Inferring Chemical Reaction Patterns Using Rule Composition in Graph Grammars

Jakob Lykke Andersen¹, Christoph Flamm², Daniel Merkle¹, Peter F. Stadler²⁻⁷

¹ Department for Mathematics and Computer Science, University of Southern Denmark, Campusvej 55, DK-5230 Odense M, Denmark ² Institute for Theoretical Chemistry, University of Vienna, Währingerstraße 17, A-1090 Wien, Austria. ³ Bioinformatics Group, Department of Computer Science, and Interdisciplinary Center for Bioinformatics, Härtelstraße 16-18, D-04107, Leipzig, Germany. ⁴ Max Planck Institute for Mathematics in the Sciences, Inselstraße 22 D-04103 Leipzig, Germany. ⁵ Fraunhofer Institute for Cell Therapy and Immunology, Perlickstraße 1, D-04103 Leipzig, Germany. ⁶ Center for non-coding RNA in Technology and Health, University of Copenhagen, Grønnegårdsvej 3, DK-1870 Frederiksberg C, Denmark. ⁷ Santa Fe Institute, 1399 Hyde Park Rd, Santa Fe, NM 87501, USA

 $\label{eq:constraint} Email: jakan 06 @ student.sdu.dk; CF*:xtof @ tbi.univie.ac.at; DM*:daniel @ imada.sdu.dk; PFS:studla @ bioinf.uni-leipzig.de; and a student st$

*Corresponding author

Abstract

Background: Modeling molecules as undirected graphs and chemical reactions as graph rewriting operations is a natural and convenient approach to modeling chemistry. Graph grammar rules are most naturally employed to model elementary reactions like merging, splitting, and isomerisation of molecules. It is often convenient, in particular in the analysis of larger systems, to summarize several subsequent reactions into a single composite chemical reaction.

Results: We introduce a generic approach for composing graph grammar rules to define a chemically useful rule compositions. We iteratively apply these rule compositions to elementary transformations in order to automatically infer complex transformation patterns. As an application we automatically derive the overall reaction pattern of the Formose cycle, namely two carbonyl groups that can react with a bound glycolaldehyde to a second glycolaldehyde. Rule composition also can be used to study polymerization reactions as well as more complicated iterative reaction schemes. Terpenes and the polyketides, for instance, form two naturally occurring classes of compounds of utmost pharmaceutical interest that can be understood as "generalized polymers" consisting of five-carbon (isoprene) and two-carbon units, respectively.

Conclusion: The framework of graph transformations provides a valuable set of tools to generate and investigate large networks of chemical networks. Within this formalism, rule composition is a canonical technique to obtain coarse-grained representations that reflect, in a natural way, "effective" reactions that are obtained by lumping together specific combinations of elementary reactions.

Introduction

Directed hypergraphs [1] are a suitable topological representation of (bio)chemical reaction networks where (catalytic) reactions are hyperedges connecting substrate nodes to product nodes. Such networks require an underlying Artificial Chemistry [2] that describes how molecules and reactions are modeled. If molecules are treated as edge and vertex labeled graphs, where the vertex labels correspond to atom types and the edge labels denote bond types, then structural change of molecules during chemical reactions can be modeled as graph rewrite [3]. In contrast to many other Artificial Chemistries this approach allows for respecting fundamental rules of chemical transformations like mass conservation, atomic types, and cyclic shifts of electron pairs in reactions. In general, a graph rewrite (rule) transforms a set of substrate graphs into a set of product graphs. Hence the graph rewrite formalism allows not only to delimit an entire chemical universe in an abstract but compact form but also provides a methodology for its explicit construction.

Most methods for the analysis of this network structure are directed towards this graph (or hypergraph) structure [1, 4], which is described by the stoichiometric matrix \mathbf{S} of the chemical system. Since \mathbf{S} is essentially the incidence matrix of the directed hypergraph, algebraic approaches such as Metabolic Flux Analysis and Flux Balance Analysis [5] have a natural interpretation in terms of the hypergraph. Indeed typical results are sets of possibly weighted reactions (i.e., hyperedges) such as elementary flux modes [6], extreme pathways [7], minimal metabolic behaviors [8] or a collection of reactions that maximize the production of a desired product in metabolic engineering. The net reaction of a given pathway is simply the linear combination of the participating hyperedges.

In the setting of generative models of chemistry, each concrete reaction is not only associated with its stoichiometry but also with the transformation rule operating on the molecules that are involved in a particular reaction. Importantly, these rules are formulated in terms of *reaction mechanisms* that readily generalize to large sets of structurally related molecules. It is thus of interest to derive not only the stoichiometric net reaction of a pathway but also the corresponding "effective transformation rule". Instead of attempting to address this issue *a posteriori*, we focus here on the possibility of composing the elementary rules of chemical transformations to new effective rules that encapsulate entire pathways.

