edge $e \in E_{(=1)}$ is always used directly.

We allow a possible elimination of edges in $E_{(0/1)}$ as an optional rule in constructing a target chemical graph from a seed graph, even though such an operation has not been included in the process R3. A subdivision S plays a role of a cyclical-base in R2. A target chemical graph $\mathbb{C} = (H, \alpha, \beta)$ will contain S as a subgraph of the interior H^{int} of \mathbb{C} .

Interior-specification

A graph H^* that serves as the interior H^{int} of a target chemical graph \mathbb{C} will be constructed as follows. First construct a subdivision S of a seed graph G_C by replacing each edge $e = uu' \in E_{(\geq 2)} \cup E_{(\geq 1)}$ with a pure u, u' -path P_e . Next construct a supergraph H^* of S by attaching a leaf path Q_v at each vertex $v \in V_C$ or at an internal vertex $v \in V(P_e) \setminus \{u, u'\}$ of each pure u, u' -path P_e for some edge $e = uu' \in E_{(\geq 2)} \cup E_{(\geq 1)}$, where possibly $Q_v = (v), E(Q_v) = \emptyset$ (i.e., we do not attach any new edges to v). We introduce the following rules for specifying the size of H^* , the length $|E(P_e)|$ of a pure path P_e , the length $|E(Q_v)|$ of a leaf path Q_v , the number of leaf paths Q_v and a bond-multiplicity of each interior-edge, where we call the set of prescribed constants an *interior-specification* σ_{int} :

- Lower and upper bounds $n_{\text{LB}}^{\text{int}}, n_{\text{UB}}^{\text{int}} \in \mathbb{Z}_+$ on the number of interior-vertices of a target chemical graph \mathbb{C} .
- For each edge $e = uu' \in E_{(\geq 2)} \cup E_{(\geq 1)}$,
 - a lower bound $\ell_{\text{LB}}(e)$ and an upper bound $\ell_{\text{UB}}(e)$ on the length $|E(P_e)|$ of a pure u, u' -path P_e . (For a notational convenience, set $\ell_{\text{LB}}(e) := 0, \ell_{\text{UB}}(e) := 1, e \in E_{(0/1)}$ and $\ell_{\text{LB}}(e) := 1, \ell_{\text{UB}}(e) := 1, e \in E_{(=1)}$.)
 - a lower bound $\text{bl}_{\text{LB}}(e)$ and an upper bound $\text{bl}_{\text{UB}}(e)$ on the number of leaf paths Q_v attached at internal vertices v of a pure u, u' -path P_e .
 - a lower bound $\text{ch}_{\text{LB}}(e)$ and an upper bound $\text{ch}_{\text{UB}}(e)$ on the maximum length $|E(Q_v)|$ of a leaf path Q_v attached at an internal vertex $v \in V(P_e) \setminus \{u, u'\}$ of a pure u, u' -path P_e .
- For each vertex $v \in V_C$,
 - a lower bound $\text{ch}_{\text{LB}}(v)$ and an upper bound $\text{ch}_{\text{UB}}(v)$ on the number of leaf paths Q_v attached to v , where $0 \leq \text{ch}_{\text{LB}}(v) \leq \text{ch}_{\text{UB}}(v) \leq 1$.
 - a lower bound $\text{ch}_{\text{LB}}(v)$ and an upper bound $\text{ch}_{\text{UB}}(v)$ on the length $|E(Q_v)|$ of a leaf path Q_v attached to v .
- For each edge $e = uu' \in E_C$, a lower bound $\text{bd}_{m,\text{LB}}(e)$ and an upper bound $\text{bd}_{m,\text{UB}}(e)$ on the number of edges with bond-multiplicity $m \in [2, 3]$ in u, u' -path P_e , where we regard $P_e, e \in E_{(0/1)} \cup E_{(=1)}$ as single edge e .

We call a graph H^* that satisfies an interior-specification σ_{int} a σ_{int} -*extension* of G_C , where the bond-multiplicity of each edge has been determined.

Table 5: Example 1 of an interior-specification σ_{int} .

$n_{\text{LB}}^{\text{int}} = 20$	$n_{\text{UB}}^{\text{int}} = 28$					
	a_1	a_2	a_3	a_4	a_5	a_6
$\ell_{\text{LB}}(a_i)$	2	2	2	3	2	1
$\ell_{\text{UB}}(a_i)$	3	4	3	5	4	4
$\text{bl}_{\text{LB}}(a_i)$	0	0	0	1	1	0
$\text{bl}_{\text{UB}}(a_i)$	1	1	0	2	1	0
$\text{ch}_{\text{LB}}(a_i)$	0	1	0	4	3	0
$\text{ch}_{\text{UB}}(a_i)$	3	3	1	6	5	2

	u_1	u_2	u_3	u_4	u_5	u_6	u_7	u_8	u_9	u_{10}	u_{11}	u_{12}	u_{23}
$\text{bl}_{\text{LB}}(u_i)$	0	0	0	0	0	0	0	0	0	0	0	0	0
$\text{bl}_{\text{UB}}(u_i)$	1	1	1	1	1	0	0	0	0	0	0	0	0
$\text{ch}_{\text{LB}}(u_i)$	0	0	0	0	1	0	0	0	0	0	0	0	0
$\text{ch}_{\text{UB}}(u_i)$	1	0	0	0	3	0	1	1	0	1	2	4	1

	a_1	a_2	a_3	a_4	a_5	a_6	a_7	a_8	a_9	a_{10}	a_{11}	a_{12}	a_{13}	a_{14}	a_{15}	a_{16}	a_{17}
$\text{bd}_{2,\text{LB}}(a_i)$	0	0	0	1	0	0	0	0	0	0	0	1	0	0	0	0	0
$\text{bd}_{2,\text{UB}}(a_i)$	1	1	0	2	2	0	0	0	0	0	0	1	0	0	0	0	0
$\text{bd}_{3,\text{LB}}(a_i)$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
$\text{bd}_{3,\text{UB}}(a_i)$	0	0	0	0	1	0	0	0	0	0	0	0	0	0	0	0	0

Table 5 shows an example of an interior-specification σ_{int} to the seed graph G_C in Figure 5.

Figure 10 illustrates an example of an σ_{int} -extension H^* of seed graph G_C in Figure 5 under the interior-specification σ_{int} in Table 5.

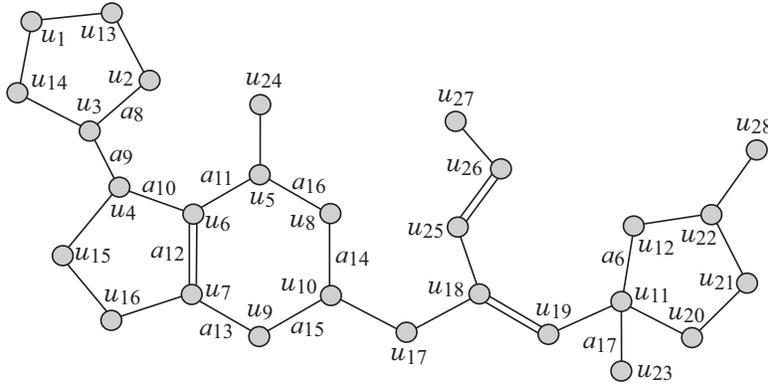


Figure 10: An illustration of a graph H^* that is obtained from the seed graph G_C in Figure 5 under the interior-specification σ_{int} in Table 5, where the vertices newly introduced by pure paths P_{a_i} and leaf paths Q_{v_i} are depicted with white squares and circles, respectively.

Chemical-specification

Let H^* be a graph that serves as the interior H^{int} of a target chemical graph \mathbb{C} , where the bond-multiplicity of each edge in H^* has been determined. Finally we introduce a set of rules for constructing a target chemical graph \mathbb{C} from H^* by choosing a chemical element $\mathbf{a} \in \Lambda$ and assigning a ρ -fringe-tree ψ to each interior-vertex $v \in V^{\text{int}}$. We introduce the following rules for specifying the size of \mathbb{C} , a set of chemical rooted trees that are allowed to use as ρ -fringe-trees and lower and upper bounds on the frequency of a chemical element, a chemical symbol, and an edge-configuration, where we call the set of prescribed constants a *chemical specification* σ_{ce} :

- Lower and upper bounds $n_{\text{LB}}, n^* \in \mathbb{Z}_+$ on the number of vertices, where $n_{\text{LB}}^{\text{int}} \leq n_{\text{LB}} \leq n^*$.
- Subsets $\mathcal{F}(v) \subseteq \mathcal{F}(D_\pi), v \in V_{\mathbb{C}}$ and $\mathcal{F}_E \subseteq \mathcal{F}(D_\pi)$ of chemical rooted trees ψ with $\text{ht}(\langle \psi \rangle) \leq \rho$, where we require that every ρ -fringe-tree $\mathbb{C}[v]$ rooted at a vertex $v \in V_{\mathbb{C}}$ (resp., at an internal vertex v not in $V_{\mathbb{C}}$) in \mathbb{C} belongs to $\mathcal{F}(v)$ (resp., \mathcal{F}_E). Let $\mathcal{F}^* := \mathcal{F}_E \cup \bigcup_{v \in V_{\mathbb{C}}} \mathcal{F}(v)$ and Λ^{ex} denote the set of chemical elements assigned to non-root vertices over all chemical rooted trees in \mathcal{F}^* .
- A subset $\Lambda^{\text{int}} \subseteq \Lambda^{\text{int}}(D_\pi)$, where we require that every chemical element $\alpha(v)$ assigned to an interior-vertex v in \mathbb{C} belongs to Λ^{int} . Let $\Lambda := \Lambda^{\text{int}} \cup \Lambda^{\text{ex}}$ and $\text{na}_{\mathbf{a}}(\mathbb{C})$ (resp., $\text{na}_{\mathbf{a}}^{\text{int}}(\mathbb{C})$ and $\text{na}_{\mathbf{a}}^{\text{ex}}(\mathbb{C})$) denote the number of vertices (resp., interior-vertices and exterior-vertices) v such that $\alpha(v) = \mathbf{a}$ in \mathbb{C} .
- A set $\Lambda_{\text{dg}}^{\text{int}} \subseteq \Lambda \times [1, 4]$ of chemical symbols and a set $\Gamma^{\text{int}} \subseteq \Gamma^{\text{int}}(D_\pi)$ of edge-configurations (μ, μ', m) with $\mu \leq \mu'$, where we require that the edge-configuration $\text{ec}(e)$ of an interior-edge e in \mathbb{C} belongs to Γ^{int} . We do not distinguish (μ, μ', m) and (μ', μ, m) .
- Define $\Gamma_{\text{ac}}^{\text{int}}$ to be the set of adjacency-configurations such that $\Gamma_{\text{ac}}^{\text{int}} := \{(\mathbf{a}, \mathbf{b}, m) \mid (\mathbf{a}\mathbf{d}, \mathbf{b}\mathbf{d}', m) \in \Gamma^{\text{int}}\}$. Let $\text{ac}_\nu^{\text{int}}(\mathbb{C}), \nu \in \Gamma_{\text{ac}}^{\text{int}}$ denote the number of interior-edges e such that $\text{ac}(e) = \nu$ in \mathbb{C} .
- Subsets $\Lambda^*(v) \subseteq \{\mathbf{a} \in \Lambda^{\text{int}} \mid \text{val}(\mathbf{a}) \geq 2\}, v \in V_{\mathbb{C}}$, we require that every chemical element $\alpha(v)$ assigned to a vertex $v \in V_{\mathbb{C}}$ in the seed graph belongs to $\Lambda^*(v)$.
- Lower and upper bound functions $\text{na}_{\text{LB}}, \text{na}_{\text{UB}} : \Lambda \rightarrow [1, n^*]$ and $\text{na}_{\text{LB}}^{\text{int}}, \text{na}_{\text{UB}}^{\text{int}} : \Lambda^{\text{int}} \rightarrow [1, n^*]$ on the number of interior-vertices v such that $\alpha(v) = \mathbf{a}$ in \mathbb{C} .
- Lower and upper bound functions $\text{ns}_{\text{LB}}^{\text{int}}, \text{ns}_{\text{UB}}^{\text{int}} : \Lambda_{\text{dg}}^{\text{int}} \rightarrow [1, n^*]$ on the number of interior-vertices v such that $\text{cs}(v) = \mu$ in \mathbb{C} .
- Lower and upper bound functions $\text{ac}_{\text{LB}}^{\text{int}}, \text{ac}_{\text{UB}}^{\text{int}} : \Gamma_{\text{ac}}^{\text{int}} \rightarrow \mathbb{Z}_+$ on the number of interior-edges e such that $\text{ac}(e) = \nu$ in \mathbb{C} .
- Lower and upper bound functions $\text{ec}_{\text{LB}}^{\text{int}}, \text{ec}_{\text{UB}}^{\text{int}} : \Gamma^{\text{int}} \rightarrow \mathbb{Z}_+$ on the number of interior-edges e such that $\text{ec}(e) = \gamma$ in \mathbb{C} .
- Lower and upper bound functions $\text{fc}_{\text{LB}}, \text{fc}_{\text{UB}} : \mathcal{F}^* \rightarrow [0, n^*]$ on the number of interior-vertices v such that $\mathbb{C}[v]$ is r-isomorphic to $\psi \in \mathcal{F}^*$ in \mathbb{C} .

- Lower and upper bound functions $\text{ac}_{\text{LB}}^{\text{lf}}, \text{ac}_{\text{UB}}^{\text{lf}} : \Gamma_{\text{ac}}^{\text{lf}} \rightarrow [0, n^*]$ on the number of leaf-edges uv in ac_{C} with adjacency-configuration ν .

We call a chemical graph \mathbb{C} that satisfies a chemical specification σ_{ce} a $(\sigma_{\text{int}}, \sigma_{\text{ce}})$ -extension of G_{C} , and denote by $\mathcal{G}(G_{\text{C}}, \sigma_{\text{int}}, \sigma_{\text{ce}})$ the set of all $(\sigma_{\text{int}}, \sigma_{\text{ce}})$ -extensions of G_{C} .

Table 6 shows an example of a chemical-specification σ_{ce} to the seed graph G_{C} in Figure 5.

Table 6: Example 2 of a chemical-specification σ_{ce} .

$n_{\text{LB}} = 30, n^* = 50.$																		
branch-parameter: $\rho = 2$																		
Each of sets $\mathcal{F}(v), v \in V_{\text{C}}$ and \mathcal{F}_E is set to be the set \mathcal{F} of chemical rooted trees ψ with $\text{ht}(\langle \psi \rangle) \leq \rho = 2$ in Figure 5(b).																		
$\Lambda = \{\text{H}, \text{C}, \text{N}, \text{O}, \text{S}_{(2)}, \text{S}_{(6)}, \text{P} = \text{P}_{(5)}\}$									$\Lambda_{\text{dg}} = \{\text{C2}, \text{C3}, \text{C4}, \text{N2}, \text{N3}, \text{O2}, \text{S}_{(2)}2, \text{S}_{(6)}3, \text{P4}\}$									
$\Gamma_{\text{ac}}^{\text{int}}$	$\nu_1 = (\text{C}, \text{C}, 1), \nu_2 = (\text{C}, \text{C}, 2), \nu_3 = (\text{C}, \text{N}, 1), \nu_4 = (\text{C}, \text{O}, 1), \nu_5 = (\text{C}, \text{S}_{(2)}, 1), \nu_6 = (\text{C}, \text{S}_{(6)}, 1), \nu_7 = (\text{C}, \text{P}, 1)$																	
Γ^{int}	$\gamma_1 = (\text{C2}, \text{C2}, 1), \gamma_2 = (\text{C2}, \text{C3}, 1), \gamma_3 = (\text{C2}, \text{C3}, 2), \gamma_4 = (\text{C2}, \text{C4}, 1), \gamma_5 = (\text{C3}, \text{C3}, 1), \gamma_6 = (\text{C3}, \text{C3}, 2), \gamma_7 = (\text{C3}, \text{C4}, 1), \gamma_8 = (\text{C2}, \text{N2}, 1), \gamma_9 = (\text{C3}, \text{N2}, 1), \gamma_{10} = (\text{C3}, \text{O2}, 1), \gamma_{11} = (\text{C2}, \text{C2}, 2), \gamma_{12} = (\text{C2}, \text{O2}, 1), \gamma_{13} = (\text{C3}, \text{N3}, 1), \gamma_{14} = (\text{C4}, \text{S}_{(2)}2, 2), \gamma_{15} = (\text{C2}, \text{S}_{(6)}3, 1), \gamma_{16} = (\text{C3}, \text{S}_{(6)}3, 1), \gamma_{17} = (\text{C2}, \text{P4}, 2), \gamma_{18} = (\text{C3}, \text{P4}, 1)$																	
$\Lambda^*(u_1) = \Lambda^*(u_8) = \{\text{C}, \text{N}\}, \Lambda^*(u_9) = \{\text{C}, \text{O}\}, \Lambda^*(u) = \{\text{C}\}, u \in V_{\text{C}} \setminus \{u_1, u_8, u_9\}$																		
	H	C	N	O	S ₍₂₎	S ₍₆₎	P		C	N	O	S ₍₂₎	S ₍₆₎	P				
$\text{na}_{\text{LB}}(\mathbf{a})$	40	27	1	1	0	0	0		$\text{na}_{\text{LB}}^{\text{int}}(\mathbf{a})$	9	1	0	0	0	0			
$\text{na}_{\text{UB}}(\mathbf{a})$	65	37	4	8	1	1	1		$\text{na}_{\text{UB}}^{\text{int}}(\mathbf{a})$	23	4	5	1	1	1			
	C2	C3	C4	N2	N3	O2	S ₍₂₎ 2	S ₍₆₎ 3	P4									
$\text{ns}_{\text{LB}}^{\text{int}}(\mu)$	3	5	0	0	0	0	0	0	0									
$\text{ns}_{\text{UB}}^{\text{int}}(\mu)$	8	15	2	2	3	5	1	1	1									
	ν_1	ν_2	ν_3	ν_4	ν_5	ν_6	ν_7											
$\text{ac}_{\text{LB}}^{\text{int}}(\nu)$	0	0	0	0	0	0	0											
$\text{ac}_{\text{UB}}^{\text{int}}(\nu)$	30	10	10	10	1	1	1											
	γ_1	γ_2	γ_3	γ_4	γ_5	γ_6	γ_7	γ_8	γ_9	γ_{10}	γ_{11}	γ_{12}	γ_{13}	γ_{14}	γ_{15}	γ_{16}	γ_{17}	γ_{18}
$\text{ec}_{\text{LB}}^{\text{int}}(\gamma)$	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0	0
$\text{ec}_{\text{UB}}^{\text{int}}(\gamma)$	4	15	4	4	10	5	4	4	6	4	4	4	2	2	2	2	2	2
	$\psi \in \{\psi_i \mid i = 1, 6, 11\} \quad \psi \in \mathcal{F}^* \setminus \{\psi_i \mid i = 1, 6, 11\}$																	
$\text{fc}_{\text{LB}}(\psi)$	1						0											
$\text{fc}_{\text{UB}}(\psi)$	10						3											
	$\nu \in \{(\text{C}, \text{C}, 1), (\text{C}, \text{C}, 2)\} \quad \nu \in \Gamma_{\text{ac}}^{\text{lf}} \setminus \{(\text{C}, \text{C}, 1), (\text{C}, \text{C}, 2)\}$																	
$\text{ac}_{\text{LB}}^{\text{lf}}(\nu)$	0						0											
$\text{ac}_{\text{UB}}^{\text{lf}}(\nu)$	10						8											

Figure 3 illustrates an example \mathbb{C} of a $(\sigma_{\text{int}}, \sigma_{\text{ce}})$ -extension of G_{C} obtained from the σ_{int} -extension H^* in Figure 10 under the chemical-specification σ_{ce} in Table 6. Note that $r(\mathbb{C}) = r(H^*) = r(G_{\text{C}}) - 1 = 4$ holds since the edge in $E_{(0/1)}$ is discarded in H^* .

C Test Instances for Stages 4 and 5

We prepared the following instances (a)-(d) for conducting experiments of Stages 4 and 5 in Phase 2.

In Stages 4 and 5, we use five properties $\pi \in \{\text{HC}, \text{VD}, \text{OPTR}, \text{IHCLIQ}, \text{VIS}\}$ and define a set $\Lambda(\pi)$ of chemical elements as follows:

$$\begin{aligned} \Lambda(\text{HC}) &= \{\text{H}, \text{C}, \text{N}, \text{O}, \text{S}_{(2)}, \text{S}_{(6)}, \text{Cl}\}, & \Lambda(\text{VD}) &= \{\text{H}, \text{C}, \text{N}, \text{O}, \text{N}, \text{Cl}, \text{P}_{(3)}, \text{P}_{(5)}\}, \\ \Lambda(\text{OPTR}) &= \{\text{H}, \text{C}, \text{N}, \text{O}, \text{S}_{(2)}, \text{F}\}, & \Lambda(\text{IHCLIQ}) &= \{\text{H}, \text{C}, \text{N}, \text{O}, \text{S}_{(2)}, \text{S}_{(6)}, \text{Cl}\} \text{ and} \\ \Lambda(\text{VIS}) &= \{\text{H}, \text{C}, \text{O}, \text{Si}\}. \end{aligned}$$

- (a) $I_{\mathbf{a}} = (G_{\mathbf{C}}, \sigma_{\text{int}}, \sigma_{\text{ce}})$: The instance introduced in Section B to explain the target specification. For each property π , we replace $\Lambda = \{\text{H}, \text{C}, \text{N}, \text{O}, \text{S}_{(2)}, \text{S}_{(6)}, \text{P}_{(5)}\}$ in Table 6 with $\Lambda(\pi) \cap \{\text{S}_{(2)}, \text{S}_{(6)}, \text{P}_{(5)}\}$ and remove from the σ_{ce} all chemical symbols, edge-configurations and fringe-configurations that cannot be constructed from the replaced element set (i.e., those containing a chemical element in $\{\text{S}_{(2)}, \text{S}_{(6)}, \text{P}_{(5)}\} \setminus \Lambda(\pi)$).
- (b) $I_{\mathbf{b}}^i = (G_{\mathbf{C}}^i, \sigma_{\text{int}}^i, \sigma_{\text{ce}}^i)$, $i = 1, 2, 3, 4$: An instance for inferring chemical graphs with rank at most 2. In the four instances $I_{\mathbf{b}}^i$, $i = 1, 2, 3, 4$, the following specifications in $(\sigma_{\text{int}}, \sigma_{\text{ce}})$ are common.