The motivation comes from the observation that string grammars are meaningfully characterized and understood by investigating the transformation rules. Consider, as a trivial example, the contextfree grammar G with the starting symbol S and the rules $S \to aS'a$, $S' \to aS'a \mid B$ and $B \to \epsilon \mid bB$. Inspecting this grammar we see that we can summarize the effect of the productions as $B \to b^k, k \ge 0$, and $S \to a^n B a^n, n \ge 1$. The language generated from \mathcal{G} is thus $\{a^n b^* a^n | n \ge 1\}$. Here we explore whether a similar reasoning, namely the systematic combination of transformation rules, can help to characterize the language of molecules that is generated by a particular graph rewriting chemistry. Similar to the example from term rewriting above, we should at the very least be able to recognize the regularities in polymerization reactions. We shall see below, however, that the rule based approach holds much higher promises.

Chemical reactions can be readily composed to "overall reactions" such as the net transformation of metabolic pathways. This observation is used implicitly in flux balance analysis at the level of the stoichiometric matrix. Recently, [9] considered the composition of concerte chemical reactions, i.e., transformations of complete molecules, as a means of reconstructing metabolic pathways. In this contribution we take a different point of view: instead of asking for concrete overall reactions, we are concerned with the composition of the underlying reaction mechanisms themselves. As we will see, these can be applied to arbitary molecule contexts. We therefore address two issues: First we establish the formal conditions under which chemical transformation rules *can* be meaningfully composed. To this end, we discuss rule composition within the framework of concurrency theory in the following section. We then investigate the specific restrictions that apply to chemical systems, leading to a constructive approach for inferring composite rules.

The basic computational task we envision starts from an unordered set \mathcal{R} of reactions such as those forming a particular metabolic reaction pathway. To derive the effective transformation rule describing the pathway we need to find the correct ordering π in which the transformation rules p_i , underlying the individual chemical reactions ρ_i , have to be composed. We illustrate this approach in some detail using the Formose reaction as an example in the Results section.

Graph Grammars and Rule Composition

Graph grammars, or graph rewriting systems, are proper generalizations of term rewrite systems. A wide variety of formal frameworks have been explored, including several different algebraic ones rooted in category theory. We base our conceptual developments on the *double pushout* (DPO) formulation of graph transformations. For the comprehensive treatise of this framework we refer to [10]. In the following sections we first outline the basic setup and then introduce full and partial rule composition. Alternative approaches to graph rewriting in the context of (artificial) chemistry have been based on the single pushout (SPO) model of graph transformations, see e.g. [11, 12]. We briefly discuss the rather technical difference between the DPO and SPO framework in Appendix A, where we also briefly outline our reasons for choosing DPO.

Double Pushout and Concurrency

The DPO formulation of graph transformations considers transformation rules of form $p = (L \stackrel{l}{\leftarrow} K \stackrel{r}{\rightarrow} R)$ where L, R, and K are called the left graph, right graph, and context graph, respectively. The maps land r are graph morphisms. The rule p transforms G to H, in symbols $G \stackrel{p,m}{\Longrightarrow} H$ if there is a pushout graph D and a "matching morphism" $m : L \to G$ such that following diagram is valid:

The existence of D is equivalent to the so-called *gluing condition*, which determines whether the rule p is applicable to a match in G. In the following we will also write $G \stackrel{p}{\Rightarrow} H$ and $G \Rightarrow H$ for derivations, if the specific match or transformation rule is unimportant or clear from the context.

Concurrency theory provides a canonical framework for the composition of two graph transformations. Given two rules $p_i = (L_i \stackrel{l_i}{\leftarrow} K_i \stackrel{r_i}{\to} R_i),$ i = 1, 2, a composition $(L \stackrel{q_l}{\leftarrow} K \stackrel{q_r}{\to} R) = p_1 *_E p_2$ can be defined whenever a dependency graph E exists so that in the following diagram:

the cycles (1) and (2) are pushouts, and (3) is a pullback, see e.g., [13]. We then have $q_l = s_1 \circ w_1$ and $q_r = t_2 \circ w_2$. The concurrency theorem [14] ensures that for any sequence of consecutive direct transformations $G \xrightarrow{p_1,m_1} H \xrightarrow{p_2,m_2} G'$ a graph E, a corresponding E-concurrent rule $p_1 *_E p_2$, and a morphism m can be found such that $G \xrightarrow{p_1 *_E p_2 *_m} G'$.

In order to use graph transformation as a model for chemical reactions additional conditions must be enforced. Most importantly, atoms are neither created, nor destroyed, nor transformed to other types. Thus only graph morphisms whose restriction to the vertex sets are bijective are valid in our context. In particular, the matching morphism m always corresponds to a subgraph isomorphism in our context. The context graph K thus is (isomorphic to) a subgraph of both L and R, describing the part of Lthat remains unchanged in R. Conservation of atoms means that the vertex sets of L, K, and R are linked by bijections known as the atom-mapping. When the atom mapping is clear, thus, we do not need to represent the context explicitly.