Set $\Lambda := \Lambda(\pi)$ for a given property $\pi \in \{\text{HC}, \text{VD}, \text{OPTR}, \text{IHCLIQ}, \text{VIS}\}$, set $\Lambda_{\text{dg}}^{\text{int}}$ to be the set of all possible symbols in $\Lambda \times [1, 4]$ that appear in the data set D_{π} and set Γ^{int} to be the set of all edge-configurations that appear in the data set D_{π} . Set $\Lambda^*(v) := \Lambda$, $v \in V_{\mathbf{C}}$.

The lower bounds $\ell_{\text{LB}}, \text{bl}_{\text{LB}}, \text{ch}_{\text{LB}}, \text{bd}_{2,\text{LB}}, \text{bd}_{3,\text{LB}}, \text{na}_{\text{LB}}, \text{na}_{\text{LB}}^{\text{int}}, \text{ns}_{\text{LB}}^{\text{int}}, \text{ac}_{\text{LB}}^{\text{int}}, \text{ec}_{\text{LB}}^{\text{int}}$ and $\text{ac}_{\text{LB}}^{\text{lf}}$ are all set to be 0.

Set upper bounds $\text{na}_{\text{UB}}(\mathbf{a}) := n^*$, $\text{na} \in \{\text{H}, \text{C}\}$, $\text{na}_{\text{UB}}(\mathbf{a}) := 5$, $\text{na} \in \{\text{O}, \text{N}\}$, $\text{na}_{\text{UB}}(\mathbf{a}) := 2$, $\text{na} \in \Lambda \setminus \{\text{H}, \text{C}, \text{O}, \text{N}\}$. The other upper bounds $\ell_{\text{UB}}, \text{bl}_{\text{UB}}, \text{ch}_{\text{UB}}, \text{bd}_{2,\text{UB}}, \text{bd}_{3,\text{UB}}, \text{na}_{\text{UB}}^{\text{int}}, \text{ns}_{\text{UB}}^{\text{int}}, \text{ac}_{\text{UB}}^{\text{int}}, \text{ec}_{\text{UB}}^{\text{int}}$ and $\text{ac}_{\text{UB}}^{\text{lf}}$ are all set to be an upper bound n^* on $n(G^*)$.

For each property π , let $\mathcal{F}(D_{\pi})$ denote the set of 2-fringe-trees in the compounds in D_{π} , and select a subset $\mathcal{F}_{\pi}^i \subseteq \mathcal{F}(D_{\pi})$ with $|\mathcal{F}_{\pi}^i| = 45 - 5i$, $i \in [1, 5]$. For each instance $I_{\mathbf{b}}^i$, set $\mathcal{F}_E := \mathcal{F}(v) := \mathcal{F}_{\pi}^i$, $v \in V_{\mathbf{C}}$ and $\text{fc}_{\text{LB}}(\psi) := 0$, $\text{fc}_{\text{UB}}(\psi) := 10$, $\psi \in \mathcal{F}_{\pi}^i$.

Instance $I_{\mathbf{b}}^1$ is given by the rank-1 seed graph $G_{\mathbf{C}}^1$ in Figure 6(i) and Instances $I_{\mathbf{b}}^i$, $i = 2, 3, 4$ are given by the rank-2 seed graph $G_{\mathbf{C}}^i$, $i = 2, 3, 4$ in Figure 6(ii)-(iv).

- (i) For instance $I_{\mathbf{b}}^1$, select as a seed graph the monocyclic graph $G_{\mathbf{C}}^1 = (V_{\mathbf{C}}, E_{\mathbf{C}} = E_{(\geq 2)} \cup E_{(\geq 1)})$ in Figure 6(i), where $V_{\mathbf{C}} = \{u_1, u_2\}$, $E_{(\geq 2)} = \{a_1\}$ and $E_{(\geq 1)} = \{a_2\}$. Set $\text{n}_{\text{LB}}^{\text{int}} := 5$, $\text{n}_{\text{UB}}^{\text{int}} := 15$, $n_{\text{LB}} := 35$ and $n^* := 38$. We include a linear constraint $\ell(a_1) \leq \ell(a_2)$ and $5 \leq \ell(a_1) + \ell(a_2) \leq 15$ as part of the side constraint.
- (ii) For instance $I_{\mathbf{b}}^2$, select as a seed graph the graph $G_{\mathbf{C}}^2 = (V_{\mathbf{C}}, E_{\mathbf{C}} = E_{(\geq 2)} \cup E_{(\geq 1)} \cup E_{(=1)})$ in Figure 6(ii), where $V_{\mathbf{C}} = \{u_1, u_2, u_3, u_4\}$, $E_{(\geq 2)} = \{a_1, a_2\}$, $E_{(\geq 1)} = \{a_3\}$ and $E_{(=1)} = \{a_4, a_5\}$. Set $\text{n}_{\text{LB}}^{\text{int}} := 25$, $\text{n}_{\text{UB}}^{\text{int}} := 30$, $n_{\text{LB}} := 45$ and $n^* := 50$. We include a linear constraint $\ell(a_1) \leq \ell(a_2)$ and $\ell(a_1) + \ell(a_2) + \ell(a_3) \leq 15$.

- (iii) For instance I_b^3 , select as a seed graph the graph $G_C^3 = (V_C, E_C = E_{(\geq 2)} \cup E_{(\geq 1)} \cup E_{(=1)})$ in Figure 6(iii), where $V_C = \{u_1, u_2, u_3, u_4\}$, $E_{(\geq 2)} = \{a_1\}$, $E_{(\geq 1)} = \{a_2, a_3\}$ and $E_{(=1)} = \{a_4, a_5\}$. Set $n_{LB}^{int} := 25$, $n_{UB}^{int} := 30$, $n_{LB} := 45$ and $n^* := 50$. We include linear constraints $\ell(a_1) \leq \ell(a_2) + \ell(a_3)$, $\ell(a_2) \leq \ell(a_3)$ and $\ell(a_1) + \ell(a_2) + \ell(a_3) \leq 15$.
- (iv) For instance I_b^4 , select as a seed graph the graph $G_C^4 = (V_C, E_C = E_{(\geq 2)} \cup E_{(\geq 1)} \cup E_{(=1)})$ in Figure 6(iv), where $V_C = \{u_1, u_2, u_3, u_4\}$, $E_{(\geq 1)} = \{a_1, a_2, a_3\}$ and $E_{(=1)} = \{a_4, a_5\}$. Set $n_{LB}^{int} := 25$, $n_{UB}^{int} := 30$, $n_{LB} := 45$ and $n^* := 50$. We include linear constraints $\ell(a_2) \leq \ell(a_1) + 1$, $\ell(a_2) \leq \ell(a_3) + 1$, $\ell(a_1) \leq \ell(a_3)$ and $\ell(a_1) + \ell(a_2) + \ell(a_3) \leq 15$.

We define instances in (c) and (d) in order to find chemical graphs that have an intermediate structure of given two chemical cyclic graphs $G_A = (H_A = (V_A, E_A), \alpha_A, \beta_A)$ and $G_B = (H_B = (V_B, E_B), \alpha_B, \beta_B)$. Let Λ_A^{int} and $\Lambda_{dg,A}^{int}$ denote the sets of chemical elements and chemical symbols of the interior-vertices in G_A , Γ_A^{int} denote the sets of edge-configurations of the interior-edges in G_A , and \mathcal{F}_A denote the set of 2-fringe-trees in G_A . Analogously define sets Λ_B^{int} , $\Lambda_{dg,B}^{int}$, Γ_B^{int} and \mathcal{F}_B in G_B .

- (c) $I_c = (G_C, \sigma_{int}, \sigma_{ce})$: An instance aimed to infer a chemical graph G^\dagger such that the core of G^\dagger is equal to the core of G_A and the frequency of each edge-configuration in the non-core of G^\dagger is equal to that of G_B . We use chemical compounds CID 24822711 and CID 59170444 in Figure 7(a) and (b) for G_A and G_B , respectively.

Set a seed graph $G_C = (V_C, E_C = E_{(=1)})$ to be the core of G_A .

Set $\Lambda := \{H, C, N, O\}$, and set Λ_{dg}^{int} to be the set of all possible chemical symbols in $\Lambda \times [1, 4]$.

Set $\Gamma^{int} := \Gamma_A^{int} \cup \Gamma_B^{int}$ and $\Lambda^*(v) := \{\alpha_A(v)\}$, $v \in V_C$.

Set $n_{LB}^{int} := \min\{n^{int}(G_A), n^{int}(G_B)\}$, $n_{UB}^{int} := \max\{n^{int}(G_A), n^{int}(G_B)\}$,

$n_{LB} := \min\{n(G_A), n(G_B)\} - 10$ and $n^* := \max\{n(G_A), n(G_B)\} + 5$.

Set lower bounds ℓ_{LB} , bl_{LB} , ch_{LB} , $bd_{2,LB}$, $bd_{3,LB}$, na_{LB} , na_{LB}^{int} , ns_{LB}^{int} , ac_{LB}^{int} and ac_{LB}^{ff} to be 0.

Set upper bounds $na_{UB}(a) := n^*$, $na \in \{H, C\}$, $na_{UB}(a) := 5$, $na \in \{O, N\}$, $na_{UB}(a) := 2$, $na \in \Lambda \setminus \{H, C, O, N\}$ and set the other upper bounds ℓ_{UB} , bl_{UB} , ch_{UB} , $bd_{2,UB}$, $bd_{3,UB}$, na_{UB}^{int} , ns_{UB}^{int} , ac_{UB}^{int} and ac_{UB}^{ff} to be n^* .

Set $ec_{LB}^{int}(\gamma)$ to be the number of core-edges in G_A with $\gamma \in \Gamma^{int}$ and $ec_{UB}^{int}(\gamma)$ to be the number interior-edges in G_A and G_B with edge-configuration γ .

Let $\mathcal{F}_B^{(p)}$, $p \in [1, 2]$ denote the set of chemical rooted trees r-isomorphic p -fringe-trees in G_B ;

Set $\mathcal{F}_E := \mathcal{F}(v) := \mathcal{F}_B^{(1)} \cup \mathcal{F}_B^{(2)}$, $v \in V_C$ and $fc_{LB}(\psi) := 0$, $fc_{UB}(\psi) := 10$, $\psi \in \mathcal{F}_B^{(1)} \cup \mathcal{F}_B^{(2)}$.

- (d) $I_d = (G_C^1, \sigma_{int}, \sigma_{ce})$: An instance aimed to infer a chemical monocyclic graph G^\dagger such that the frequency vector of edge-configurations in G^\dagger is a vector obtained by merging those of G_A and G_B . We use chemical monocyclic compounds CID 10076784 and CID 44340250 in Figure 7(c) and (d) for G_A and G_B , respectively. Set a seed graph to be the monocyclic seed graph $G_C^1 = (V_C, E_C = E_{(\geq 2)} \cup E_{(\geq 1)})$ with $V_C = \{u_1, u_2\}$, $E_{(\geq 2)} = \{a_1\}$ and $E_{(\geq 1)} = \{a_2\}$ in Figure 6(i).

Set $\Lambda := \{H, C, N, O\}$, $\Lambda_{dg}^{int} := \Lambda_{dg,A}^{int} \cup \Lambda_{dg,B}^{int}$ and $\Gamma^{int} := \Gamma_A^{int} \cup \Gamma_B^{int}$.

Set $n_{LB}^{int} := \min\{n^{int}(G_A), n^{int}(G_B)\}$, $n_{UB}^{int} := \max\{n^{int}(G_A), n^{int}(G_B)\}$,

$n_{LB} := \min\{n(G_A), n(G_B)\}$ and $n^* := \max\{n(G_A), n(G_B)\}$.

Set lower bounds $\ell_{\text{LB}}, \text{bl}_{\text{LB}}, \text{ch}_{\text{LB}}, \text{bd}_{2,\text{LB}}, \text{bd}_{3,\text{LB}}, \text{na}_{\text{LB}}, \text{na}_{\text{LB}}^{\text{int}}, \text{ns}_{\text{LB}}^{\text{int}}, \text{ac}_{\text{LB}}^{\text{int}}$ and $\text{ac}_{\text{LB}}^{\text{f}}$ to be 0. Set upper bounds $\text{na}_{\text{UB}}(\mathbf{a}) := n^*, \text{na} \in \{\mathbf{H}, \mathbf{C}\}, \text{na}_{\text{UB}}(\mathbf{a}) := 5, \text{na} \in \{\mathbf{0}, \mathbf{N}\}, \text{na}_{\text{UB}}(\mathbf{a}) := 2, \text{na} \in \Lambda \setminus \{\mathbf{H}, \mathbf{C}, \mathbf{0}, \mathbf{N}\}$ and set the other upper bounds $\ell_{\text{UB}}, \text{bl}_{\text{UB}}, \text{ch}_{\text{UB}}, \text{bd}_{2,\text{UB}}, \text{bd}_{3,\text{UB}}, \text{na}_{\text{UB}}^{\text{int}}, \text{ns}_{\text{UB}}^{\text{int}}, \text{ac}_{\text{UB}}^{\text{int}}$ and $\text{ac}_{\text{UB}}^{\text{f}}$ to be n^* .

For each edge-configuration $\gamma \in \Gamma^{\text{int}}$, let $x_A^*(\gamma^{\text{int}})$ (resp., $x_B^*(\gamma^{\text{int}})$) denote the number of interior-edges with γ in G_A (resp., G_B), $\gamma \in \Gamma^{\text{int}}$ and set

$$x_{\min}^*(\gamma) := \min\{x_A^*(\gamma), x_B^*(\gamma)\}, x_{\max}^*(\gamma) := \max\{x_A^*(\gamma), x_B^*(\gamma)\},$$

$$\text{ec}_{\text{LB}}^{\text{int}}(\gamma) := \lfloor (3/4)x_{\min}^*(\gamma) + (1/4)x_{\max}^*(\gamma) \rfloor \text{ and}$$

$$\text{ec}_{\text{UB}}^{\text{int}}(\gamma) := \lceil (1/4)x_{\min}^*(\gamma) + (3/4)x_{\max}^*(\gamma) \rceil.$$

Set $\mathcal{F}_E := \mathcal{F}(v) := \mathcal{F}_A \cup \mathcal{F}_B$, $v \in V_C$ and $\text{fc}_{\text{LB}}(\psi) := 0, \text{fc}_{\text{UB}}(\psi) := 10, \psi \in \mathcal{F}_A \cup \mathcal{F}_B$.

We include a linear constraint $\ell(a_1) \leq \ell(a_2)$ and $5 \leq \ell(a_1) + \ell(a_2) \leq 15$ as part of the side constraint.

D All Constraints in an MILP Formulation for Chemical Graphs

We define a standard encoding of a finite set A of elements to be a bijection $\sigma : A \rightarrow [1, |A|]$, where we denote by $[A]$ the set $[1, |A|]$ of integers and by $[\mathbf{e}]$ the encoded element $\sigma(\mathbf{e})$. Let ϵ denote *null*, a fictitious chemical element that does not belong to any set of chemical elements, chemical symbols, adjacency-configurations and edge-configurations in the following formulation. Given a finite set A , let A_ϵ denote the set $A \cup \{\epsilon\}$ and define a standard encoding of A_ϵ to be a bijection $\sigma : A \rightarrow [0, |A|]$ such that $\sigma(\epsilon) = 0$, where we denote by $[A_\epsilon]$ the set $[0, |A|]$ of integers and by $[\mathbf{e}]$ the encoded element $\sigma(\mathbf{e})$, where $[\epsilon] = 0$.

Let $\sigma = (G_C, \sigma_{\text{int}}, \sigma_{\text{ce}})$ be a target specification, ρ denote the branch-parameter in the specification σ and \mathbb{C} denote a chemical graph in $\mathcal{G}(G_C, \sigma_{\text{int}}, \sigma_{\text{ce}})$.

D.1 Selecting a Cyclical-base

Recall that

$$E_{(=1)} = \{e \in E_C \mid \ell_{\text{LB}}(e) = \ell_{\text{UB}}(e) = 1\}; \quad E_{(0/1)} = \{e \in E_C \mid \ell_{\text{LB}}(e) = 0, \ell_{\text{UB}}(e) = 1\};$$

$$E_{(\geq 1)} = \{e \in E_C \mid \ell_{\text{LB}}(e) = 1, \ell_{\text{UB}}(e) \geq 2\}; \quad E_{(\geq 2)} = \{e \in E_C \mid \ell_{\text{LB}}(e) \geq 2\};$$

- Every edge $a_i \in E_{(=1)}$ is included in $\langle \mathbb{C} \rangle$;
- Each edge $a_i \in E_{(0/1)}$ is included in $\langle \mathbb{C} \rangle$ if necessary;
- For each edge $a_i \in E_{(\geq 2)}$, edge a_i is not included in $\langle \mathbb{C} \rangle$ and instead a path

$$P_i = (v_{\text{tail}(i)}^{\text{C}}, v_{j-1}^{\text{T}}, v_j^{\text{T}}, \dots, v_{j+t}^{\text{T}}, v_{\text{head}(i)}^{\text{C}})$$

of length at least 2 from vertex $v_{\text{tail}(i)}^{\text{C}}$ to vertex $v_{\text{head}(i)}^{\text{C}}$ visiting some vertices in V_{T} is constructed in $\langle \mathbb{C} \rangle$; and

- For each edge $a_i \in E_{(\geq 1)}$, either edge a_i is directly used in $\langle \mathbb{C} \rangle$ or the above path P_i of length at least 2 is constructed in $\langle \mathbb{C} \rangle$.

Let $t_C \triangleq |V_C|$ and denote V_C by $\{v^C_i \mid i \in [1, t_C]\}$. Regard the seed graph G_C as a digraph such that each edge a_i with end-vertices v^C_j and $v^C_{j'}$ is directed from v^C_j to $v^C_{j'}$ when $j < j'$. For each directed edge $a_i \in E_C$, let $\text{head}(i)$ and $\text{tail}(i)$ denote the head and tail of $e^C(i)$; i.e., $a_i = (v^C_{\text{tail}(i)}, v^C_{\text{head}(i)})$.

Define

$$k_C \triangleq |E_{(\geq 2)} \cup E_{(\geq 1)}|, \quad \widetilde{k}_C \triangleq |E_{(\geq 2)}|,$$

and denote $E_C = \{a_i \mid i \in [1, m_C]\}$, $E_{(\geq 2)} = \{a_k \mid k \in [1, \widetilde{k}_C]\}$, $E_{(\geq 1)} = \{a_k \mid k \in [\widetilde{k}_C + 1, k_C]\}$, $E_{(0/1)} = \{a_i \mid i \in [k_C + 1, k_C + |E_{(0/1)}|]\}$ and $E_{(=1)} = \{a_i \mid i \in [k_C + |E_{(0/1)}| + 1, m_C]\}$. Let $I_{(=1)}$ denote the set of indices i of edges $a_i \in E_{(=1)}$. Similarly for $I_{(0/1)}$, $I_{(\geq 1)}$ and $I_{(\geq 2)}$.

To control the construction of such a path P_i for each edge $a_k \in E_{(\geq 2)} \cup E_{(\geq 1)}$, we regard the index $k \in [1, k_C]$ of each edge $a_k \in E_{(\geq 2)} \cup E_{(\geq 1)}$ as the ‘‘color’’ of the edge. To introduce necessary linear constraints that can construct such a path P_k properly in our MILP, we assign the color k to the vertices $v^T_{j-1}, v^T_j, \dots, v^T_{j+t}$ in V_T when the above path P_k is used in $\langle \mathbb{C} \rangle$.

For each index $s \in [1, t_C]$, let $I_C(s)$ denote the set of edges $e \in E_C$ incident to vertex v^C_s , and $E_{(=1)}^+(s)$ (resp., $E_{(=1)}^-(s)$) denote the set of edges $a_i \in E_{(=1)}$ such that the tail (resp., head) of a_i is vertex v^C_s . Similarly for $E_{(0/1)}^+(s)$, $E_{(0/1)}^-(s)$, $E_{(\geq 1)}^+(s)$, $E_{(\geq 1)}^-(s)$, $E_{(\geq 2)}^+(s)$ and $E_{(\geq 2)}^-(s)$. Let $I_C(s)$ denote the set of indices i of edges $a_i \in I_C(s)$. Similarly for $I_{(=1)}^+(s)$, $I_{(=1)}^-(s)$, $I_{(0/1)}^+(s)$, $I_{(0/1)}^-(s)$, $I_{(\geq 1)}^+(s)$, $I_{(\geq 1)}^-(s)$, $I_{(\geq 2)}^+(s)$ and $I_{(\geq 2)}^-(s)$. Note that $[1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$ and $[\widetilde{k}_C + 1, m_C] = I_{(\geq 1)} \cup I_{(0/1)} \cup I_{(=1)}$.

constants:

- $t_C = |V_C|$, $\widetilde{k}_C = |E_{(\geq 2)}|$, $k_C = |E_{(\geq 2)} \cup E_{(\geq 1)}|$, $t_T = n_{\text{UB}}^{\text{int}} - |V_C|$, $m_C = |E_C|$. Note that $a_i \in E_C \setminus (E_{(\geq 2)} \cup E_{(\geq 1)})$ holds $i \in [k_C + 1, m_C]$;
- $\ell_{\text{LB}}(k), \ell_{\text{UB}}(k) \in [1, t_T]$, $k \in [1, k_C]$: lower and upper bounds on the length of path P_k ;
- $r_{G_C} \in [1, m_C]$: the rank $r(G_C)$ of seed graph G_C ;

variables:

- $e^C(i) \in [0, 1]$, $i \in [1, m_C]$: $e^C(i)$ represents edge $a_i \in E_C$, $i \in [1, m_C]$ ($e^C(i) = 1$, $i \in I_{(=1)}$; $e^C(i) = 0$, $i \in I_{(\geq 2)}$) ($e^C(i) = 1 \Leftrightarrow$ edge a_i is used in $\langle \mathbb{C} \rangle$);
- $v^T(i) \in [0, 1]$, $i \in [1, t_T]$: $v^T(i) = 1 \Leftrightarrow$ vertex v^T_i is used in $\langle \mathbb{C} \rangle$;
- $e^T(i) \in [0, 1]$, $i \in [1, t_T + 1]$: $e^T(i)$ represents edge $e^T_i = (v^T_{i-1}, v^T_i) \in E_T$, where e^T_1 and $e^T_{t_T+1}$ are fictitious edges ($e^T(i) = 1 \Leftrightarrow$ edge e^T_i is used in $\langle \mathbb{C} \rangle$);
- $\chi^T(i) \in [0, k_C]$, $i \in [1, t_T]$: $\chi^T(i)$ represents the color assigned to vertex v^T_i ($\chi^T(i) = k > 0 \Leftrightarrow$ vertex v^T_i is assigned color k ; $\chi^T(i) = 0$ means that vertex v^T_i is not used in $\langle \mathbb{C} \rangle$);
- $\text{clr}^T(k) \in [\ell_{\text{LB}}(k) - 1, \ell_{\text{UB}}(k) - 1]$, $k \in [1, k_C]$, $\text{clr}^T(0) \in [0, t_T]$: the number of vertices $v^T_i \in V_T$ with color c ;

- $\delta_x^T(k) \in [0, 1]$, $k \in [0, k_C]$: $\delta_x^T(k) = 1 \Leftrightarrow \chi^T(i) = k$ for some $i \in [1, t_T]$;
- $\chi^T(i, k) \in [0, 1]$, $i \in [1, t_T]$, $k \in [0, k_C]$ ($\chi^T(i, k) = 1 \Leftrightarrow \chi^T(i) = k$);
- $\widetilde{\text{deg}}_C^+(i) \in [0, 4]$, $i \in [1, t_C]$: the out-degree of vertex v^C_i with the used edges e^C in E_C ;
- $\widetilde{\text{deg}}_C^-(i) \in [0, 4]$, $i \in [1, t_C]$: the in-degree of vertex v^C_i with the used edges e^C in E_C ;
- rank: the rank $r(\mathbb{C})$ of a target chemical graph \mathbb{C} ;

constraints:

$$\text{rank} = r_{G_C} - \sum_{i \in I_{(0/1)}} (1 - e^C(i)), \quad (6)$$

$$e^C(i) = 1, \quad i \in I_{(=1)}, \quad (7)$$

$$e^C(i) = 0, \quad \text{clr}^T(i) \geq 1, \quad i \in I_{(\geq 2)}, \quad (8)$$

$$e^C(i) + \text{clr}^T(i) \geq 1, \quad \text{clr}^T(i) \leq t_T \cdot (1 - e^C(i)), \quad i \in I_{(\geq 1)}, \quad (9)$$

$$\sum_{c \in I_{(\geq 1)}^-(i) \cup I_{(0/1)}^-(i) \cup I_{(=1)}^-(i)} e^C(c) = \widetilde{\text{deg}}_C^-(i), \quad \sum_{c \in I_{(\geq 1)}^+(i) \cup I_{(0/1)}^+(i) \cup I_{(=1)}^+(i)} e^C(c) = \widetilde{\text{deg}}_C^+(i), \quad i \in [1, t_C], \quad (10)$$

$$\chi^T(i, 0) = 1 - v^T(i), \quad \sum_{k \in [0, k_C]} \chi^T(i, k) = 1, \quad \sum_{k \in [0, k_C]} k \cdot \chi^T(i, k) = \chi^T(i), \quad i \in [1, t_T], \quad (11)$$

$$\sum_{i \in [1, t_T]} \chi^T(i, k) = \text{clr}^T(k), \quad t_T \cdot \delta_x^T(k) \geq \sum_{i \in [1, t_T]} \chi^T(i, k) \geq \delta_x^T(k), \quad k \in [0, k_C], \quad (12)$$

$$v^T(i-1) \geq v^T(i), \quad k_C \cdot (v^T(i-1) - e^T(i)) \geq \chi^T(i-1) - \chi^T(i) \geq v^T(i-1) - e^T(i), \quad i \in [2, t_T]. \quad (13)$$

D.2 Constraints for Including Leaf Paths

Let \tilde{t}_C denote the number of vertices $u \in V_C$ such that $\text{bl}_{\text{UB}}(u) = 1$ and assume that $V_C = \{u_1, u_2, \dots, u_p\}$ so that

$$\text{bl}_{\text{UB}}(u_i) = 1, \quad i \in [1, \tilde{t}_C] \quad \text{and} \quad \text{bl}_{\text{UB}}(u_i) = 0, \quad i \in [\tilde{t}_C + 1, t_C].$$

Define the set of colors for the vertex set $\{u_i \mid i \in [1, \tilde{t}_C]\} \cup V_T$ to be $[1, c_F]$ with

$$c_F \triangleq \tilde{t}_C + t_T = |\{u_i \mid i \in [1, \tilde{t}_C]\} \cup V_T|.$$

Let each vertex v^C_i , $i \in [1, \tilde{t}_C]$ (resp., $v^T_i \in V_T$) correspond to a color $i \in [1, c_F]$ (resp., $i + \tilde{t}_C \in [1, c_F]$). When a path $P = (u, v^F_j, v^F_{j+1}, \dots, v^F_{j+t})$ from a vertex $u \in V_C \cup V_T$ is used in $\langle \mathbb{C} \rangle$, we assign the color $i \in [1, c_F]$ of the vertex u to the vertices $v^F_j, v^F_{j+1}, \dots, v^F_{j+t} \in V_F$.

constants:

- c_F : the maximum number of different colors assigned to the vertices in V_F ;
- n^* : an upper bound on the number $n(\mathbb{C})$ of non-hydrogen atoms in \mathbb{C} ;
- $n_{LB}^{\text{int}}, n_{UB}^{\text{int}} \in [2, n^*]$: lower and upper bounds on the number of interior-vertices in \mathbb{C} ;
- $\text{bl}_{LB}(i) \in [0, 1]$, $i \in [1, \tilde{t}_C]$: a lower bound on the number of leaf ρ -branches in the leaf path rooted at a vertex v^C_i ;
- $\text{bl}_{LB}(k), \text{bl}_{UB}(k) \in [0, \ell_{UB}(k) - 1]$, $k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$: lower and upper bounds on the number of leaf ρ -branches in the trees rooted at internal vertices of a pure path P_k for an edge $a_k \in E_{(\geq 1)} \cup E_{(\geq 2)}$;

variables:

- $n_G^{\text{int}} \in [n_{LB}^{\text{int}}, n_{UB}^{\text{int}}]$: the number of interior-vertices in \mathbb{C} ;
- $v^F(i) \in [0, 1]$, $i \in [1, t_F]$: $v^F(i) = 1 \Leftrightarrow$ vertex v^F_i is used in \mathbb{C} ;
- $e^F(i) \in [0, 1]$, $i \in [1, t_F + 1]$: $e^F(i)$ represents edge $e^F_i = v^F_{i-1}v^F_i$, where e^F_1 and $e^F_{t_F+1}$ are fictitious edges ($e^F(i) = 1 \Leftrightarrow$ edge e^F_i is used in \mathbb{C});
- $\chi^F(i) \in [0, c_F]$, $i \in [1, t_F]$: $\chi^F(i)$ represents the color assigned to vertex v^F_i ($\chi^F(i) = c \Leftrightarrow$ vertex v^F_i is assigned color c);
- $\text{clr}^F(c) \in [0, t_F]$, $c \in [0, c_F]$: the number of vertices v^F_i with color c ;
- $\delta_\chi^F(c) \in [\text{bl}_{LB}(c), 1]$, $c \in [1, \tilde{t}_C]$: $\delta_\chi^F(c) = 1 \Leftrightarrow \chi^F(i) = c$ for some $i \in [1, t_F]$;
- $\delta_\chi^F(c) \in [0, 1]$, $c \in [\tilde{t}_C + 1, c_F]$: $\delta_\chi^F(c) = 1 \Leftrightarrow \chi^F(i) = c$ for some $i \in [1, t_F]$;
- $\chi^F(i, c) \in [0, 1]$, $i \in [1, t_F]$, $c \in [0, c_F]$: $\chi^F(i, c) = 1 \Leftrightarrow \chi^F(i) = c$;
- $\text{bl}(k, i) \in [0, 1]$, $k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$, $i \in [1, t_T]$: $\text{bl}(k, i) = 1 \Leftrightarrow$ path P_k contains vertex v^T_i as an internal vertex and the ρ -fringe-tree rooted at v^T_i contains a leaf ρ -branch;

constraints:

$$\chi^F(i, 0) = 1 - v^F(i), \quad \sum_{c \in [0, c_F]} \chi^F(i, c) = 1, \quad \sum_{c \in [0, c_F]} c \cdot \chi^F(i, c) = \chi^F(i), \quad i \in [1, t_F], \quad (14)$$

$$\sum_{i \in [1, t_F]} \chi^F(i, c) = \text{clr}^F(c), \quad t_F \cdot \delta_\chi^F(c) \geq \sum_{i \in [1, t_F]} \chi^F(i, c) \geq \delta_\chi^F(c), \quad c \in [0, c_F], \quad (15)$$

$$e^F(1) = e^F(t_F + 1) = 0, \quad (16)$$

$$\begin{aligned}
& v^{\text{F}}(i-1) \geq v^{\text{F}}(i), \\
c_{\text{F}} \cdot (v^{\text{F}}(i-1) - e^{\text{F}}(i)) & \geq \chi^{\text{F}}(i-1) - \chi^{\text{F}}(i) \geq v^{\text{F}}(i-1) - e^{\text{F}}(i), \quad i \in [2, t_{\text{F}}], \quad (17)
\end{aligned}$$

$$\text{bl}(k, i) \geq \delta_{\chi}^{\text{F}}(\tilde{t}_{\text{C}} + i) + \chi^{\text{T}}(i, k) - 1, \quad k \in [1, k_{\text{C}}], i \in [1, t_{\text{T}}], \quad (18)$$

$$\sum_{k \in [1, k_{\text{C}}], i \in [1, t_{\text{T}}]} \text{bl}(k, i) \leq \sum_{i \in [1, t_{\text{T}}]} \delta_{\chi}^{\text{F}}(\tilde{t}_{\text{C}} + i), \quad (19)$$

$$\text{bl}_{\text{LB}}(k) \leq \sum_{i \in [1, t_{\text{T}}]} \text{bl}(k, i) \leq \text{bl}_{\text{UB}}(k), \quad k \in [1, k_{\text{C}}], \quad (20)$$

$$t_{\text{C}} + \sum_{i \in [1, t_{\text{T}}]} v^{\text{T}}(i) + \sum_{i \in [1, t_{\text{F}}]} v^{\text{F}}(i) = n_{\text{G}}^{\text{int}}. \quad (21)$$

D.3 Constraints for Including Fringe-trees

Recall that $\mathcal{F}(D_{\pi})$ denotes the set of chemical rooted trees ψ r-isomorphic to a chemical rooted tree in $\mathcal{T}(\mathbb{C})$ over all chemical graphs $\mathbb{C} \in D_{\pi}$, where possibly a chemical rooted tree $\psi \in \mathcal{F}(D_{\pi})$ consists of a single chemical element $\mathbf{a} \in \Lambda \setminus \{\text{H}\}$.

To express the condition that the ρ -fringe-tree is chosen from a rooted tree C_i , T_i or F_i , we introduce the following set of variables and constraints.

constants:

- n_{LB} : a lower bound on the number $n(\mathbb{C})$ of non-hydrogen atoms in \mathbb{C} , where $n_{\text{LB}}, n^* \geq n_{\text{LB}}^{\text{int}}$;
- $\text{ch}_{\text{LB}}(i), \text{ch}_{\text{UB}}(i) \in [0, n^*]$, $i \in [1, t_{\text{T}}]$: lower and upper bounds on $\text{ht}(\langle T_i \rangle)$ of the tree T_i rooted at a vertex v_{C}^i ;
- $\text{ch}_{\text{LB}}(k), \text{ch}_{\text{UB}}(k) \in [0, n^*]$, $k \in [1, k_{\text{C}}] = I_{(\geq 2)} \cup I_{(\geq 1)}$: lower and upper bounds on the maximum height $\text{ht}(\langle T \rangle)$ of the tree $T \in \mathcal{F}(P_k)$ rooted at an internal vertex of a path P_k for an edge $a_k \in E_{(\geq 1)} \cup E_{(\geq 2)}$;
- Prepare a coding of the set $\mathcal{F}(D_{\pi})$ and let $[\psi]$ denote the coded integer of an element ψ in $\mathcal{F}(D_{\pi})$;
- Sets $\mathcal{F}(v) \subseteq \mathcal{F}(D_{\pi})$, $v \in V_{\text{C}}$ and $\mathcal{F}_E \subseteq \mathcal{F}(D_{\pi})$ of chemical rooted trees T with $\text{ht}(T) \in [1, \rho]$;
- Define $\mathcal{F}^* := \bigcup_{v \in V_{\text{C}}} \mathcal{F}(v) \cup \mathcal{F}_E$, $\mathcal{F}_i^{\text{C}} := \mathcal{F}(v_{\text{C}}^i)$, $i \in [1, t_{\text{C}}]$, $\mathcal{F}_i^{\text{T}} := \mathcal{F}_E$, $i \in [1, t_{\text{T}}]$ and $\mathcal{F}_i^{\text{F}} := \mathcal{F}_E$, $i \in [1, t_{\text{F}}]$;

- $\text{fc}_{\text{LB}}(\psi), \text{fc}_{\text{UB}}(\psi) \in [0, n^*], \psi \in \mathcal{F}^*$: lower and upper bound functions on the number of interior-vertices v such that $\mathbb{C}[v]$ is r-isomorphic to ψ in \mathbb{C} ;
- $\mathcal{F}_i^X[p], p \in [1, \rho], X \in \{\text{C}, \text{T}, \text{F}\}$: the set of chemical rooted trees $T \in \mathcal{F}_i^X$ with $\text{ht}(\langle T \rangle) = p$;
- $n_{\bar{\text{H}}}([\psi]) \in [0, 3^\rho], \psi \in \mathcal{F}^*$: the number $n(\langle \psi \rangle)$ of non-root hydrogen vertices in a chemical rooted tree ψ ;
- $\text{ht}_{\bar{\text{H}}}([\psi]) \in [0, \rho], \psi \in \mathcal{F}^*$: the height $\text{ht}(\langle \psi \rangle)$ of the hydrogen-suppressed chemical rooted tree $\langle \psi \rangle$;
- $\text{deg}_{\text{r}}^{\bar{\text{H}}}([\psi]) \in [0, 3], \psi \in \mathcal{F}^*$: the number $\text{deg}_{\text{r}}(\langle \psi \rangle)$ of non-hydrogen children of the root r of a chemical rooted tree ψ ;
- $\text{deg}_{\text{r}}^{\text{hyd}}([\psi]) \in [0, 3], \psi \in \mathcal{F}^*$: the number $\text{deg}_{\text{r}}(\psi) - \text{deg}_{\text{r}}(\langle \psi \rangle)$ of hydrogen children of the root r of a chemical rooted tree ψ ;
- $v_{\text{ion}}(\psi) \in [-3, +3], \psi \in \mathcal{F}^*$: the ion-valence of the root in ψ ;
- $\text{ac}_{\nu}^{\text{lf}}(\psi), \nu \in \Gamma_{\text{ac}}^{\text{lf}}$: the frequency of leaf-edges with adjacency-configuration ν in ψ ;
- $\text{ac}_{\text{LB}}^{\text{lf}}, \text{ac}_{\text{UB}}^{\text{lf}} : \Gamma_{\text{ac}}^{\text{lf}} \rightarrow [0, n^*]$: lower and upper bound functions on the number of leaf-edges uv in ac_{C} with adjacency-configuration ν ;

variables:

- $n_{\text{G}} \in [n_{\text{LB}}, n^*]$: the number $n(\mathbb{C})$ of non-hydrogen atoms in \mathbb{C} ;
- $v^{\text{X}}(i) \in [0, 1], i \in [1, t_{\text{X}}], X \in \{\text{T}, \text{F}\}$: $v^{\text{X}}(i) = 1 \Leftrightarrow$ vertex v^{X}_i is used in \mathbb{C} ;
- $\delta_{\text{fr}}^{\text{X}}(i, [\psi]) \in [0, 1], i \in [1, t_{\text{X}}], \psi \in \mathcal{F}_i^{\text{X}}, X \in \{\text{C}, \text{T}, \text{F}\}$: $\delta_{\text{fr}}^{\text{X}}(i, [\psi]) = 1 \Leftrightarrow \psi$ is the ρ -fringe-tree rooted at vertex v^{X}_i in \mathbb{C} ;
- $\text{fc}([\psi]) \in [\text{fc}_{\text{LB}}(\psi), \text{fc}_{\text{UB}}(\psi)], \psi \in \mathcal{F}^*$: the number of interior-vertices v such that $\mathbb{C}[v]$ is r-isomorphic to ψ in \mathbb{C} ;
- $\text{ac}^{\text{lf}}([\nu]) \in [\text{ac}_{\text{LB}}^{\text{lf}}(\nu), \text{ac}_{\text{UB}}^{\text{lf}}(\nu)], \nu \in \Gamma_{\text{ac}}^{\text{lf}}$: the number of leaf-edge with adjacency-configuration ν in \mathbb{C} ;
- $\text{deg}_{\text{X}}^{\text{ex}}(i) \in [0, 3], i \in [1, t_{\text{X}}], X \in \{\text{C}, \text{T}, \text{F}\}$: the number of non-hydrogen children of the root of the ρ -fringe-tree rooted at vertex v^{X}_i in \mathbb{C} ;
- $\text{hyddeg}^{\text{X}}(i) \in [0, 4], i \in [1, t_{\text{X}}], X \in \{\text{C}, \text{T}, \text{F}\}$: the number of hydrogen atoms adjacent to vertex v^{X}_i (i.e., $\text{hyddeg}(v^{\text{X}}_i)$) in $\mathbb{C} = (H, \alpha, \beta)$;
- $\text{eledeg}_{\text{X}}(i) \in [-3, +3], i \in [1, t_{\text{X}}], X \in \{\text{C}, \text{T}, \text{F}\}$: the ion-valence $v_{\text{ion}}(\psi)$ of vertex v^{X}_i (i.e., $\text{eledeg}_{\text{X}}(i) = v_{\text{ion}}(\psi)$ for the ρ -fringe-tree ψ rooted at v^{X}_i) in $\mathbb{C} = (H, \alpha, \beta)$;
- $h^{\text{X}}(i) \in [0, \rho], i \in [1, t_{\text{X}}], X \in \{\text{C}, \text{T}, \text{F}\}$: the height $\text{ht}(\langle T \rangle)$ of the hydrogen-suppressed chemical rooted tree $\langle T \rangle$ of the ρ -fringe-tree T rooted at vertex v^{X}_i in \mathbb{C} ;

- $\sigma(k, i) \in [0, 1]$, $k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$, $i \in [1, t_T]$: $\sigma(k, i) = 1 \Leftrightarrow$ the ρ -fringe-tree T_v rooted at vertex $v = v_i^T$ with color k has the largest height $\text{ht}(\langle \mathcal{T}_v \rangle)$ among such trees $T_v, v \in V_T$;

constraints:

$$\begin{aligned} \sum_{\psi \in \mathcal{F}_i^C} \delta_{\text{fr}}^C(i, [\psi]) &= 1, & i \in [1, t_C], \\ \sum_{\psi \in \mathcal{F}_i^X} \delta_{\text{fr}}^X(i, [\psi]) &= v^X(i), & i \in [1, t_X], X \in \{T, F\}, \end{aligned} \quad (22)$$

$$\begin{aligned} \sum_{\psi \in \mathcal{F}_i^X} \text{deg}_{\text{r}}^{\bar{H}}([\psi]) \cdot \delta_{\text{fr}}^X(i, [\psi]) &= \text{deg}_X^{\text{ex}}(i), \\ \sum_{\psi \in \mathcal{F}_i^X} \text{deg}_{\text{r}}^{\text{hyd}}([\psi]) \cdot \delta_{\text{fr}}^X(i, [\psi]) &= \text{hyddeg}^X(i), \\ \sum_{\psi \in \mathcal{F}_i^X} v_{\text{ion}}([\psi]) \cdot \delta_{\text{fr}}^X(i, [\psi]) &= \text{eledeg}_X(i), & i \in [1, t_X], X \in \{C, T, F\}, \end{aligned} \quad (23)$$

$$\sum_{\psi \in \mathcal{F}_i^F[\rho]} \delta_{\text{fr}}^F(i, [\psi]) \geq v^F(i) - e^F(i+1), \quad i \in [1, t_F] \ (e^F(t_F+1) = 0), \quad (24)$$

$$\sum_{\psi \in \mathcal{F}_i^X} \text{ht}_{\bar{H}}([\psi]) \cdot \delta_{\text{fr}}^X(i, [\psi]) = h^X(i), \quad i \in [1, t_X], X \in \{C, T, F\}, \quad (25)$$

$$\sum_{\substack{\psi \in \mathcal{F}_i^X \\ i \in [1, t_X], X \in \{C, T, F\}}} n_{\bar{H}}([\psi]) \cdot \delta_{\text{fr}}^X(i, [\psi]) + \sum_{i \in [1, t_X], X \in \{T, F\}} v^X(i) + t_C = n_G, \quad (26)$$

$$\sum_{i \in [1, t_X], X \in \{C, T, F\}} \delta_{\text{fr}}^X(i, [\psi]) = \text{fc}([\psi]), \quad \psi \in \mathcal{F}^*, \quad (27)$$

$$\sum_{\psi \in \mathcal{F}_i^X, i \in [1, t_X], X \in \{C, T, F\}} \text{ac}_{\nu}^{\text{lf}}(\psi) \cdot \delta_{\text{fr}}^X(i, [\psi]) = \text{ac}^{\text{lf}}([\nu]), \quad \nu \in \Gamma_{\text{ac}}^{\text{lf}}, \quad (28)$$

$$\begin{aligned} h^C(i) &\geq \text{ch}_{\text{LB}}(i) - n^* \cdot \delta_{\chi}^F(i), \quad \text{clr}^F(i) + \rho \geq \text{ch}_{\text{LB}}(i), \\ h^C(i) &\leq \text{ch}_{\text{UB}}(i), \quad \text{clr}^F(i) + \rho \leq \text{ch}_{\text{UB}}(i) + n^* \cdot (1 - \delta_{\chi}^F(i)), & i \in [1, \tilde{t}_C], \end{aligned} \quad (29)$$

$$\text{ch}_{\text{LB}}(i) \leq h^C(i) \leq \text{ch}_{\text{UB}}(i), \quad i \in [\tilde{t}_C + 1, t_C], \quad (30)$$

$$\begin{aligned}
h^T(i) &\leq \text{ch}_{\text{UB}}(k) + n^* \cdot (\delta_{\chi}^{\text{F}}(\tilde{t}_{\text{C}} + i) + 1 - \chi^{\text{T}}(i, k)), \\
\text{clr}^{\text{F}}(\tilde{t}_{\text{C}} + i) + \rho &\leq \text{ch}_{\text{UB}}(k) + n^* \cdot (2 - \delta_{\chi}^{\text{F}}(\tilde{t}_{\text{C}} + i) - \chi^{\text{T}}(i, k)), \quad k \in [1, k_{\text{C}}], i \in [1, t_{\text{T}}], \quad (31)
\end{aligned}$$

$$\sum_{i \in [1, t_{\text{T}}]} \sigma(k, i) = \delta_{\chi}^{\text{T}}(k), \quad k \in [1, k_{\text{C}}], \quad (32)$$

$$\begin{aligned}
\chi^{\text{T}}(i, k) &\geq \sigma(k, i), \\
h^T(i) &\geq \text{ch}_{\text{LB}}(k) - n^* \cdot (\delta_{\chi}^{\text{F}}(\tilde{t}_{\text{C}} + i) + 1 - \sigma(k, i)), \\
\text{clr}^{\text{F}}(\tilde{t}_{\text{C}} + i) + \rho &\geq \text{ch}_{\text{LB}}(k) - n^* \cdot (2 - \delta_{\chi}^{\text{F}}(\tilde{t}_{\text{C}} + i) - \sigma(k, i)), \quad k \in [1, k_{\text{C}}], i \in [1, t_{\text{T}}]. \quad (33)
\end{aligned}$$

D.4 Descriptor for the Number of Specified Degree

We include constraints to compute descriptors for degrees in \mathbb{C} .

variables:

- $\text{deg}^{\text{X}}(i) \in [0, 4]$, $i \in [1, t_{\text{X}}]$, $\text{X} \in \{\text{C}, \text{T}, \text{F}\}$: the number of non-hydrogen atoms adjacent to vertex $v = v^{\text{X}}_i$ (i.e., $\text{deg}_{\mathbb{C}}(v) = \text{deg}_H(v) - \text{hyddeg}_{\mathbb{C}}(v)$) in $\mathbb{C} = (H, \alpha, \beta)$;
- $\text{deg}_{\text{CT}}(i) \in [0, 4]$, $i \in [1, t_{\text{C}}]$: the number of edges from vertex v^{C}_i to vertices v^{T}_j , $j \in [1, t_{\text{T}}]$;
- $\text{deg}_{\text{TC}}(i) \in [0, 4]$, $i \in [1, t_{\text{C}}]$: the number of edges from vertices v^{T}_j , $j \in [1, t_{\text{T}}]$ to vertex v^{C}_i ;
- $\delta_{\text{dg}}^{\text{C}}(i, d) \in [0, 1]$, $i \in [1, t_{\text{C}}]$, $d \in [1, 4]$, $\delta_{\text{dg}}^{\text{X}}(i, d) \in [0, 1]$, $i \in [1, t_{\text{X}}]$, $d \in [0, 4]$, $\text{X} \in \{\text{T}, \text{F}\}$:
 $\delta_{\text{dg}}^{\text{X}}(i, d) = 1 \Leftrightarrow \text{deg}^{\text{X}}(i) + \text{hyddeg}^{\text{X}}(i) = d$;
- $\text{dg}(d) \in [\text{dg}_{\text{LB}}(d), \text{dg}_{\text{UB}}(d)]$, $d \in [1, 4]$: the number of interior-vertices v with $\text{deg}_H(v^{\text{X}}_i) = d$ in $\mathbb{C} = (H, \alpha, \beta)$;
- $\text{deg}_{\text{C}}^{\text{int}}(i) \in [1, 4]$, $i \in [1, t_{\text{C}}]$, $\text{deg}_{\text{X}}^{\text{int}}(i) \in [0, 4]$, $i \in [1, t_{\text{X}}]$, $\text{X} \in \{\text{T}, \text{F}\}$: the interior-degree $\text{deg}_{H^{\text{int}}}(v^{\text{X}}_i)$ in the interior $H^{\text{int}} = (V^{\text{int}}(\mathbb{C}), E^{\text{int}}(\mathbb{C}))$ of \mathbb{C} ; i.e., the number of interior-edges incident to vertex v^{X}_i ;
- $\delta_{\text{dg}, \text{C}}^{\text{int}}(i, d) \in [0, 1]$, $i \in [1, t_{\text{C}}]$, $d \in [1, 4]$, $\delta_{\text{dg}, \text{X}}^{\text{int}}(i, d) \in [0, 1]$, $i \in [1, t_{\text{X}}]$, $d \in [0, 4]$, $\text{X} \in \{\text{T}, \text{F}\}$:
 $\delta_{\text{dg}, \text{X}}^{\text{int}}(i, d) = 1 \Leftrightarrow \text{deg}_{\text{X}}^{\text{int}}(i) = d$;
- $\text{dg}^{\text{int}}(d) \in [\text{dg}_{\text{LB}}(d), \text{dg}_{\text{UB}}(d)]$, $d \in [1, 4]$: the number of interior-vertices v with the interior-degree $\text{deg}_{H^{\text{int}}}(v) = d$ in the interior $H^{\text{int}} = (V^{\text{int}}(\mathbb{C}), E^{\text{int}}(\mathbb{C}))$ of $\mathbb{C} = (H, \alpha, \beta)$.

constraints:

$$\sum_{k \in I_{(\geq 2)}^+(i) \cup I_{(\geq 1)}^+(i)} \delta_{\chi}^{\text{T}}(k) = \text{deg}_{\text{CT}}(i), \quad \sum_{k \in I_{(\geq 2)}^-(i) \cup I_{(\geq 1)}^-(i)} \delta_{\chi}^{\text{T}}(k) = \text{deg}_{\text{TC}}(i), \quad i \in [1, t_{\text{C}}], \quad (34)$$

$$\widetilde{\deg}_C^-(i) + \widetilde{\deg}_C^+(i) + \deg_{CT}(i) + \deg_{TC}(i) + \delta_\chi^F(i) = \deg_C^{\text{int}}(i), \quad i \in [1, \widetilde{t}_C], \quad (35)$$

$$\widetilde{\deg}_C^-(i) + \widetilde{\deg}_C^+(i) + \deg_{CT}(i) + \deg_{TC}(i) = \deg_C^{\text{int}}(i), \quad i \in [\widetilde{t}_C + 1, t_C], \quad (36)$$

$$\deg_C^{\text{int}}(i) + \deg_C^{\text{ex}}(i) = \deg_C^{\text{C}}(i), \quad i \in [1, t_C], \quad (37)$$

$$\sum_{\psi \in \mathcal{F}_i^{\text{C}}[\rho]} \delta_{\text{fr}}^{\text{C}}(i, [\psi]) \geq 2 - \deg_C^{\text{int}}(i) \quad i \in [1, t_C], \quad (38)$$

$$\begin{aligned} 2v^{\text{T}}(i) + \delta_\chi^{\text{F}}(\widetilde{t}_C + i) &= \deg_{\text{T}}^{\text{int}}(i), \\ \deg_{\text{T}}^{\text{int}}(i) + \deg_{\text{T}}^{\text{ex}}(i) &= \deg_{\text{T}}^{\text{T}}(i), \end{aligned} \quad i \in [1, t_{\text{T}}] \quad (e^{\text{T}}(1) = e^{\text{T}}(t_{\text{T}} + 1) = 0), \quad (39)$$

$$\begin{aligned} v^{\text{F}}(i) + e^{\text{F}}(i + 1) &= \deg_{\text{F}}^{\text{int}}(i), \\ \deg_{\text{F}}^{\text{int}}(i) + \deg_{\text{F}}^{\text{ex}}(i) &= \deg_{\text{F}}^{\text{F}}(i), \end{aligned} \quad i \in [1, t_{\text{F}}] \quad (e^{\text{F}}(1) = e^{\text{F}}(t_{\text{F}} + 1) = 0), \quad (40)$$

$$\begin{aligned} \sum_{d \in [0,4]} \delta_{\text{dg}}^{\text{X}}(i, d) &= 1, \quad \sum_{d \in [1,4]} d \cdot \delta_{\text{dg}}^{\text{X}}(i, d) = \deg^{\text{X}}(i) + \text{hyddeg}^{\text{X}}(i), \\ \sum_{d \in [0,4]} \delta_{\text{dg},\text{X}}^{\text{int}}(i, d) &= 1, \quad \sum_{d \in [1,4]} d \cdot \delta_{\text{dg},\text{X}}^{\text{int}}(i, d) = \deg_{\text{X}}^{\text{int}}(i), \end{aligned} \quad i \in [1, t_{\text{X}}], \text{X} \in \{\text{T}, \text{C}, \text{F}\}, \quad (41)$$

$$\begin{aligned} \sum_{i \in [1, t_{\text{C}}]} \delta_{\text{dg}}^{\text{C}}(i, d) + \sum_{i \in [1, t_{\text{T}}]} \delta_{\text{dg}}^{\text{T}}(i, d) + \sum_{i \in [1, t_{\text{F}}]} \delta_{\text{dg}}^{\text{F}}(i, d) &= \text{dg}(d), \\ \sum_{i \in [1, t_{\text{C}}]} \delta_{\text{dg},\text{C}}^{\text{int}}(i, d) + \sum_{i \in [1, t_{\text{T}}]} \delta_{\text{dg},\text{T}}^{\text{int}}(i, d) + \sum_{i \in [1, t_{\text{F}}]} \delta_{\text{dg},\text{F}}^{\text{int}}(i, d) &= \text{dg}^{\text{int}}(d), \end{aligned} \quad d \in [1, 4]. \quad (42)$$

D.5 Assigning Multiplicity

We prepare an integer variable $\beta(e)$ for each edge e in the scheme graph SG to denote the bond-multiplicity of e in a selected graph H and include necessary constraints for the variables to satisfy in H .

constants:

- $\beta_{\text{T}}([\psi])$: the sum $\beta_\psi(r)$ of bond-multiplicities of edges incident to the root r of a chemical rooted tree $\psi \in \mathcal{F}^*$;

variables:

- $\beta^X(i) \in [0, 3]$, $i \in [2, t_X]$, $X \in \{T, F\}$: the bond-multiplicity of edge e^X_i in \mathbb{C} ;
- $\beta^C(i) \in [0, 3]$, $i \in [\widetilde{k}_C + 1, m_C] = I_{(\geq 1)} \cup I_{(0/1)} \cup I_{(=1)}$: the bond-multiplicity of edge $a_i \in E_{(\geq 1)} \cup E_{(0/1)} \cup E_{(=1)}$ in \mathbb{C} ;
- $\beta^{CT}(k), \beta^{TC}(k) \in [0, 3]$, $k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$: the bond-multiplicity of the first (resp., last) edge of the pure path P_k in \mathbb{C} ;
- $\beta^{*F}(c) \in [0, 3]$, $c \in [1, c_F = \widetilde{t}_C + t_T]$: the bond-multiplicity of the first edge of the leaf path Q_c rooted at vertex v^C_c , $c \leq \widetilde{t}_C$ or $v^T_{c-\widetilde{t}_C}$, $c > \widetilde{t}_C$ in \mathbb{C} ;
- $\beta_{\text{ex}}^X(i) \in [0, 4]$, $i \in [1, t_X]$, $X \in \{C, T, F\}$: the sum $\beta_{\mathbb{C}[v]}(v)$ of bond-multiplicities of edges in the ρ -fringe-tree $\mathbb{C}[v]$ rooted at interior-vertex $v = v^X_i$;
- $\delta_\beta^X(i, m) \in [0, 1]$, $i \in [2, t_X]$, $m \in [0, 3]$, $X \in \{T, F\}$: $\delta_\beta^X(i, m) = 1 \Leftrightarrow \beta^X(i) = m$;
- $\delta_\beta^C(i, m) \in [0, 1]$, $i \in [\widetilde{k}_C, m_C] = I_{(\geq 1)} \cup I_{(0/1)} \cup I_{(=1)}$, $m \in [0, 3]$: $\delta_\beta^C(i, m) = 1 \Leftrightarrow \beta^C(i) = m$;
- $\delta_\beta^{CT}(k, m), \delta_\beta^{TC}(k, m) \in [0, 1]$, $k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$, $m \in [0, 3]$: $\delta_\beta^{CT}(k, m) = 1$ (resp., $\delta_\beta^{TC}(k, m) = 1$) $\Leftrightarrow \beta^{CT}(k) = m$ (resp., $\beta^{TC}(k) = m$);
- $\delta_\beta^{*F}(c, m) \in [0, 1]$, $c \in [1, c_F]$, $m \in [0, 3]$, $X \in \{C, T\}$: $\delta_\beta^{*F}(c, m) = 1 \Leftrightarrow \beta^{*F}(c) = m$;
- $\text{bd}^{\text{int}}(m) \in [0, 2n_{\text{UB}}^{\text{int}}]$, $m \in [1, 3]$: the number of interior-edges with bond-multiplicity m in \mathbb{C} ;
- $\text{bd}_X(m) \in [0, 2n_{\text{UB}}^{\text{int}}]$, $X \in \{C, T, CT, TC\}$, $\text{bd}_X(m) \in [0, 2n_{\text{UB}}^{\text{int}}]$, $X \in \{F, CF, TF\}$, $m \in [1, 3]$: the number of interior-edges $e \in E_X$ with bond-multiplicity m in \mathbb{C} ;

constraints:

$$e^C(i) \leq \beta^C(i) \leq 3e^C(i), i \in [\widetilde{k}_C + 1, m_C] = I_{(\geq 1)} \cup I_{(0/1)} \cup I_{(=1)}, \quad (43)$$

$$e^X(i) \leq \beta^X(i) \leq 3e^X(i), \quad i \in [2, t_X], X \in \{T, F\}, \quad (44)$$

$$\delta_\chi^T(k) \leq \beta^{CT}(k) \leq 3\delta_\chi^T(k), \quad \delta_\chi^T(k) \leq \beta^{TC}(k) \leq 3\delta_\chi^T(k), \quad k \in [1, k_C], \quad (45)$$

$$\delta_\chi^F(c) \leq \beta^{*F}(c) \leq 3\delta_\chi^F(c), \quad c \in [1, c_F], \quad (46)$$

$$\sum_{m \in [0, 3]} \delta_\beta^X(i, m) = 1, \quad \sum_{m \in [0, 3]} m \cdot \delta_\beta^X(i, m) = \beta^X(i), \quad i \in [2, t_X], X \in \{T, F\}, \quad (47)$$

$$\sum_{m \in [0,3]} \delta_{\beta}^{\text{C}}(i, m) = 1, \quad \sum_{m \in [0,3]} m \cdot \delta_{\beta}^{\text{C}}(i, m) = \beta^{\text{C}}(i), \quad i \in [\widetilde{k}_{\text{C}} + 1, m_{\text{C}}], \quad (48)$$

$$\begin{aligned} \sum_{m \in [0,3]} \delta_{\beta}^{\text{CT}}(k, m) = 1, \quad \sum_{m \in [0,3]} m \cdot \delta_{\beta}^{\text{CT}}(k, m) &= \beta^{\text{CT}}(k), & k \in [1, k_{\text{C}}], \\ \sum_{m \in [0,3]} \delta_{\beta}^{\text{TC}}(k, m) = 1, \quad \sum_{m \in [0,3]} m \cdot \delta_{\beta}^{\text{TC}}(k, m) &= \beta^{\text{TC}}(k), & k \in [1, k_{\text{C}}], \\ \sum_{m \in [0,3]} \delta_{\beta}^{\text{F}}(c, m) = 1, \quad \sum_{m \in [0,3]} m \cdot \delta_{\beta}^{\text{F}}(c, m) &= \beta^{\text{F}}(c), & c \in [1, c_{\text{F}}], \end{aligned} \quad (49)$$

$$\sum_{\psi \in \mathcal{F}_i^{\text{X}}} \beta_{\text{r}}([\psi]) \cdot \delta_{\text{fr}}^{\text{X}}(i, [\psi]) = \beta_{\text{ex}}^{\text{X}}(i), \quad i \in [1, t_{\text{X}}], \text{X} \in \{\text{C}, \text{T}, \text{F}\}, \quad (50)$$

$$\begin{aligned} \sum_{i \in [\widetilde{k}_{\text{C}} + 1, m_{\text{C}}]} \delta_{\beta}^{\text{C}}(i, m) &= \text{bd}_{\text{C}}(m), & \sum_{i \in [2, t_{\text{T}}]} \delta_{\beta}^{\text{T}}(i, m) &= \text{bd}_{\text{T}}(m), \\ \sum_{k \in [1, k_{\text{C}}]} \delta_{\beta}^{\text{CT}}(k, m) &= \text{bd}_{\text{CT}}(m), & \sum_{k \in [1, k_{\text{C}}]} \delta_{\beta}^{\text{TC}}(k, m) &= \text{bd}_{\text{TC}}(m), \\ \sum_{i \in [2, t_{\text{F}}]} \delta_{\beta}^{\text{F}}(i, m) &= \text{bd}_{\text{F}}(m), & \sum_{c \in [1, \widetilde{t}_{\text{C}}]} \delta_{\beta}^{\text{F}}(c, m) &= \text{bd}_{\text{CF}}(m), \\ & & \sum_{c \in [\widetilde{t}_{\text{C}} + 1, c_{\text{F}}]} \delta_{\beta}^{\text{F}}(c, m) &= \text{bd}_{\text{TF}}(m), \end{aligned}$$

$$\text{bd}_{\text{C}}(m) + \text{bd}_{\text{T}}(m) + \text{bd}_{\text{F}}(m) + \text{bd}_{\text{CT}}(m) + \text{bd}_{\text{TC}}(m) + \text{bd}_{\text{TF}}(m) + \text{bd}_{\text{CF}}(m) = \text{bd}^{\text{int}}(m), \quad m \in [1, 3]. \quad (51)$$

D.6 Assigning Chemical Elements and Valence Condition

We include constraints so that each vertex v in a selected graph H satisfies the valence condition; i.e., $\beta_{\text{C}}(v) = \text{val}(\alpha(v)) + \text{eledeg}_{\text{C}}(v)$, where $\text{eledeg}_{\text{C}}(v) = v_{\text{ion}}(\psi)$ for the ρ -fringe-tree $\mathbb{C}[v]$ r-isomorphic to ψ . With these constraints, a chemical graph $\mathbb{C} = (H, \alpha, \beta)$ on a selected subgraph H will be constructed.

constants:

- Subsets $\Lambda^{\text{int}} \subseteq \Lambda \setminus \{\text{H}\}$, $\Lambda^{\text{ex}} \subseteq \Lambda$ of chemical elements, where we denote by $[\mathbf{e}]$ (resp., $[\mathbf{e}]^{\text{int}}$ and $[\mathbf{e}]^{\text{ex}}$) of a standard encoding of an element \mathbf{e} in the set Λ (resp., $\Lambda_{\epsilon}^{\text{int}}$ and $\Lambda_{\epsilon}^{\text{ex}}$);
- A valence function: $\text{val} : \Lambda \rightarrow [1, 6]$;
- Subsets $\Lambda^*(i) \subseteq \Lambda^{\text{int}}$, $i \in [1, t_{\text{C}}]$;

- $\text{na}_{\text{LB}}(\mathbf{a}), \text{na}_{\text{UB}}(\mathbf{a}) \in [0, n^*]$, $\mathbf{a} \in \Lambda$: lower and upper bounds on the number of vertices v with $\alpha(v) = \mathbf{a}$;
- $\text{na}_{\text{LB}}^{\text{int}}(\mathbf{a}), \text{na}_{\text{UB}}^{\text{int}}(\mathbf{a}) \in [0, n^*]$, $\mathbf{a} \in \Lambda^{\text{int}}$: lower and upper bounds on the number of interior-vertices v with $\alpha(v) = \mathbf{a}$;
- $\alpha_r([\psi]) \in [\Lambda^{\text{ex}}]$, $\in \mathcal{F}^*$: the chemical element $\alpha(r)$ of the root r of ψ ;
- $\text{na}_{\mathbf{a}}^{\text{ex}}([\psi]) \in [0, n^*]$, $\mathbf{a} \in \Lambda^{\text{ex}}$, $\psi \in \mathcal{F}^*$: the frequency of chemical element \mathbf{a} in the set of non-rooted vertices in ψ , where possibly $\mathbf{a} = \text{H}$;
- M : an upper bound for the average $\overline{\text{ms}}(\mathbb{C})$ of mass* over all atoms in \mathbb{C} ;

variables:

- $\beta^{\text{CT}}(i), \beta^{\text{TC}}(i) \in [0, 3]$, $i \in [1, t_{\text{T}}]$: the bond-multiplicity of edge $e^{\text{CT}}_{j,i}$ (resp., $e^{\text{TC}}_{j,i}$) if one exists;
- $\beta^{\text{CF}}(i), \beta^{\text{TF}}(i) \in [0, 3]$, $i \in [1, t_{\text{F}}]$: the bond-multiplicity of $e^{\text{CF}}_{j,i}$ (resp., $e^{\text{TF}}_{j,i}$) if one exists;
- $\alpha^{\text{X}}(i) \in [\Lambda_{\epsilon}^{\text{int}}]$, $\delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{int}}) \in [0, 1]$, $\mathbf{a} \in \Lambda_{\epsilon}^{\text{int}}$, $i \in [1, t_{\text{X}}]$, $\text{X} \in \{\text{C}, \text{T}, \text{F}\}$: $\alpha^{\text{X}}(i) = [\mathbf{a}]^{\text{int}} \geq 1$ (resp., $\alpha^{\text{X}}(i) = 0$) $\Leftrightarrow \delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{int}}) = 1$ (resp., $\delta_{\alpha}^{\text{X}}(i, 0) = 0$) $\Leftrightarrow \alpha(v^{\text{X}}_i) = \mathbf{a} \in \Lambda$ (resp., vertex v^{X}_i is not used in \mathbb{C});
- $\delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{int}}) \in [0, 1]$, $i \in [1, t_{\text{X}}]$, $\mathbf{a} \in \Lambda^{\text{int}}$, $\text{X} \in \{\text{C}, \text{T}, \text{F}\}$: $\delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{t}}) = 1 \Leftrightarrow \alpha(v^{\text{X}}_i) = \mathbf{a}$;
- $\text{na}([\mathbf{a}]) \in [\text{na}_{\text{LB}}(\mathbf{a}), \text{na}_{\text{UB}}(\mathbf{a})]$, $\mathbf{a} \in \Lambda$: the number of vertices $v \in V(H)$ with $\alpha(v) = \mathbf{a}$, where possibly $\mathbf{a} = \text{H}$;
- $\text{na}^{\text{int}}([\mathbf{a}]^{\text{int}}) \in [\text{na}_{\text{LB}}^{\text{int}}(\mathbf{a}), \text{na}_{\text{UB}}^{\text{int}}(\mathbf{a})]$, $\mathbf{a} \in \Lambda$, $\text{X} \in \{\text{C}, \text{T}, \text{F}\}$: the number of interior-vertices $v \in V(\mathbb{C})$ with $\alpha(v) = \mathbf{a}$;
- $\text{na}_{\text{X}}^{\text{ex}}([\mathbf{a}]^{\text{ex}}), \text{na}^{\text{ex}}([\mathbf{a}]^{\text{ex}}) \in [0, \text{na}_{\text{UB}}(\mathbf{a})]$, $\mathbf{a} \in \Lambda$, $\text{X} \in \{\text{C}, \text{T}, \text{F}\}$: the number of exterior-vertices rooted at vertices $v \in V_{\text{X}}$ and the number of exterior-vertices v such that $\alpha(v) = \mathbf{a}$;

constraints:

$$\begin{aligned}
\beta^{\text{CT}}(k) - 3(e^{\text{T}}(i) - \chi^{\text{T}}(i, k) + 1) &\leq \beta^{\text{CT}}(i) \leq \beta^{\text{CT}}(k) + 3(e^{\text{T}}(i) - \chi^{\text{T}}(i, k) + 1), i \in [1, t_{\text{T}}], \\
\beta^{\text{TC}}(k) - 3(e^{\text{T}}(i+1) - \chi^{\text{T}}(i, k) + 1) &\leq \beta^{\text{TC}}(i) \leq \beta^{\text{TC}}(k) + 3(e^{\text{T}}(i+1) - \chi^{\text{T}}(i, k) + 1), i \in [1, t_{\text{T}}], \\
&k \in [1, k_{\text{C}}],
\end{aligned} \tag{52}$$

$$\begin{aligned}
\beta^{*\text{F}}(c) - 3(e^{\text{F}}(i) - \chi^{\text{F}}(i, c) + 1) &\leq \beta^{\text{CF}}(i) \leq \beta^{*\text{F}}(c) + 3(e^{\text{F}}(i) - \chi^{\text{F}}(i, c) + 1), i \in [1, t_{\text{F}}], \quad c \in [1, \tilde{t}_{\text{C}}], \\
\beta^{*\text{F}}(c) - 3(e^{\text{F}}(i) - \chi^{\text{F}}(i, c) + 1) &\leq \beta^{\text{TF}}(i) \leq \beta^{*\text{F}}(c) + 3(e^{\text{F}}(i) - \chi^{\text{F}}(i, c) + 1), i \in [1, t_{\text{F}}], \quad c \in [\tilde{t}_{\text{C}} + 1, c_{\text{F}}],
\end{aligned} \tag{53}$$

$$\begin{aligned}
\sum_{\mathbf{a} \in \Lambda^{\text{int}}} \delta_{\alpha}^{\text{C}}(i, [\mathbf{a}]^{\text{int}}) &= 1, & \sum_{\mathbf{a} \in \Lambda^{\text{int}}} [\mathbf{a}]^{\text{int}} \cdot \delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{int}}) &= \alpha^{\text{C}}(i), & i \in [1, t_{\text{C}}], \\
\sum_{\mathbf{a} \in \Lambda^{\text{int}}} \delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{int}}) &= v^{\text{X}}(i), & \sum_{\mathbf{a} \in \Lambda^{\text{int}}} [\mathbf{a}]^{\text{int}} \cdot \delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{int}}) &= \alpha^{\text{X}}(i), & i \in [1, t_{\text{X}}], \text{X} \in \{\text{T}, \text{F}\},
\end{aligned} \tag{54}$$

$$\sum_{\psi \in \mathcal{F}_i^{\text{X}}} \alpha_{\text{r}}([\psi]) \cdot \delta_{\text{fr}}^{\text{X}}(i, [\psi]) = \alpha^{\text{X}}(i), \quad i \in [1, t_{\text{X}}], \text{X} \in \{\text{C}, \text{T}, \text{F}\}, \tag{55}$$

$$\begin{aligned}
\sum_{j \in I_{\text{C}}(i)} \beta^{\text{C}}(j) + \sum_{k \in I_{(\geq 2)}^+(i) \cup I_{(\geq 1)}^+(i)} \beta^{\text{CT}}(k) + \sum_{k \in I_{(\geq 2)}^-(i) \cup I_{(\geq 1)}^-(i)} \beta^{\text{TC}}(k) \\
+ \beta^{*\text{F}}(i) + \beta_{\text{ex}}^{\text{C}}(i) - \text{eledeg}_{\text{C}}(i) &= \sum_{\mathbf{a} \in \Lambda^{\text{int}}} \text{val}(\mathbf{a}) \delta_{\alpha}^{\text{C}}(i, [\mathbf{a}]^{\text{int}}), & i \in [1, \tilde{t}_{\text{C}}],
\end{aligned} \tag{56}$$

$$\begin{aligned}
\sum_{j \in I_{\text{C}}(i)} \beta^{\text{C}}(j) + \sum_{k \in I_{(\geq 2)}^+(i) \cup I_{(\geq 1)}^+(i)} \beta^{\text{CT}}(k) + \sum_{k \in I_{(\geq 2)}^-(i) \cup I_{(\geq 1)}^-(i)} \beta^{\text{TC}}(k) \\
+ \beta_{\text{ex}}^{\text{C}}(i) - \text{eledeg}_{\text{C}}(i) &= \sum_{\mathbf{a} \in \Lambda^{\text{int}}} \text{val}(\mathbf{a}) \delta_{\alpha}^{\text{C}}(i, [\mathbf{a}]^{\text{int}}), & i \in [\tilde{t}_{\text{C}} + 1, t_{\text{C}}],
\end{aligned} \tag{57}$$

$$\begin{aligned}
\beta^{\text{T}}(i) + \beta^{\text{T}}(i+1) + \beta_{\text{ex}}^{\text{T}}(i) + \beta^{\text{CT}}(i) + \beta^{\text{TC}}(i) \\
+ \beta^{*\text{F}}(\tilde{t}_{\text{C}} + i) - \text{eledeg}_{\text{T}}(i) &= \sum_{\mathbf{a} \in \Lambda^{\text{int}}} \text{val}(\mathbf{a}) \delta_{\alpha}^{\text{T}}(i, [\mathbf{a}]^{\text{int}}), \\
i \in [1, t_{\text{T}}] \quad (\beta^{\text{T}}(1) = \beta^{\text{T}}(t_{\text{T}} + 1) = 0), &
\end{aligned} \tag{58}$$

$$\begin{aligned}
\beta^{\text{F}}(i) + \beta^{\text{F}}(i+1) + \beta^{\text{CF}}(i) + \beta^{\text{TF}}(i) \\
+ \beta_{\text{ex}}^{\text{F}}(i) - \text{eledeg}_{\text{F}}(i) &= \sum_{\mathbf{a} \in \Lambda^{\text{int}}} \text{val}(\mathbf{a}) \delta_{\alpha}^{\text{F}}(i, [\mathbf{a}]^{\text{int}}), \\
i \in [1, t_{\text{F}}] \quad (\beta^{\text{F}}(1) = \beta^{\text{F}}(t_{\text{F}} + 1) = 0), &
\end{aligned} \tag{59}$$

$$\sum_{i \in [1, t_{\text{X}}]} \delta_{\alpha}^{\text{X}}(i, [\mathbf{a}]^{\text{int}}) = \text{na}_{\text{X}}([\mathbf{a}]^{\text{int}}), \quad \mathbf{a} \in \Lambda^{\text{int}}, \text{X} \in \{\text{C}, \text{T}, \text{F}\}, \tag{60}$$

$$\sum_{\psi \in \mathcal{F}_i^{\text{X}}, i \in [1, t_{\text{X}}]} \text{na}_{\mathbf{a}}^{\text{ex}}([\psi]) \cdot \delta_{\text{fr}}^{\text{X}}(i, [\psi]) = \text{na}_{\text{X}}^{\text{ex}}([\mathbf{a}]^{\text{ex}}), \quad \mathbf{a} \in \Lambda^{\text{ex}}, \text{X} \in \{\text{C}, \text{T}, \text{F}\}, \tag{61}$$

$$\begin{aligned}
na_C([\mathbf{a}]^{\text{int}}) + na_T([\mathbf{a}]^{\text{int}}) + na_F([\mathbf{a}]^{\text{int}}) &= na^{\text{int}}([\mathbf{a}]^{\text{int}}), & \mathbf{a} \in \Lambda^{\text{int}}, \\
\sum_{X \in \{C, T, F\}} na_X^{\text{ex}}([\mathbf{a}]^{\text{ex}}) &= na^{\text{ex}}([\mathbf{a}]^{\text{ex}}), & \mathbf{a} \in \Lambda^{\text{ex}}, \\
na^{\text{int}}([\mathbf{a}]^{\text{int}}) + na^{\text{ex}}([\mathbf{a}]^{\text{ex}}) &= na([\mathbf{a}]), & \mathbf{a} \in \Lambda^{\text{int}} \cap \Lambda^{\text{ex}}, \\
na^{\text{int}}([\mathbf{a}]^{\text{int}}) &= na([\mathbf{a}]), & \mathbf{a} \in \Lambda^{\text{int}} \setminus \Lambda^{\text{ex}}, \\
na^{\text{ex}}([\mathbf{a}]^{\text{ex}}) &= na([\mathbf{a}]), & \mathbf{a} \in \Lambda^{\text{ex}} \setminus \Lambda^{\text{int}},
\end{aligned} \tag{62}$$

$$\sum_{\mathbf{a} \in \Lambda^*(i)} \delta_\alpha^C(i, [\mathbf{a}]^{\text{int}}) = 1, \quad i \in [1, t_C]. \tag{63}$$

D.7 Constraints for Bounds on the Number of Bonds

We include constraints for specification of lower and upper bounds bd_{LB} and bd_{UB} .

constants:

- $bd_{m, \text{LB}}(i), bd_{m, \text{UB}}(i) \in [0, n_{\text{UB}}^{\text{int}}], i \in [1, m_C], m \in [2, 3]$: lower and upper bounds on the number of edges $e \in E(P_i)$ with bond-multiplicity $\beta(e) = m$ in the pure path P_i for edge $e_i \in E_C$;

variables :

- $bd_T(k, i, m) \in [0, 1], k \in [1, k_C], i \in [2, t_T], m \in [2, 3]$: $bd_T(k, i, m) = 1 \Leftrightarrow$ the pure path P_k for edge $e_k \in E_C$ contains edge e_i^T with $\beta(e_i^T) = m$;

constraints:

$$bd_{m, \text{LB}}(i) \leq \delta_\beta^C(i, m) \leq bd_{m, \text{UB}}(i), i \in I_{(=1)} \cup I_{(0/1)}, m \in [2, 3], \tag{64}$$

$$bd_T(k, i, m) \geq \delta_\beta^T(i, m) + \chi^T(i, k) - 1, \quad k \in [1, k_C], i \in [2, t_T], m \in [2, 3], \tag{65}$$

$$\sum_{j \in [2, t_T]} \delta_\beta^T(j, m) \geq \sum_{k \in [1, k_C], i \in [2, t_T]} bd_T(k, i, m), \quad m \in [2, 3], \tag{66}$$

$$bd_{m, \text{LB}}(k) \leq \sum_{i \in [2, t_T]} bd_T(k, i, m) + \delta_\beta^{\text{CT}}(k, m) + \delta_\beta^{\text{TC}}(k, m) \leq bd_{m, \text{UB}}(k), \tag{67}$$

$$k \in [1, k_C], m \in [2, 3].$$

D.8 Descriptor for the Number of Adjacency-configurations

We call a tuple $(\mathbf{a}, \mathbf{b}, m) \in (\Lambda \setminus \{\mathbf{H}\}) \times (\Lambda \setminus \{\mathbf{H}\}) \times [1, 3]$ an *adjacency-configuration*. The adjacency-configuration of an edge-configuration $(\mu = \mathbf{ad}, \mu' = \mathbf{bd}', m)$ is defined to be $(\mathbf{a}, \mathbf{b}, m)$. We include constraints to compute the frequency of each adjacency-configuration in an inferred chemical graph \mathbb{C} .

constants:

- A set Γ^{int} of edge-configurations $\gamma = (\mu, \mu', m)$ with $\mu \leq \mu'$;
- Let $\bar{\gamma}$ of an edge-configuration $\gamma = (\mu, \mu', m)$ denote the edge-configuration (μ', μ, m) ;
- Let $\Gamma_{<}^{\text{int}} = \{(\mu, \mu', m) \in \Gamma^{\text{int}} \mid \mu < \mu'\}$, $\Gamma_{=}^{\text{int}} = \{(\mu, \mu', m) \in \Gamma^{\text{int}} \mid \mu = \mu'\}$ and $\Gamma_{>}^{\text{int}} = \{\bar{\gamma} \mid \gamma \in \Gamma_{<}^{\text{int}}\}$;
- Let $\Gamma_{\text{ac}, <}^{\text{int}}$, $\Gamma_{\text{ac}, =}^{\text{int}}$ and $\Gamma_{\text{ac}, >}^{\text{int}}$ denote the sets of the adjacency-configurations of edge-configurations in the sets $\Gamma_{<}^{\text{int}}$, $\Gamma_{=}^{\text{int}}$ and $\Gamma_{>}^{\text{int}}$, respectively;
- Let $\bar{\nu}$ of an adjacency-configuration $\nu = (\mathbf{a}, \mathbf{b}, m)$ denote the adjacency-configuration $(\mathbf{b}, \mathbf{a}, m)$;
- Prepare a coding of the set $\Gamma_{\text{ac}}^{\text{int}} \cup \Gamma_{\text{ac}, >}^{\text{int}}$ and let $[\nu]^{\text{int}}$ denote the coded integer of an element ν in $\Gamma_{\text{ac}}^{\text{int}} \cup \Gamma_{\text{ac}, >}^{\text{int}}$;
- Choose subsets $\tilde{\Gamma}_{\text{ac}}^{\text{C}}, \tilde{\Gamma}_{\text{ac}}^{\text{T}}, \tilde{\Gamma}_{\text{ac}}^{\text{CT}}, \tilde{\Gamma}_{\text{ac}}^{\text{TC}}, \tilde{\Gamma}_{\text{ac}}^{\text{F}}, \tilde{\Gamma}_{\text{ac}}^{\text{CF}}, \tilde{\Gamma}_{\text{ac}}^{\text{TF}} \subseteq \Gamma_{\text{ac}}^{\text{int}} \cup \Gamma_{\text{ac}, >}^{\text{int}}$; To compute the frequency of adjacency-configurations exactly, set $\tilde{\Gamma}_{\text{ac}}^{\text{C}} := \tilde{\Gamma}_{\text{ac}}^{\text{T}} := \tilde{\Gamma}_{\text{ac}}^{\text{CT}} := \tilde{\Gamma}_{\text{ac}}^{\text{TC}} := \tilde{\Gamma}_{\text{ac}}^{\text{F}} := \tilde{\Gamma}_{\text{ac}}^{\text{CF}} := \tilde{\Gamma}_{\text{ac}}^{\text{TF}} := \Gamma_{\text{ac}}^{\text{int}} \cup \Gamma_{\text{ac}, >}^{\text{int}}$;
- $\text{ac}_{\text{LB}}^{\text{int}}(\nu), \text{ac}_{\text{UB}}^{\text{int}}(\nu) \in [0, 2n_{\text{UB}}^{\text{int}}], \nu = (\mathbf{a}, \mathbf{b}, m) \in \Gamma_{\text{ac}}^{\text{int}}$: lower and upper bounds on the number of interior-edges $e = uv$ with $\alpha(u) = \mathbf{a}$, $\alpha(v) = \mathbf{b}$ and $\beta(e) = m$;

variables:

- $\text{ac}^{\text{int}}([\nu]^{\text{int}}) \in [\text{ac}_{\text{LB}}^{\text{int}}(\nu), \text{ac}_{\text{UB}}^{\text{int}}(\nu)], \nu \in \Gamma_{\text{ac}}^{\text{int}}$: the number of interior-edges with adjacency-configuration ν ;
- $\text{ac}_{\text{C}}([\nu]^{\text{int}}) \in [0, m_{\text{C}}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{C}}, \text{ac}_{\text{T}}([\nu]^{\text{int}}) \in [0, t_{\text{T}}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{T}}, \text{ac}_{\text{F}}([\nu]^{\text{int}}) \in [0, t_{\text{F}}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{F}}$: the number of edges $e^{\text{C}} \in E_{\text{C}}$ (resp., edges $e^{\text{T}} \in E_{\text{T}}$ and edges $e^{\text{F}} \in E_{\text{F}}$) with adjacency-configuration ν ;
- $\text{ac}_{\text{CT}}([\nu]^{\text{int}}) \in [0, \min\{k_{\text{C}}, t_{\text{T}}\}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{CT}}, \text{ac}_{\text{TC}}([\nu]^{\text{int}}) \in [0, \min\{k_{\text{C}}, t_{\text{T}}\}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{TC}}, \text{ac}_{\text{CF}}([\nu]^{\text{int}}) \in [0, \tilde{t}_{\text{C}}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{CF}}, \text{ac}_{\text{TF}}([\nu]^{\text{int}}) \in [0, t_{\text{T}}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{TF}}$: the number of edges $e^{\text{CT}} \in E_{\text{CT}}$ (resp., edges $e^{\text{TC}} \in E_{\text{TC}}$ and edges $e^{\text{CF}} \in E_{\text{CF}}$ and $e^{\text{TF}} \in E_{\text{TF}}$) with adjacency-configuration ν ;
- $\delta_{\text{ac}}^{\text{C}}(i, [\nu]^{\text{int}}) \in [0, 1], i \in [\tilde{k}_{\text{C}} + 1, m_{\text{C}}] = I_{(\geq 1)} \cup I_{(0/1)} \cup I_{(=1)}, \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{C}}, \delta_{\text{ac}}^{\text{T}}(i, [\nu]^{\text{int}}) \in [0, 1], i \in [2, t_{\text{T}}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{T}}, \delta_{\text{ac}}^{\text{F}}(i, [\nu]^{\text{int}}) \in [0, 1], i \in [2, t_{\text{F}}], \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{F}}$: $\delta_{\text{ac}}^{\text{X}}(i, [\nu]^{\text{int}}) = 1 \Leftrightarrow$ edge e^{X}_i has adjacency-configuration ν ;

- $\delta_{ac}^{CT}(k, [\nu]^{int}), \delta_{ac}^{TC}(k, [\nu]^{int}) \in [0, 1], k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}, \nu \in \tilde{\Gamma}_{ac}^{CT}$: $\delta_{ac}^{CT}(k, [\nu]^{int}) = 1$ (resp., $\delta_{ac}^{TC}(k, [\nu]^{int}) = 1$) \Leftrightarrow edge $e^{CT}_{tail(k),j}$ (resp., $e^{TC}_{head(k),j}$) for some $j \in [1, t_T]$ has adjacency-configuration ν ;
- $\delta_{ac}^{CF}(c, [\nu]^{int}) \in [0, 1], c \in [1, \tilde{t}_C], \nu \in \tilde{\Gamma}_{ac}^{CF}$: $\delta_{ac}^{CF}(c, [\nu]^{int}) = 1 \Leftrightarrow$ edge $e^{CF}_{c,i}$ for some $i \in [1, t_F]$ has adjacency-configuration ν ;
- $\delta_{ac}^{TF}(i, [\nu]^{int}) \in [0, 1], i \in [1, t_T], \nu \in \tilde{\Gamma}_{ac}^{TF}$: $\delta_{ac}^{TF}(i, [\nu]^{int}) = 1 \Leftrightarrow$ edge $e^{TF}_{i,j}$ for some $j \in [1, t_F]$ has adjacency-configuration ν ;
- $\alpha^{CT}(k), \alpha^{TC}(k) \in [0, |\Lambda^{int}|], k \in [1, k_C]$: $\alpha(v)$ of the edge $(v^C_{tail(k)}, v) \in E_{CT}$ (resp., $(v, v^C_{head(k)}) \in E_{TC}$) if any;
- $\alpha^{CF}(c) \in [0, |\Lambda^{int}|], c \in [1, \tilde{t}_C]$: $\alpha(v)$ of the edge $(v^C_c, v) \in E_{CF}$ if any;
- $\alpha^{TF}(i) \in [0, |\Lambda^{int}|], i \in [1, t_T]$: $\alpha(v)$ of the edge $(v^T_i, v) \in E_{TF}$ if any;
- $\Delta_{ac}^{C+}(i), \Delta_{ac}^{C-}(i) \in [0, |\Lambda^{int}|], i \in [\tilde{k}_C+1, m_C], \Delta_{ac}^{T+}(i), \Delta_{ac}^{T-}(i) \in [0, |\Lambda^{int}|], i \in [2, t_T], \Delta_{ac}^{F+}(i), \Delta_{ac}^{F-}(i) \in [0, |\Lambda^{int}|], i \in [2, t_F]$: $\Delta_{ac}^{X+}(i) = \Delta_{ac}^{X-}(i) = 0$ (resp., $\Delta_{ac}^{X+}(i) = \alpha(u)$ and $\Delta_{ac}^{X-}(i) = \alpha(v)$) \Leftrightarrow edge $e^X_i = (u, v) \in E_X$ is used in \mathbb{C} (resp., $e^X_i \notin E(G)$);
- $\Delta_{ac}^{CT+}(k), \Delta_{ac}^{CT-}(k) \in [0, |\Lambda^{int}|], k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$: $\Delta_{ac}^{CT+}(k) = \Delta_{ac}^{CT-}(k) = 0$ (resp., $\Delta_{ac}^{CT+}(k) = \alpha(u)$ and $\Delta_{ac}^{CT-}(k) = \alpha(v)$) \Leftrightarrow edge $e^{CT}_{tail(k),j} = (u, v) \in E_{CT}$ for some $j \in [1, t_T]$ is used in \mathbb{C} (resp., otherwise);
- $\Delta_{ac}^{TC+}(k), \Delta_{ac}^{TC-}(k) \in [0, |\Lambda^{int}|], k \in [1, k_C] = I_{(\geq 2)} \cup I_{(\geq 1)}$: Analogous with $\Delta_{ac}^{CT+}(k)$ and $\Delta_{ac}^{CT-}(k)$;
- $\Delta_{ac}^{CF+}(c) \in [0, |\Lambda^{int}|], \Delta_{ac}^{CF-}(c) \in [0, |\Lambda^{int}|], c \in [1, \tilde{t}_C]$: $\Delta_{ac}^{CF+}(c) = \Delta_{ac}^{CF-}(c) = 0$ (resp., $\Delta_{ac}^{CF+}(c) = \alpha(u)$ and $\Delta_{ac}^{CF-}(c) = \alpha(v)$) \Leftrightarrow edge $e^{CF}_{c,i} = (u, v) \in E_{CF}$ for some $i \in [1, t_F]$ is used in \mathbb{C} (resp., otherwise);
- $\Delta_{ac}^{TF+}(i) \in [0, |\Lambda^{int}|], \Delta_{ac}^{TF-}(i) \in [0, |\Lambda^{int}|], i \in [1, t_T]$: Analogous with $\Delta_{ac}^{CF+}(c)$ and $\Delta_{ac}^{CF-}(c)$;

constraints:

$$\begin{aligned}
ac_C([\nu]^{int}) &= 0, & \nu &\in \Gamma_{ac}^{int} \setminus \tilde{\Gamma}_{ac}^C, \\
ac_T([\nu]^{int}) &= 0, & \nu &\in \Gamma_{ac}^{int} \setminus \tilde{\Gamma}_{ac}^T, \\
ac_F([\nu]^{int}) &= 0, & \nu &\in \Gamma_{ac}^{int} \setminus \tilde{\Gamma}_{ac}^F, \\
ac_{CT}([\nu]^{int}) &= 0, & \nu &\in \Gamma_{ac}^{int} \setminus \tilde{\Gamma}_{ac}^{CT}, \\
ac_{TC}([\nu]^{int}) &= 0, & \nu &\in \Gamma_{ac}^{int} \setminus \tilde{\Gamma}_{ac}^{TC}, \\
ac_{CF}([\nu]^{int}) &= 0, & \nu &\in \Gamma_{ac}^{int} \setminus \tilde{\Gamma}_{ac}^{CF}, \\
ac_{TF}([\nu]^{int}) &= 0, & \nu &\in \Gamma_{ac}^{int} \setminus \tilde{\Gamma}_{ac}^{TF},
\end{aligned}$$

(68)

$$\begin{aligned}
\sum_{(\mathbf{a}, \mathbf{b}, m) = \nu \in \Gamma_{\text{ac}}^{\text{int}}} \text{ac}_C([\nu]^{\text{int}}) &= \sum_{i \in [\widetilde{k}_C + 1, m_C]} \delta_\beta^C(i, m), & m \in [1, 3], \\
\sum_{(\mathbf{a}, \mathbf{b}, m) = \nu \in \Gamma_{\text{ac}}^{\text{int}}} \text{ac}_T([\nu]^{\text{int}}) &= \sum_{i \in [2, t_T]} \delta_\beta^T(i, m), & m \in [1, 3], \\
\sum_{(\mathbf{a}, \mathbf{b}, m) = \nu \in \Gamma_{\text{ac}}^{\text{int}}} \text{ac}_F([\nu]^{\text{int}}) &= \sum_{i \in [2, t_F]} \delta_\beta^F(i, m), & m \in [1, 3], \\
\sum_{(\mathbf{a}, \mathbf{b}, m) = \nu \in \Gamma_{\text{ac}}^{\text{int}}} \text{ac}_{\text{CT}}([\nu]^{\text{int}}) &= \sum_{k \in [1, k_C]} \delta_\beta^{\text{CT}}(k, m), & m \in [1, 3], \\
\sum_{(\mathbf{a}, \mathbf{b}, m) = \nu \in \Gamma_{\text{ac}}^{\text{int}}} \text{ac}_{\text{TC}}([\nu]^{\text{int}}) &= \sum_{k \in [1, k_C]} \delta_\beta^{\text{TC}}(k, m), & m \in [1, 3], \\
\sum_{(\mathbf{a}, \mathbf{b}, m) = \nu \in \Gamma_{\text{ac}}^{\text{int}}} \text{ac}_{\text{CF}}([\nu]^{\text{int}}) &= \sum_{c \in [1, \widetilde{t}_C]} \delta_\beta^{*\text{F}}(c, m), & m \in [1, 3], \\
\sum_{(\mathbf{a}, \mathbf{b}, m) = \nu \in \Gamma_{\text{ac}}^{\text{int}}} \text{ac}_{\text{TF}}([\nu]^{\text{int}}) &= \sum_{c \in [\widetilde{t}_C + 1, c_F]} \delta_\beta^{*\text{F}}(c, m), & m \in [1, 3],
\end{aligned} \tag{69}$$

$$\begin{aligned}
\sum_{\nu = (\mathbf{a}, \mathbf{b}, m) \in \widetilde{\Gamma}_{\text{ac}}^{\text{C}}} m \cdot \delta_{\text{ac}}^{\text{C}}(i, [\nu]^{\text{int}}) &= \beta^{\text{C}}(i), \\
\Delta_{\text{ac}}^{\text{C}+}(i) + \sum_{\nu = (\mathbf{a}, \mathbf{b}, m) \in \widetilde{\Gamma}_{\text{ac}}^{\text{C}}} [\mathbf{a}]^{\text{int}} \delta_{\text{ac}}^{\text{C}}(i, [\nu]^{\text{int}}) &= \alpha^{\text{C}}(\text{tail}(i)), \\
\Delta_{\text{ac}}^{\text{C}-}(i) + \sum_{\nu = (\mathbf{a}, \mathbf{b}, m) \in \widetilde{\Gamma}_{\text{ac}}^{\text{C}}} [\mathbf{b}]^{\text{int}} \delta_{\text{ac}}^{\text{C}}(i, [\nu]^{\text{int}}) &= \alpha^{\text{C}}(\text{head}(i)), \\
\Delta_{\text{ac}}^{\text{C}+}(i) + \Delta_{\text{ac}}^{\text{C}-}(i) &\leq 2|\Lambda^{\text{int}}|(1 - e^{\text{C}}(i)), & i \in [\widetilde{k}_C + 1, m_C], \\
\sum_{i \in [\widetilde{k}_C + 1, m_C]} \delta_{\text{ac}}^{\text{C}}(i, [\nu]^{\text{int}}) &= \text{ac}_C([\nu]^{\text{int}}), & \nu \in \widetilde{\Gamma}_{\text{ac}}^{\text{C}},
\end{aligned} \tag{70}$$

$$\begin{aligned}
\sum_{\nu = (\mathbf{a}, \mathbf{b}, m) \in \widetilde{\Gamma}_{\text{ac}}^{\text{T}}} m \cdot \delta_{\text{ac}}^{\text{T}}(i, [\nu]^{\text{int}}) &= \beta^{\text{T}}(i), \\
\Delta_{\text{ac}}^{\text{T}+}(i) + \sum_{\nu = (\mathbf{a}, \mathbf{b}, m) \in \widetilde{\Gamma}_{\text{ac}}^{\text{T}}} [\mathbf{a}]^{\text{int}} \delta_{\text{ac}}^{\text{T}}(i, [\nu]^{\text{int}}) &= \alpha^{\text{T}}(i - 1), \\
\Delta_{\text{ac}}^{\text{T}-}(i) + \sum_{\nu = (\mathbf{a}, \mathbf{b}, m) \in \widetilde{\Gamma}_{\text{ac}}^{\text{T}}} [\mathbf{b}]^{\text{int}} \delta_{\text{ac}}^{\text{T}}(i, [\nu]^{\text{int}}) &= \alpha^{\text{T}}(i), \\
\Delta_{\text{ac}}^{\text{T}+}(i) + \Delta_{\text{ac}}^{\text{T}-}(i) &\leq 2|\Lambda^{\text{int}}|(1 - e^{\text{T}}(i)), & i \in [2, t_T], \\
\sum_{i \in [2, t_T]} \delta_{\text{ac}}^{\text{T}}(i, [\nu]^{\text{int}}) &= \text{ac}_T([\nu]^{\text{int}}), & \nu \in \widetilde{\Gamma}_{\text{ac}}^{\text{T}},
\end{aligned} \tag{71}$$

$$\begin{aligned}
& \sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{F}}} m \cdot \delta_{\text{ac}}^{\text{F}}(i, [\nu]^{\text{int}}) = \beta^{\text{F}}(i), \\
\Delta_{\text{ac}}^{\text{F}+}(i) + & \sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{F}}} [\mathbf{a}]^{\text{int}} \delta_{\text{ac}}^{\text{F}}(i, [\nu]^{\text{int}}) = \alpha^{\text{F}}(i-1), \\
\Delta_{\text{ac}}^{\text{F}-}(i) + & \sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{F}}} [\mathbf{b}]^{\text{int}} \delta_{\text{ac}}^{\text{F}}(i, [\nu]^{\text{int}}) = \alpha^{\text{F}}(i), \\
\Delta_{\text{ac}}^{\text{F}+}(i) + \Delta_{\text{ac}}^{\text{F}-}(i) & \leq 2|\Lambda^{\text{ex}}|(1 - e^{\text{F}}(i)), \\
\sum_{i\in[2,t_{\text{F}}]} \delta_{\text{ac}}^{\text{F}}(i, [\nu]^{\text{int}}) & = \text{ac}_{\text{F}}([\nu]^{\text{int}}),
\end{aligned} \tag{72}$$

$i \in [2, t_{\text{F}}],$
 $\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{F}},$

$$\begin{aligned}
\alpha^{\text{T}}(i) + |\Lambda^{\text{int}}|(1 - \chi^{\text{T}}(i, k) + e^{\text{T}}(i)) & \geq \alpha^{\text{CT}}(k), \\
\alpha^{\text{CT}}(k) \geq \alpha^{\text{T}}(i) - |\Lambda^{\text{int}}|(1 - \chi^{\text{T}}(i, k) + e^{\text{T}}(i)), \\
\sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{CT}}} m \cdot \delta_{\text{ac}}^{\text{CT}}(k, [\nu]^{\text{int}}) & = \beta^{\text{CT}}(k), \\
\Delta_{\text{ac}}^{\text{CT}+}(k) + \sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{CT}}} [\mathbf{a}]^{\text{int}} \delta_{\text{ac}}^{\text{CT}}(k, [\nu]^{\text{int}}) & = \alpha^{\text{C}}(\text{tail}(k)), \\
\Delta_{\text{ac}}^{\text{CT}-}(k) + \sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{CT}}} [\mathbf{b}]^{\text{int}} \delta_{\text{ac}}^{\text{CT}}(k, [\nu]^{\text{int}}) & = \alpha^{\text{CT}}(k), \\
\Delta_{\text{ac}}^{\text{CT}+}(k) + \Delta_{\text{ac}}^{\text{CT}-}(k) & \leq 2|\Lambda^{\text{int}}|(1 - \delta_{\chi}^{\text{T}}(k)), \\
\sum_{k\in[1,k_{\text{C}}]} \delta_{\text{ac}}^{\text{CT}}(k, [\nu]^{\text{int}}) & = \text{ac}_{\text{CT}}([\nu]^{\text{int}}),
\end{aligned} \tag{73}$$

$i \in [1, t_{\text{T}}],$
 $k \in [1, k_{\text{C}}],$
 $\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{CT}},$

$$\begin{aligned}
\alpha^{\text{T}}(i) + |\Lambda^{\text{int}}|(1 - \chi^{\text{T}}(i, k) + e^{\text{T}}(i+1)) & \geq \alpha^{\text{TC}}(k), \\
\alpha^{\text{TC}}(k) \geq \alpha^{\text{T}}(i) - |\Lambda^{\text{int}}|(1 - \chi^{\text{T}}(i, k) + e^{\text{T}}(i+1)), \\
\sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{TC}}} m \cdot \delta_{\text{ac}}^{\text{TC}}(k, [\nu]^{\text{int}}) & = \beta^{\text{TC}}(k), \\
\Delta_{\text{ac}}^{\text{TC}+}(k) + \sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{TC}}} [\mathbf{a}]^{\text{int}} \delta_{\text{ac}}^{\text{TC}}(k, [\nu]^{\text{int}}) & = \alpha^{\text{TC}}(k), \\
\Delta_{\text{ac}}^{\text{TC}-}(k) + \sum_{\nu=(\mathbf{a},\mathbf{b},m)\in\tilde{\Gamma}_{\text{ac}}^{\text{TC}}} [\mathbf{b}]^{\text{int}} \delta_{\text{ac}}^{\text{TC}}(k, [\nu]^{\text{int}}) & = \alpha^{\text{C}}(\text{head}(k)), \\
\Delta_{\text{ac}}^{\text{TC}+}(k) + \Delta_{\text{ac}}^{\text{TC}-}(k) & \leq 2|\Lambda^{\text{int}}|(1 - \delta_{\chi}^{\text{T}}(k)), \\
\sum_{k\in[1,k_{\text{C}}]} \delta_{\text{ac}}^{\text{TC}}(k, [\nu]^{\text{int}}) & = \text{ac}_{\text{TC}}([\nu]^{\text{int}}),
\end{aligned} \tag{74}$$

$i \in [1, t_{\text{T}}],$
 $k \in [1, k_{\text{C}}],$
 $\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{TC}},$

$$\begin{aligned}
\alpha^F(i) + |\Lambda^{\text{int}}|(1 - \chi^F(i, c) + e^F(i)) &\geq \alpha^{\text{CF}}(c), \\
\alpha^{\text{CF}}(c) &\geq \alpha^F(i) - |\Lambda^{\text{int}}|(1 - \chi^F(i, c) + e^F(i)), & i \in [1, t_F], \\
\sum_{\nu=(\mathbf{a}, \mathbf{b}, m) \in \tilde{\Gamma}_{\text{ac}}^{\text{CF}}} m \cdot \delta_{\text{ac}}^{\text{CF}}(c, [\nu]^{\text{int}}) &= \beta^{*F}(c), \\
\Delta_{\text{ac}}^{\text{CF}+}(c) + \sum_{\nu=(\mathbf{a}, \mathbf{b}, m) \in \tilde{\Gamma}_{\text{ac}}^{\text{CF}}} [\mathbf{a}]^{\text{int}} \delta_{\text{ac}}^{\text{CF}}(c, [\nu]^{\text{int}}) &= \alpha^{\text{C}}(\text{head}(c)), \\
\Delta_{\text{ac}}^{\text{CF}-}(c) + \sum_{\nu=(\mathbf{a}, \mathbf{b}, m) \in \tilde{\Gamma}_{\text{ac}}^{\text{CF}}} [\mathbf{b}]^{\text{int}} \delta_{\text{ac}}^{\text{CF}}(c, [\nu]^{\text{int}}) &= \alpha^{\text{CF}}(c), \\
\Delta_{\text{ac}}^{\text{CF}+}(c) + \Delta_{\text{ac}}^{\text{CF}-}(c) &\leq 2 \max\{|\Lambda^{\text{int}}|, |\Lambda^{\text{int}}|\}(1 - \delta_{\chi}^F(c)), & c \in [1, \tilde{t}_{\text{C}}], \\
\sum_{c \in [1, \tilde{t}_{\text{C}}]} \delta_{\text{ac}}^{\text{CF}}(c, [\nu]^{\text{int}}) &= \text{ac}_{\text{CF}}([\nu]^{\text{int}}), & \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{CF}}, \quad (75)
\end{aligned}$$

$$\begin{aligned}
\alpha^F(j) + |\Lambda^{\text{int}}|(1 - \chi^F(j, i + \tilde{t}_{\text{C}}) + e^F(j)) &\geq \alpha^{\text{TF}}(i), \\
\alpha^{\text{TF}}(i) &\geq \alpha^F(j) - |\Lambda^{\text{int}}|(1 - \chi^F(j, i + \tilde{t}_{\text{C}}) + e^F(j)), & j \in [1, t_F], \\
\sum_{\nu=(\mathbf{a}, \mathbf{b}, m) \in \tilde{\Gamma}_{\text{ac}}^{\text{TF}}} m \cdot \delta_{\text{ac}}^{\text{TF}}(i, [\nu]^{\text{int}}) &= \beta^{*F}(i + \tilde{t}_{\text{C}}), \\
\Delta_{\text{ac}}^{\text{TF}+}(i) + \sum_{\nu=(\mathbf{a}, \mathbf{b}, m) \in \tilde{\Gamma}_{\text{ac}}^{\text{TF}}} [\mathbf{a}]^{\text{int}} \delta_{\text{ac}}^{\text{TF}}(i, [\nu]^{\text{int}}) &= \alpha^{\text{T}}(i), \\
\Delta_{\text{ac}}^{\text{TF}-}(i) + \sum_{\nu=(\mathbf{a}, \mathbf{b}, m) \in \tilde{\Gamma}_{\text{ac}}^{\text{TF}}} [\mathbf{b}]^{\text{int}} \delta_{\text{ac}}^{\text{TF}}(i, [\nu]^{\text{int}}) &= \alpha^{\text{TF}}(i), \\
\Delta_{\text{ac}}^{\text{TF}+}(i) + \Delta_{\text{ac}}^{\text{TF}-}(i) &\leq 2 \max\{|\Lambda^{\text{int}}|, |\Lambda^{\text{int}}|\}(1 - \delta_{\chi}^F(i + \tilde{t}_{\text{C}})), & i \in [1, t_T], \\
\sum_{i \in [1, t_T]} \delta_{\text{ac}}^{\text{TF}}(i, [\nu]^{\text{int}}) &= \text{ac}_{\text{TF}}([\nu]^{\text{int}}), & \nu \in \tilde{\Gamma}_{\text{ac}}^{\text{TF}}, \quad (76)
\end{aligned}$$

$$\begin{aligned}
\sum_{\mathbf{X} \in \{\text{C}, \text{T}, \text{F}, \text{CT}, \text{TC}, \text{CF}, \text{TF}\}} (\text{ac}_{\mathbf{X}}([\nu]^{\text{int}}) + \text{ac}_{\mathbf{X}}([\bar{\nu}]^{\text{int}})) &= \text{ac}^{\text{int}}([\nu]^{\text{int}}), & \nu \in \Gamma_{\text{ac}, <}^{\text{int}}, \\
\sum_{\mathbf{X} \in \{\text{C}, \text{T}, \text{F}, \text{CT}, \text{TC}, \text{CF}, \text{TF}\}} \text{ac}_{\mathbf{X}}([\nu]^{\text{int}}) &= \text{ac}^{\text{int}}([\nu]^{\text{int}}), & \nu \in \Gamma_{\text{ac}, =}^{\text{int}}. \quad (77)
\end{aligned}$$

D.9 Descriptor for the Number of Chemical Symbols

We include constraints for computing the frequency of each chemical symbol in Λ_{dg} . Let $\text{cs}(v)$ denote the chemical symbol of an interior-vertex v in a chemical graph \mathbb{C} to be inferred; i.e., $\text{cs}(v) = \mu = \mathbf{ad} \in \Lambda_{\text{dg}}$ such that $\alpha(v) = \mathbf{a}$ and $\text{deg}_{\langle \mathbb{C} \rangle}(v) = \text{deg}_H(v) - \text{deg}_{\mathbb{C}}^{\text{hyd}}(v) = d$ in $\mathbb{C} = (H, \alpha, \beta)$.

constants:

- A set $\Lambda_{\text{dg}}^{\text{int}}$ of chemical symbols;

- Prepare a coding of each of the two sets $\Lambda_{\text{dg}}^{\text{int}}$ and let $[\mu]^{\text{int}}$ denote the coded integer of an element $\mu \in \Lambda_{\text{dg}}^{\text{int}}$;
- Choose subsets $\tilde{\Lambda}_{\text{dg}}^{\text{C}}, \tilde{\Lambda}_{\text{dg}}^{\text{T}}, \tilde{\Lambda}_{\text{dg}}^{\text{F}} \subseteq \Lambda_{\text{dg}}^{\text{int}}$: To compute the frequency of chemical symbols exactly, set $\tilde{\Lambda}_{\text{dg}}^{\text{C}} := \tilde{\Lambda}_{\text{dg}}^{\text{T}} := \tilde{\Lambda}_{\text{dg}}^{\text{F}} := \Lambda_{\text{dg}}^{\text{int}}$;

variables:

- $\text{ns}^{\text{int}}([\mu]^{\text{int}}) \in [0, \text{n}_{\text{UB}}^{\text{int}}]$, $\mu \in \Lambda_{\text{dg}}^{\text{int}}$: the number of interior-vertices v with $\text{cs}(v) = \mu$;
- $\delta_{\text{ns}}^{\text{X}}(i, [\mu]^{\text{int}}) \in [0, 1]$, $i \in [1, t_{\text{X}}]$, $\mu \in \Lambda_{\text{dg}}^{\text{int}}$, $\text{X} \in \{\text{C}, \text{T}, \text{F}\}$;

constraints:

$$\begin{aligned} \sum_{\mu \in \tilde{\Lambda}_{\text{dg}}^{\text{X}} \cup \{\epsilon\}} \delta_{\text{ns}}^{\text{X}}(i, [\mu]^{\text{int}}) &= 1, & \sum_{\mu = \text{ad} \in \tilde{\Lambda}_{\text{dg}}^{\text{X}}} [\mathbf{a}]^{\text{int}} \cdot \delta_{\text{ns}}^{\text{X}}(i, [\mu]^{\text{int}}) &= \alpha^{\text{X}}(i), \\ \sum_{\mu = \text{ad} \in \tilde{\Lambda}_{\text{dg}}^{\text{X}}} d \cdot \delta_{\text{ns}}^{\text{X}}(i, [\mu]^{\text{int}}) &= \text{deg}^{\text{X}}(i), \\ & & i \in [1, t_{\text{X}}], \text{X} \in \{\text{C}, \text{T}, \text{F}\}, \end{aligned} \quad (78)$$

$$\sum_{i \in [1, t_{\text{C}}]} \delta_{\text{ns}}^{\text{C}}(i, [\mu]^{\text{int}}) + \sum_{i \in [1, t_{\text{T}}]} \delta_{\text{ns}}^{\text{T}}(i, [\mu]^{\text{int}}) + \sum_{i \in [1, t_{\text{F}}]} \delta_{\text{ns}}^{\text{F}}(i, [\mu]^{\text{int}}) = \text{ns}^{\text{int}}([\mu]^{\text{int}}), \quad \mu \in \Lambda_{\text{dg}}^{\text{int}}. \quad (79)$$

D.10 Descriptor for the Number of Edge-configurations

We include constraints to compute the frequency of each edge-configuration in an inferred chemical graph \mathbb{C} .

constants:

- A set Γ^{int} of edge-configurations $\gamma = (\mu, \mu', m)$ with $\mu \leq \mu'$;
- Let $\Gamma_{<}^{\text{int}} = \{(\mu, \mu', m) \in \Gamma^{\text{int}} \mid \mu < \mu'\}$, $\Gamma_{=}^{\text{int}} = \{(\mu, \mu', m) \in \Gamma^{\text{int}} \mid \mu = \mu'\}$ and $\Gamma_{>}^{\text{int}} = \{(\mu', \mu, m) \mid (\mu, \mu', m) \in \Gamma_{<}^{\text{int}}\}$;
- Prepare a coding of the set $\Gamma^{\text{int}} \cup \Gamma_{>}^{\text{int}}$ and let $[\gamma]^{\text{int}}$ denote the coded integer of an element γ in $\Gamma^{\text{int}} \cup \Gamma_{>}^{\text{int}}$;
- Choose subsets $\tilde{\Gamma}_{\text{ec}}^{\text{C}}, \tilde{\Gamma}_{\text{ec}}^{\text{T}}, \tilde{\Gamma}_{\text{ec}}^{\text{CT}}, \tilde{\Gamma}_{\text{ec}}^{\text{TC}}, \tilde{\Gamma}_{\text{ec}}^{\text{F}}, \tilde{\Gamma}_{\text{ec}}^{\text{CF}}, \tilde{\Gamma}_{\text{ec}}^{\text{TF}} \subseteq \Gamma^{\text{int}} \cup \Gamma_{>}^{\text{int}}$; To compute the frequency of edge-configurations exactly, set $\tilde{\Gamma}_{\text{ec}}^{\text{C}} := \tilde{\Gamma}_{\text{ec}}^{\text{T}} := \tilde{\Gamma}_{\text{ec}}^{\text{CT}} := \tilde{\Gamma}_{\text{ec}}^{\text{TC}} := \tilde{\Gamma}_{\text{ec}}^{\text{F}} := \tilde{\Gamma}_{\text{ec}}^{\text{CF}} := \tilde{\Gamma}_{\text{ec}}^{\text{TF}} := \Gamma^{\text{int}} \cup \Gamma_{>}^{\text{int}}$;
- $\text{ec}_{\text{LB}}^{\text{int}}(\gamma), \text{ec}_{\text{UB}}^{\text{int}}(\gamma) \in [0, 2\text{n}_{\text{UB}}^{\text{int}}]$, $\gamma = (\mu, \mu', m) \in \Gamma^{\text{int}}$: lower and upper bounds on the number of interior-edges $e = uv$ with $\text{cs}(u) = \mu$, $\text{cs}(v) = \mu'$ and $\beta(e) = m$;

variables:

- $ec^{\text{int}}([\gamma]^{\text{int}}) \in [ec_{\text{LB}}^{\text{int}}(\gamma), ec_{\text{UB}}^{\text{int}}(\gamma)], \gamma \in \Gamma^{\text{int}}$: the number of interior-edges with edge-configuration γ ;
- $ec_{\text{C}}([\gamma]^{\text{int}}) \in [0, m_{\text{C}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{C}}, ec_{\text{T}}([\gamma]^{\text{int}}) \in [0, t_{\text{T}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{T}}, ec_{\text{F}}([\gamma]^{\text{int}}) \in [0, t_{\text{F}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{F}}$: the number of edges $e^{\text{C}} \in E_{\text{C}}$ (resp., edges $e^{\text{T}} \in E_{\text{T}}$ and edges $e^{\text{F}} \in E_{\text{F}}$) with edge-configuration γ ;
- $ec_{\text{CT}}([\gamma]^{\text{int}}) \in [0, \min\{k_{\text{C}}, t_{\text{T}}\}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{CT}}, ec_{\text{TC}}([\gamma]^{\text{int}}) \in [0, \min\{k_{\text{C}}, t_{\text{T}}\}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{CT}}, ec_{\text{CF}}([\gamma]^{\text{int}}) \in [0, \tilde{t}_{\text{C}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{CF}}, ec_{\text{TF}}([\gamma]^{\text{int}}) \in [0, t_{\text{T}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{TF}}$: the number of edges $e^{\text{CT}} \in E_{\text{CT}}$ (resp., edges $e^{\text{TC}} \in E_{\text{TC}}$ and edges $e^{\text{CF}} \in E_{\text{CF}}$ and $e^{\text{TF}} \in E_{\text{TF}}$) with edge-configuration γ ;
- $\delta_{\text{ec}}^{\text{C}}(i, [\gamma]^{\text{int}}) \in [0, 1], i \in [\widetilde{k_{\text{C}}} + 1, m_{\text{C}}] = I_{(\geq 1)} \cup I_{(0/1)} \cup I_{(=1)}, \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{C}}, \delta_{\text{ec}}^{\text{T}}(i, [\gamma]^{\text{int}}) \in [0, 1], i \in [2, t_{\text{T}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{T}}, \delta_{\text{ec}}^{\text{F}}(i, [\gamma]^{\text{int}}) \in [0, 1], i \in [2, t_{\text{F}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{F}}$: $\delta_{\text{ec}}^{\text{X}}(i, [\gamma]^{\text{t}}) = 1 \Leftrightarrow$ edge e^{X}_i has edge-configuration γ ;
- $\delta_{\text{ec,C}}^{\text{CT}}(k, [\gamma]^{\text{int}}), \delta_{\text{ec,C}}^{\text{TC}}(k, [\gamma]^{\text{int}}) \in [0, 1], k \in [1, k_{\text{C}}] = I_{(\geq 2)} \cup I_{(\geq 1)}, \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{CT}}$: $\delta_{\text{ec,C}}^{\text{CT}}(k, [\gamma]^{\text{int}}) = 1$ (resp., $\delta_{\text{ec,C}}^{\text{TC}}(k, [\gamma]^{\text{int}}) = 1$) \Leftrightarrow edge $e^{\text{CT}}_{\text{tail}(k),j}$ (resp., $e^{\text{TC}}_{\text{head}(k),j}$) for some $j \in [1, t_{\text{T}}]$ has edge-configuration γ ;
- $\delta_{\text{ec,C}}^{\text{CF}}(c, [\gamma]^{\text{int}}) \in [0, 1], c \in [1, \tilde{t}_{\text{C}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{CF}}$: $\delta_{\text{ec,C}}^{\text{CF}}(c, [\gamma]^{\text{int}}) = 1 \Leftrightarrow$ edge $e^{\text{CF}}_{c,i}$ for some $i \in [1, t_{\text{F}}]$ has edge-configuration γ ;
- $\delta_{\text{ec,T}}^{\text{TF}}(i, [\gamma]^{\text{int}}) \in [0, 1], i \in [1, t_{\text{T}}], \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{TF}}$: $\delta_{\text{ec,T}}^{\text{TF}}(i, [\gamma]^{\text{int}}) = 1 \Leftrightarrow$ edge $e^{\text{TF}}_{i,j}$ for some $j \in [1, t_{\text{F}}]$ has edge-configuration γ ;
- $\deg_{\text{T}}^{\text{CT}}(k), \deg_{\text{T}}^{\text{TC}}(k) \in [0, 4], k \in [1, k_{\text{C}}]$: $\deg_{\langle \text{C} \rangle}(v)$ of an end-vertex $v \in V_{\text{T}}$ of the edge $(v^{\text{C}}_{\text{tail}(k)}, v) \in E_{\text{CT}}$ (resp., $(v, v^{\text{C}}_{\text{head}(k)}) \in E_{\text{TC}}$) if any;
- $\deg_{\text{F}}^{\text{CF}}(c) \in [0, 4], c \in [1, \tilde{t}_{\text{C}}]$: $\deg_{\langle \text{C} \rangle}(v)$ of an end-vertex $v \in V_{\text{F}}$ of the edge $(v^{\text{C}}_c, v) \in E_{\text{CF}}$ if any;
- $\deg_{\text{F}}^{\text{TF}}(i) \in [0, 4], i \in [1, t_{\text{T}}]$: $\deg_{\langle \text{C} \rangle}(v)$ of an end-vertex $v \in V_{\text{F}}$ of the edge $(v^{\text{T}}_i, v) \in E_{\text{TF}}$ if any;
- $\Delta_{\text{ec}}^{\text{C}+}(i), \Delta_{\text{ec}}^{\text{C}-}(i) \in [0, 4], i \in [\widetilde{k_{\text{C}}} + 1, m_{\text{C}}], \Delta_{\text{ec}}^{\text{T}+}(i), \Delta_{\text{ec}}^{\text{T}-}(i) \in [0, 4], i \in [2, t_{\text{T}}], \Delta_{\text{ec}}^{\text{F}+}(i), \Delta_{\text{ec}}^{\text{F}-}(i) \in [0, 4], i \in [2, t_{\text{F}}]$: $\Delta_{\text{ec}}^{\text{X}+}(i) = \Delta_{\text{ec}}^{\text{X}-}(i) = 0$ (resp., $\Delta_{\text{ec}}^{\text{X}+}(i) = \deg_{\langle \text{C} \rangle}(u)$ and $\Delta_{\text{ec}}^{\text{X}-}(i) = \deg_{\langle \text{C} \rangle}(v)$) \Leftrightarrow edge $e^{\text{X}}_i = (u, v) \in E_{\text{X}}$ is used in $\langle \text{C} \rangle$ (resp., $e^{\text{X}}_i \notin E(\langle \text{C} \rangle)$);
- $\Delta_{\text{ec}}^{\text{CT}+}(k), \Delta_{\text{ec}}^{\text{CT}-}(k) \in [0, 4], k \in [1, k_{\text{C}}] = I_{(\geq 2)} \cup I_{(\geq 1)}$: $\Delta_{\text{ec}}^{\text{CT}+}(k) = \Delta_{\text{ec}}^{\text{CT}-}(k) = 0$ (resp., $\Delta_{\text{ec}}^{\text{CT}+}(k) = \deg_{\langle \text{C} \rangle}(u)$ and $\Delta_{\text{ec}}^{\text{CT}-}(k) = \deg_{\langle \text{C} \rangle}(v)$) \Leftrightarrow edge $e^{\text{CT}}_{\text{tail}(k),j} = (u, v) \in E_{\text{CT}}$ for some $j \in [1, t_{\text{T}}]$ is used in $\langle \text{C} \rangle$ (resp., otherwise);
- $\Delta_{\text{ec}}^{\text{TC}+}(k), \Delta_{\text{ec}}^{\text{TC}-}(k) \in [0, 4], k \in [1, k_{\text{C}}] = I_{(\geq 2)} \cup I_{(\geq 1)}$: Analogous with $\Delta_{\text{ec}}^{\text{CT}+}(k)$ and $\Delta_{\text{ec}}^{\text{CT}-}(k)$;
- $\Delta_{\text{ec}}^{\text{CF}+}(c), \Delta_{\text{ec}}^{\text{CF}-}(c) \in [0, 4], c \in [1, \tilde{t}_{\text{C}}]$: $\Delta_{\text{ec}}^{\text{CF}+}(c) = \Delta_{\text{ec}}^{\text{CF}-}(c) = 0$ (resp., $\Delta_{\text{ec}}^{\text{CF}+}(c) = \deg_{\langle \text{C} \rangle}(u)$ and $\Delta_{\text{ec}}^{\text{CF}-}(c) = \deg_{\langle \text{C} \rangle}(v)$) \Leftrightarrow edge $e^{\text{CF}}_{c,j} = (u, v) \in E_{\text{CF}}$ for some $j \in [1, t_{\text{F}}]$ is used in $\langle \text{C} \rangle$ (resp., otherwise);
- $\Delta_{\text{ec}}^{\text{TF}+}(i), \Delta_{\text{ec}}^{\text{TF}-}(i) \in [0, 4], i \in [1, t_{\text{T}}]$: Analogous with $\Delta_{\text{ec}}^{\text{CF}+}(c)$ and $\Delta_{\text{ec}}^{\text{CF}-}(c)$;

constraints:

$$\begin{aligned}
\text{ec}_C([\gamma]^{\text{int}}) &= 0, & \gamma &\in \Gamma^{\text{int}} \setminus \tilde{\Gamma}_{\text{ec}}^{\text{C}}, \\
\text{ec}_T([\gamma]^{\text{int}}) &= 0, & \gamma &\in \Gamma^{\text{int}} \setminus \tilde{\Gamma}_{\text{ec}}^{\text{T}}, \\
\text{ec}_F([\gamma]^{\text{int}}) &= 0, & \gamma &\in \Gamma^{\text{int}} \setminus \tilde{\Gamma}_{\text{ec}}^{\text{F}}, \\
\text{ec}_{\text{CT}}([\gamma]^{\text{int}}) &= 0, & \gamma &\in \Gamma^{\text{int}} \setminus \tilde{\Gamma}_{\text{ec}}^{\text{CT}}, \\
\text{ec}_{\text{TC}}([\gamma]^{\text{int}}) &= 0, & \gamma &\in \Gamma^{\text{int}} \setminus \tilde{\Gamma}_{\text{ec}}^{\text{TC}}, \\
\text{ec}_{\text{CF}}([\gamma]^{\text{int}}) &= 0, & \gamma &\in \Gamma^{\text{int}} \setminus \tilde{\Gamma}_{\text{ec}}^{\text{CF}}, \\
\text{ec}_{\text{TF}}([\gamma]^{\text{int}}) &= 0, & \gamma &\in \Gamma^{\text{int}} \setminus \tilde{\Gamma}_{\text{ec}}^{\text{TF}},
\end{aligned} \tag{80}$$

$$\begin{aligned}
\sum_{(\mu, \mu', m) = \gamma \in \Gamma^{\text{int}}} \text{ec}_C([\gamma]^{\text{int}}) &= \sum_{i \in [\tilde{k}_C + 1, m_C]} \delta_\beta^{\text{C}}(i, m), & m &\in [1, 3], \\
\sum_{(\mu, \mu', m) = \gamma \in \Gamma^{\text{int}}} \text{ec}_T([\gamma]^{\text{int}}) &= \sum_{i \in [2, t_T]} \delta_\beta^{\text{T}}(i, m), & m &\in [1, 3], \\
\sum_{(\mu, \mu', m) = \gamma \in \Gamma^{\text{int}}} \text{ec}_F([\gamma]^{\text{int}}) &= \sum_{i \in [2, t_F]} \delta_\beta^{\text{F}}(i, m), & m &\in [1, 3], \\
\sum_{(\mu, \mu', m) = \gamma \in \Gamma^{\text{int}}} \text{ec}_{\text{CT}}([\gamma]^{\text{int}}) &= \sum_{k \in [1, k_C]} \delta_\beta^{\text{CT}}(k, m), & m &\in [1, 3], \\
\sum_{(\mu, \mu', m) = \gamma \in \Gamma^{\text{int}}} \text{ec}_{\text{TC}}([\gamma]^{\text{int}}) &= \sum_{k \in [1, k_C]} \delta_\beta^{\text{TC}}(k, m), & m &\in [1, 3], \\
\sum_{(\mu, \mu', m) = \gamma \in \Gamma^{\text{int}}} \text{ec}_{\text{CF}}([\gamma]^{\text{int}}) &= \sum_{c \in [1, \tilde{t}_C]} \delta_\beta^{*\text{F}}(c, m), & m &\in [1, 3], \\
\sum_{(\mu, \mu', m) = \gamma \in \Gamma^{\text{int}}} \text{ec}_{\text{TF}}([\gamma]^{\text{int}}) &= \sum_{c \in [\tilde{t}_C + 1, c_F]} \delta_\beta^{*\text{F}}(c, m), & m &\in [1, 3],
\end{aligned} \tag{81}$$

$$\begin{aligned}
\sum_{\gamma = (\text{ad}, \text{bd}', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{C}}} [(\mathbf{a}, \mathbf{b}, m)]^{\text{int}} \cdot \delta_{\text{ec}}^{\text{C}}(i, [\gamma]^{\text{int}}) &= \sum_{\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{C}}} [\nu]^{\text{int}} \cdot \delta_{\text{ac}}^{\text{C}}(i, [\nu]^{\text{int}}), \\
\Delta_{\text{ec}}^{\text{C}+}(i) + \sum_{\gamma = (\text{ad}, \mu', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{C}}} d \cdot \delta_{\text{ec}}^{\text{C}}(i, [\gamma]^{\text{int}}) &= \text{deg}^{\text{C}}(\text{tail}(i)), \\
\Delta_{\text{ec}}^{\text{C}-}(i) + \sum_{\gamma = (\mu, \text{bd}, m) \in \tilde{\Gamma}_{\text{ec}}^{\text{C}}} d \cdot \delta_{\text{ec}}^{\text{C}}(i, [\gamma]^{\text{int}}) &= \text{deg}^{\text{C}}(\text{head}(i)), \\
\Delta_{\text{ec}}^{\text{C}+}(i) + \Delta_{\text{ec}}^{\text{C}-}(i) &\leq 8(1 - e^{\text{C}}(i)), & i &\in [\tilde{k}_C + 1, m_C], \\
\sum_{i \in [\tilde{k}_C + 1, m_C]} \delta_{\text{ec}}^{\text{C}}(i, [\gamma]^{\text{int}}) &= \text{ec}_C([\gamma]^{\text{int}}), & \gamma &\in \tilde{\Gamma}_{\text{ec}}^{\text{C}},
\end{aligned} \tag{82}$$

$$\begin{aligned}
\sum_{\gamma=(\mathbf{ad}, \mathbf{bd}', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{T}}} [(\mathbf{a}, \mathbf{b}, m)]^{\text{int}} \cdot \delta_{\text{ec}}^{\text{T}}(i, [\gamma]^{\text{int}}) &= \sum_{\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{T}}} [\nu]^{\text{int}} \cdot \delta_{\text{ac}}^{\text{T}}(i, [\nu]^{\text{int}}), \\
\Delta_{\text{ec}}^{\text{T}+}(i) + \sum_{\gamma=(\mathbf{ad}, \mu', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{T}}} d \cdot \delta_{\text{ec}}^{\text{T}}(i, [\gamma]^{\text{int}}) &= \deg^{\text{T}}(i-1), \\
\Delta_{\text{ec}}^{\text{T}-}(i) + \sum_{\gamma=(\mu, \mathbf{bd}, m) \in \tilde{\Gamma}_{\text{ec}}^{\text{T}}} d \cdot \delta_{\text{ec}}^{\text{T}}(i, [\gamma]^{\text{int}}) &= \deg^{\text{T}}(i), \\
\Delta_{\text{ec}}^{\text{T}+}(i) + \Delta_{\text{ec}}^{\text{T}-}(i) &\leq 8(1 - e^{\text{T}}(i)), & i \in [2, t_{\text{T}}], \\
\sum_{i \in [2, t_{\text{T}}]} \delta_{\text{ec}}^{\text{T}}(i, [\gamma]^{\text{int}}) &= \text{ec}_{\text{T}}([\gamma]^{\text{int}}), & \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{T}}, \quad (83)
\end{aligned}$$

$$\begin{aligned}
\sum_{\gamma=(\mathbf{ad}, \mathbf{bd}', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{F}}} [(\mathbf{a}, \mathbf{b}, m)]^{\text{int}} \cdot \delta_{\text{ec}}^{\text{F}}(i, [\gamma]^{\text{int}}) &= \sum_{\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{F}}} [\nu]^{\text{int}} \cdot \delta_{\text{ac}}^{\text{F}}(i, [\nu]^{\text{int}}), \\
\Delta_{\text{ec}}^{\text{F}+}(i) + \sum_{\gamma=(\mathbf{ad}, \mu', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{F}}} d \cdot \delta_{\text{ec}}^{\text{F}}(i, [\gamma]^{\text{int}}) &= \deg^{\text{F}}(i-1), \\
\Delta_{\text{ec}}^{\text{F}-}(i) + \sum_{\gamma=(\mu, \mathbf{bd}, m) \in \tilde{\Gamma}_{\text{ec}}^{\text{F}}} d \cdot \delta_{\text{ec}}^{\text{F}}(i, [\gamma]^{\text{int}}) &= \deg^{\text{F}}(i, 0), \\
\Delta_{\text{ec}}^{\text{F}+}(i) + \Delta_{\text{ec}}^{\text{F}-}(i) &\leq 8(1 - e^{\text{F}}(i)), & i \in [2, t_{\text{F}}], \\
\sum_{i \in [2, t_{\text{F}}]} \delta_{\text{ec}}^{\text{F}}(i, [\gamma]^{\text{int}}) &= \text{ec}_{\text{F}}([\gamma]^{\text{int}}), & \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{F}}, \quad (84)
\end{aligned}$$

$$\begin{aligned}
\deg^{\text{T}}(i) + 4(1 - \chi^{\text{T}}(i, k) + e^{\text{T}}(i)) &\geq \deg_{\text{T}}^{\text{CT}}(k), \\
\deg_{\text{T}}^{\text{CT}}(k) &\geq \deg^{\text{T}}(i) - 4(1 - \chi^{\text{T}}(i, k) + e^{\text{T}}(i)), & i \in [1, t_{\text{T}}], \\
\sum_{\gamma=(\mathbf{ad}, \mathbf{bd}', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{CT}}} [(\mathbf{a}, \mathbf{b}, m)]^{\text{int}} \cdot \delta_{\text{ec}, \text{C}}^{\text{CT}}(k, [\gamma]^{\text{int}}) &= \sum_{\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{CT}}} [\nu]^{\text{int}} \cdot \delta_{\text{ac}}^{\text{CT}}(k, [\nu]^{\text{int}}), \\
\Delta_{\text{ec}}^{\text{CT}+}(k) + \sum_{\gamma=(\mathbf{ad}, \mu', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{CT}}} d \cdot \delta_{\text{ec}, \text{C}}^{\text{CT}}(k, [\gamma]^{\text{int}}) &= \deg^{\text{C}}(\text{tail}(k)), \\
\Delta_{\text{ec}}^{\text{CT}-}(k) + \sum_{\gamma=(\mu, \mathbf{bd}, m) \in \tilde{\Gamma}_{\text{ec}}^{\text{CT}}} d \cdot \delta_{\text{ec}, \text{C}}^{\text{CT}}(k, [\gamma]^{\text{int}}) &= \deg_{\text{T}}^{\text{CT}}(k), \\
\Delta_{\text{ec}}^{\text{CT}+}(k) + \Delta_{\text{ec}}^{\text{CT}-}(k) &\leq 8(1 - \delta_{\chi}^{\text{T}}(k)), & k \in [1, k_{\text{C}}], \\
\sum_{k \in [1, k_{\text{C}}]} \delta_{\text{ec}, \text{C}}^{\text{CT}}(k, [\gamma]^{\text{int}}) &= \text{ec}_{\text{CT}}([\gamma]^{\text{int}}), & \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{CT}}, \quad (85)
\end{aligned}$$

$$\begin{aligned}
& \deg^T(i) + 4(1 - \chi^T(i, k) + e^T(i + 1)) \geq \deg_{\Gamma}^{\text{TC}}(k), \\
& \deg_{\Gamma}^{\text{TC}}(k) \geq \deg^T(i) - 4(1 - \chi^T(i, k) + e^T(i + 1)), & i \in [1, t_{\Gamma}], \\
& \sum_{\gamma=(\mathbf{ad}, \mathbf{bd}', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{TC}}} [(\mathbf{a}, \mathbf{b}, m)]^{\text{int}} \cdot \delta_{\text{ec}, \text{C}}^{\text{TC}}(k, [\gamma]^{\text{int}}) = \sum_{\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{TC}}} [\nu]^{\text{int}} \cdot \delta_{\text{ac}}^{\text{TC}}(k, [\nu]^{\text{int}}), \\
& \Delta_{\text{ec}}^{\text{TC}+}(k) + \sum_{\gamma=(\mathbf{ad}, \mu', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{TC}}} d \cdot \delta_{\text{ec}, \text{C}}^{\text{TC}}(k, [\gamma]^{\text{int}}) = \deg_{\Gamma}^{\text{TC}}(k), \\
& \Delta_{\text{ec}}^{\text{TC}-}(k) + \sum_{\gamma=(\mu, \mathbf{bd}, m) \in \tilde{\Gamma}_{\text{ec}}^{\text{TC}}} d \cdot \delta_{\text{ec}, \text{C}}^{\text{TC}}(k, [\gamma]^{\text{int}}) = \deg^{\text{C}}(\text{head}(k)), \\
& \Delta_{\text{ec}}^{\text{TC}+}(k) + \Delta_{\text{ec}}^{\text{TC}-}(k) \leq 8(1 - \delta_{\chi}^T(k)), & k \in [1, k_{\text{C}}], \\
& \sum_{k \in [1, k_{\text{C}}]} \delta_{\text{ec}, \text{C}}^{\text{TC}}(k, [\gamma]^{\text{int}}) = \text{ec}_{\text{TC}}([\gamma]^{\text{int}}), & \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{TC}}, \quad (86)
\end{aligned}$$

$$\begin{aligned}
& \deg^F(i) + 4(1 - \chi^F(i, c) + e^F(i)) \geq \deg_{\Gamma}^{\text{CF}}(c), \\
& \deg_{\Gamma}^{\text{CF}}(c) \geq \deg^F(i) - 4(1 - \chi^F(i, c) + e^F(i)), & i \in [1, t_{\Gamma}], \\
& \sum_{\gamma=(\mathbf{ad}, \mathbf{bd}', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{CF}}} [(\mathbf{a}, \mathbf{b}, m)]^{\text{int}} \cdot \delta_{\text{ec}, \text{C}}^{\text{CF}}(c, [\gamma]^{\text{int}}) = \sum_{\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{CF}}} [\nu]^{\text{int}} \cdot \delta_{\text{ac}}^{\text{CF}}(c, [\nu]^{\text{int}}), \\
& \Delta_{\text{ec}}^{\text{CF}+}(c) + \sum_{\gamma=(\mathbf{ad}, \mu', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{CF}}} d \cdot \delta_{\text{ec}, \text{C}}^{\text{CF}}(c, [\gamma]^{\text{int}}) = \deg^{\text{C}}(c), \\
& \Delta_{\text{ec}}^{\text{CF}-}(c) + \sum_{\gamma=(\mu, \mathbf{bd}, m) \in \tilde{\Gamma}_{\text{ec}}^{\text{CF}}} d \cdot \delta_{\text{ec}, \text{C}}^{\text{CF}}(c, [\gamma]^{\text{int}}) = \deg_{\Gamma}^{\text{CF}}(c), \\
& \Delta_{\text{ec}}^{\text{CF}+}(c) + \Delta_{\text{ec}}^{\text{CF}-}(c) \leq 8(1 - \delta_{\chi}^F(c)), & c \in [1, \tilde{t}_{\text{C}}], \\
& \sum_{c \in [1, \tilde{t}_{\text{C}}]} \delta_{\text{ec}, \text{C}}^{\text{CF}}(c, [\gamma]^{\text{int}}) = \text{ec}_{\text{CF}}([\gamma]^{\text{int}}), & \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{CF}}, \quad (87)
\end{aligned}$$

$$\begin{aligned}
& \deg^F(j) + 4(1 - \chi^F(j, i + \tilde{t}_{\text{C}}) + e^F(j)) \geq \deg_{\Gamma}^{\text{TF}}(i), \\
& \deg_{\Gamma}^{\text{TF}}(i) \geq \deg^F(j) - 4(1 - \chi^F(j, i + \tilde{t}_{\text{C}}) + e^F(j)), & j \in [1, t_{\Gamma}], \\
& \sum_{\gamma=(\mathbf{ad}, \mathbf{bd}', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{TF}}} [(\mathbf{a}, \mathbf{b}, m)]^{\text{int}} \cdot \delta_{\text{ec}, \text{T}}^{\text{TF}}(i, [\gamma]^{\text{int}}) = \sum_{\nu \in \tilde{\Gamma}_{\text{ac}}^{\text{TF}}} [\nu]^{\text{int}} \cdot \delta_{\text{ac}}^{\text{TF}}(i, [\nu]^{\text{int}}), \\
& \Delta_{\text{ec}}^{\text{TF}+}(i) + \sum_{\gamma=(\mathbf{ad}, \mu', m) \in \tilde{\Gamma}_{\text{ec}}^{\text{TF}}} d \cdot \delta_{\text{ec}, \text{T}}^{\text{TF}}(i, [\gamma]^{\text{int}}) = \deg^{\text{T}}(i), \\
& \Delta_{\text{ec}}^{\text{TF}-}(i) + \sum_{\gamma=(\mu, \mathbf{bd}, m) \in \tilde{\Gamma}_{\text{ec}}^{\text{TF}}} d \cdot \delta_{\text{ec}, \text{T}}^{\text{TF}}(i, [\gamma]^{\text{int}}) = \deg_{\Gamma}^{\text{TF}}(i), \\
& \Delta_{\text{ec}}^{\text{TF}+}(i) + \Delta_{\text{ec}}^{\text{TF}-}(i) \leq 8(1 - \delta_{\chi}^F(i + \tilde{t}_{\text{C}})), & i \in [1, t_{\Gamma}], \\
& \sum_{i \in [1, t_{\Gamma}]} \delta_{\text{ec}, \text{T}}^{\text{TF}}(i, [\gamma]^{\text{int}}) = \text{ec}_{\text{TF}}([\gamma]^{\text{int}}), & \gamma \in \tilde{\Gamma}_{\text{ec}}^{\text{TF}}, \quad (88)
\end{aligned}$$

$$\begin{aligned}
\sum_{X \in \{C, T, F, CT, TC, CF, TF\}} (\text{ec}_X([\gamma]^{\text{int}}) + \text{ec}_X([\bar{\gamma}]^{\text{int}})) &= \text{ec}^{\text{int}}([\gamma]^{\text{int}}), & \gamma \in \Gamma_{<}^{\text{int}}, \\
\sum_{X \in \{C, T, F, CT, TC, CF, TF\}} \text{ec}_X([\gamma]^{\text{int}}) &= \text{ec}^{\text{int}}([\gamma]^{\text{int}}), & \gamma \in \Gamma_{=}^{\text{int}}.
\end{aligned} \tag{89}$$

D.11 Constraints for Prediction Functions

This section introduces an MILP that simulates the computation process of a prediction function constructed with elastic linear regression.

Let $x = (x(1), x(2), \dots, x(K)) \in \mathbb{R}^K$ denote the feature function $f(\mathbb{C})$ of a chemical graph \mathbb{C} . Let $c_{\min}(j)$ (resp., $c_{\max}(j)$) denote the minimum (resp., maximum) values of the j -th descriptor in a data set D_π for a chemical property π . Let $\text{atm}_{\text{LB}} \in \mathbb{Z}_+$ (resp., $\text{atm}_{\text{UB}} \in \mathbb{Z}_+$) be a lower bound (resp., an upper bound) on the number of atoms in a chemical graph \mathbb{C} to be inferred. Let $\text{mass}(\mathbf{a})$ denote the observed mass of a chemical element $\mathbf{a} \in \Lambda$, and define $\text{mass}^*(\mathbf{a}) \triangleq \lfloor 10 \cdot \text{mass}(\mathbf{a}) \rfloor$. Let $\text{Ms}_{\text{LB}} \in \mathbb{Z}_+$ (resp., $\text{Ms}_{\text{UB}} \in \mathbb{Z}_+$) be a lower bound (resp., an upper bound) on the sum $\frac{1}{|V(H)|} \sum_{v \in V(H)} \text{mass}^*(\alpha(v))$ in a chemical graph \mathbb{C} to be inferred. Let j_{ms} denote the index $j \in [1, K]$ such that the j -th descriptor $\text{dcp}_j(\mathbb{C})$ is the average mass $\overline{\text{ms}}(\mathbb{C}) = \frac{1}{|V(H)|} \sum_{v \in V(H)} \text{mass}^*(\alpha(v))$. Assume that all other descriptors $\text{dcp}_j(\mathbb{C}), j \in [1, K] \setminus \{j_{\text{ms}}\}$ are integers.

Let $\eta_{\Psi, w, b}$ be a prediction function obtained by elastic linear regression, where $\Psi = \{\phi_j \mid j \in [0, K]\}$ and (w, b) is a hyperplane.

We first normalize each of the sets of descriptors $x(j), j \in [1, K]$ and the set of observed values $a(\mathbb{C})$ before we apply the prediction function to compute a predicted value $\eta_{\Psi, w, b}(f(\mathbb{C}))$ of a chemical graph \mathbb{C} , where the set $\{a_i \mid i \in [1, m]\}$ of observed values in the data set D_π is converted into a set $\{\phi_0^{-1}(\frac{a_i - \underline{a}}{\bar{a} - \underline{a}}) \mid i \in [1, m]\}$, where \underline{a} (resp., \bar{a}) denotes the minimum (resp., maximum) value of $a(\mathbb{C})$ over the chemical graphs $\mathbb{C} \in D_\pi$.

Let \underline{y}^* and \bar{y}^* be lower and upper bounds on the predicted value $\eta_{\Psi, w, b}(f(\mathbb{C}))$ of a target chemical graph \mathbb{C} , respectively.

We first converted them into $\phi_0^{-1}(\frac{\underline{y}^* - \underline{a}}{\bar{a} - \underline{a}})$ and $\phi_0^{-1}(\frac{\bar{y}^* - \underline{a}}{\bar{a} - \underline{a}})$. We denote by $\psi_j(t)$ the function $\phi_j((t - c_{\min}(j)) / (c_{\max}(j) - c_{\min}(j)))$. We pre-compute the values $\psi_j(s)$ for all integers $s \in [c_{\min}(j), c_{\max}(j)]$, $j \in [1, K] \setminus \{j_{\text{ms}}\}$ (resp., $\phi_{j_{\text{ms}}}(\frac{s/i - c_{\min}(j_{\text{ms}})}{c_{\max}(j_{\text{ms}}) - c_{\min}(j_{\text{ms}})})$ for all integers $s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}]$ and $i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]$) as constants.

An MILP that simulates the computation process of a prediction function $\eta_{\Psi, w, b}$ is described as follows.

$\mathcal{M}(x, y; \mathcal{C}_1)$:

constants:

- A hyperplane (w, b) with $w \in \mathbb{R}^K$ and $b \in \mathbb{R}$;
- Activation functions $\phi_j : \mathbb{R} \rightarrow \mathbb{R}, j \in [0, K]$;
- Real values $\underline{y}^*, \bar{y}^* \in \mathbb{R}$ such that $\underline{y}^* < \bar{y}^*$; Set $\underline{y}^{**} := \phi_0^{-1}(\frac{\underline{y}^* - \underline{a}}{\bar{a} - \underline{a}})$ and $\bar{y}^{**} := \phi_0^{-1}(\frac{\bar{y}^* - \underline{a}}{\bar{a} - \underline{a}})$.
- $c_{\min}(j), c_{\max}(j) \in \mathbb{R}, j \in [1, K]$: the minimum and maximum values of the j -th descriptor in the data set D_π , respectively;

- Reals $\Delta(j, s) \in \mathbb{R}, j \in [1, K] \setminus \{j_{\text{ms}}\}, s \in [c_{\min}(j), c_{\max}(j)]$: $\Delta(j, c_{\min}(j)) := \psi_j(c_{\min}(j))$ and $\Delta(j, s) := \psi_j(s) - \psi_j(s-1), s \in [c_{\min}(j) + 1, c_{\max}(j)]$;
- $\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}} \in \mathbb{Z}_+$: lower and upper bounds on the number of atoms in a chemical graph \mathbb{C} to be inferred; Set $\text{atm}_{\text{LB}} := n_{\text{LB}} + \text{na}_{\text{LB}}(\mathbf{H})$ and $\text{atm}_{\text{UB}} := n^* + \text{na}_{\text{UB}}(\mathbf{H})$;
- $\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}} \in \mathbb{Z}_+$: lower and upper bounds on the sum $\sum_{v \in V(H)} \text{mass}^*(\alpha(v))$; For example, set $\text{Ms}_{\text{LB}} := \lfloor \min\{\text{mass}^*([\mathbf{a}]) \mid \mathbf{a} \in \Lambda, \text{val}(a) = 1\} \cdot (3n_{\text{LB}}/4) + \min\{\text{mass}^*([\mathbf{a}]) \mid \mathbf{a} \in \Lambda, \text{val}(a) \geq 2\} \cdot (n_{\text{LB}}/4) + \text{mass}^*([\mathbf{H}])\text{na}_{\text{LB}}(\mathbf{H}) \rfloor$ and $\text{Ms}_{\text{UB}} := n^* \max\{\text{mass}^*([\mathbf{a}]) \mid \mathbf{a} \in \Lambda\} + \text{mass}^*([\mathbf{H}])\text{na}_{\text{UB}}(\mathbf{H})$;
- $M \in \mathbb{R}_+$: an upper bound on $\zeta(x(j_{\text{ms}}))$; For example, set $M := 2\phi_{j_{\text{ms}}}(1)$;
- Define $\zeta(s, i) \triangleq \phi_{j_{\text{ms}}}\left(\frac{s/i - c_{\min}(j_{\text{ms}})}{c_{\max}(j_{\text{ms}}) - c_{\min}(j_{\text{ms}})}\right), s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}], i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]$, where $\psi_{j_{\text{ms}}}(s/i) = \zeta(s, i)$;
- Reals $\Delta_{\text{Ms}}(s, i) \in \mathbb{R}, s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}], i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]$: $\Delta_{\text{Ms}}(\text{Ms}_{\text{LB}}, i) := \zeta(\text{Ms}_{\text{LB}}, i)$ and $\Delta_{\text{Ms}}(s, i) := \zeta(s, i) - \zeta(s-1, i), s \in [\text{Ms}_{\text{LB}} + 1, \text{Ms}_{\text{UB}}], i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]$;
- A real $\varepsilon(j) > 0, j \in [1, K]$: a tolerance. For example, set $\varepsilon(j) := \frac{1}{10^5} \min\{\Delta(j, s) \mid s \in [c_{\min}(j), c_{\max}(j)]\}, j \in [1, K] \setminus \{j_{\text{ms}}\}$ and $\varepsilon(j_{\text{ms}}) := \frac{1}{10^5} \min\{\Delta_{\text{Ms}}(s, i) \mid s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}], i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]\}$;

variables:

- Real variables $\hat{x}(j) \in \mathbb{R}, j \in [1, K]$: $\hat{x}(j)$ represents $\psi_j(x(j))$;
- Integer variables $x(j) \in [c_{\min}(j), c_{\max}(j)], j \in [1, K] \setminus \{j_{\text{ms}}\}$: $x(j)$ represents the j -th descriptor in an MILP $\mathcal{M}(x, g; \mathcal{C}_2)$;
- A real variable $x(j_{\text{ms}}) \in \mathbb{R}_+$ with $c_{\min}(j_{\text{ms}}) \leq x(j_{\text{ms}}) \leq c_{\max}(j_{\text{ms}})$: $x(j_{\text{ms}})$ represents the average mass $\overline{\text{ms}}(\mathbb{C})$ in an MILP $\mathcal{M}(x, g; \mathcal{C}_2)$;
- Binary variables $\delta(j, s) \in [0, 1], j \in [1, K] \setminus \{j_{\text{ms}}\}, s \in [c_{\min}(j), c_{\max}(j)]$: $\delta(j, s) = 1 \Leftrightarrow x(j) \geq s$;
- Binary variables $\delta_{\text{atm}}(i) \in [0, 1], i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]$: $\delta_{\text{atm}}(i) = 1 \Leftrightarrow |V(H)| = i$;
- Binary variables $\delta_{\text{Ms}}(s) \in [0, 1], s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}]$: $\delta_{\text{Ms}}(s) = 1 \Leftrightarrow \sum_{v \in V(H)} \text{mass}^*(\alpha(v)) \geq s$;

constraints:

$$\underline{y}^{**} \leq \sum_{j \in [1, K]} w(j) \hat{x}(j) + b \leq \bar{y}^{**}, \quad (90)$$

$$\begin{aligned} & \sum_{s \in [c_{\min}(j), c_{\max}(j)]} \delta(j, s) + c_{\min}(j) - 1 = x(j), \\ & \delta(j, s) \geq \delta(j, s+1), \quad s \in [c_{\min}(j), c_{\max}(j) - 1], \\ & \sum_{s \in [c_{\min}(j), c_{\max}(j)]} \Delta(j, s) \delta(j, s) - \varepsilon(j) \leq \hat{x}(j) \leq \sum_{s \in [c_{\min}(j), c_{\max}(j)]} \Delta(j, s) \delta(j, s) + \varepsilon(j), \quad j \in [1, K] \setminus \{j_{\text{ms}}\}, \end{aligned} \quad (91)$$

$$\begin{aligned}
\sum_{i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]} \delta_{\text{atm}}(i) &= 1, \\
\sum_{i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}]} i \cdot \delta_{\text{atm}}(i) &= n_G + \text{na}^{\text{ex}}([\text{H}]^{\text{ex}}), \\
\sum_{\mathbf{a} \in \Lambda} \text{mass}^*(\mathbf{a}) \cdot \text{na}([\mathbf{a}]) &= \sum_{s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}]} \delta_{\text{Ms}}(s) + \text{Ms}_{\text{LB}} - 1,
\end{aligned} \tag{92}$$

$$\delta_{\text{Ms}}(s) \geq \delta_{\text{Ms}}(s+1), \quad s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}} - 1], \tag{93}$$

$$\begin{aligned}
\sum_{s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}]} \Delta_{\text{Ms}}(s, i) \delta_{\text{Ms}}(s) - M \cdot (1 - \delta_{\text{atm}}(i)) - \epsilon(j_{\text{ms}}) &\leq \widehat{x}(j_{\text{ms}}) \leq \\
\sum_{s \in [\text{Ms}_{\text{LB}}, \text{Ms}_{\text{UB}}]} \Delta_{\text{Ms}}(s, i) \delta_{\text{Ms}}(s) + M \cdot (1 - \delta_{\text{atm}}(i)) + \epsilon(j_{\text{ms}}), &\quad i \in [\text{atm}_{\text{LB}}, \text{atm}_{\text{UB}}].
\end{aligned} \tag{94}$$