It is important to note that the existence of the matching morphism $m: L \to G$ alone is not sufficient to guarantee the applicability of the transformation. In our context, we require in addition that the transformation rule does not attempt to introduce an edge in R that has been present already before the transformation is applied. Formally, the gluing condition requires that $(l(x), l(y)) \notin L$ and $(r(x), r(y)) \in R$ implies $(m(l(x)), m(r(y))) \notin G$.

Full Rule Composition

In the following we will be concerned only with special, chemically motivated, types of rule compositions. In the simplest case the dependency graph E is isomorphic to R_1 , later we will also consider a more general setting in which E is isomorphic to the disjoint union of R_1 and some connected components of L_2 . For the ease of notation from now on we only refer to a rule composition, and not to a composition of morphisms as in the Graph Grammar section, i.e., $p_1 *_E p_2$ will be denoted as $p_2 \circ p_1$ (note the order of the arguments changes). If $E \cong R_1$, then $L_2 \cong e_2(L_2)$ is a subgraph of R_1 . Omitting the explicit references to the subgraph matching morphism e_2 we can simply view L_2 as subgraph of R_1 as illustrated in Figure 1.

The rule composition thus amounts to a rewrit-



Fig. 1: Full composition of two rules requires that L_2 is (isomorphic to) a subgraph of R_1 .

ing $R_1 \xrightarrow{p_2,e_2} R$, while the left side L_1 is preserved. We will use the notation $p_2 \circ p_1$ and $G' \xrightarrow{p_2 \circ p_1} G'$ for this restricted type of rule composition, and call it *full* composition as the complete left side of p_2 is a subgraph of R_1 . Note that L_2 may fit into R_1 in more than one way so that there may be more than one composite rule. Formally, the alternative compositions are distinguished by different matching morphisms e_2 in the diagram (2); we will return to this point below.

Partial Rule Composition

An important issue for the application to chemical reactions is that the graphs involved in the rules are in general not connected. Typical chemical reactions combine molecules, split molecules or transfer groups of atoms from one molecule to another. The transformation rules for all these reactions therefore require multiple connected components. For the purpose of dealing with these rules, we introduce the following notation for graphs and derivations.

Let Q be a graph with #Q connected components $Q_i, i = 1, \ldots, \#Q$. It will be convenient to treat Q as the multiset of its components. A typical chemical graph derivation, corresponding to a bi-molecular reaction can be written in the form $\{G^1, G^2\} \xrightarrow{p,m} \{H^1, H^2, H^3\}$, where we take the notation to imply that all graphs G^i and H^j are connected. We will furthermore insist that representations of chemical reactions are minimal in the following sense: If the left graph of the rule $p = (L \leftarrow K \rightarrow R)$ matches entirely within G^1 , i.e., $m(L) \cap G^2 = \emptyset$, then G^2 can be omitted. (In a chemical rewriting grammar, then, one of the H^i must be isomorphic to G^2 , becoming redundant as well.) More formally, we say that a derivation $\{G^1, G^2, \dots, G^{\#G}\} \xrightarrow{p,m} \{H^1, H^2, \dots, H^{\#H}\}$ is proper if

$$\forall i, j : G_i \cong H_j \Rightarrow G_i \cap m(L) \neq \emptyset$$

That is, a proper derivation cannot be simplified. If the derivation $G \xrightarrow{p} H$ is proper then $\#G \leq \#L$. The inequality comes from that fact that multiple components of L may easily be matched to a single component of G while each component of L must match within a component of G.

The conditions for the \circ composition of rules are a bit too strict for our applications. We thus relax them respect the component structure of left and right graphs. More precisely, we require that E is isomorphic to a disjoint union of a copy of R_1 and some connected components of L_2 so that for every connected component L_2^i of L_2 holds that either $e_2(L_2^i) \subseteq e_1(R_1)$ or $e_2(L_2^i)$ is a connected component of E isomorphic to L_2^i . For a rule composition of this type to be well defined we need that $\exists i$ such that $e_2(L_2^i) \subseteq e_1(R_1)$ holds. We remark that the latter condition could be relaxed further to lead to additional compositions for which left and right sides are disjoint unions.

The composition of $p_1 = (L_1, K_1, R_1)$ and $p_2 = (L_2, K_2, R_2)$ now yields $p_2 \circ p_1 = (\{L_1, L_2^2\}, K_3, R_3)$ (cmp. Figure 2). Note that right graph R_3 cannot no longer be regarded simply as a rewritten version of R_1 because rule p_2 now adds additional vertices to both the left and the right graph. The composite context K_3 contains only subsets of K_1 and K_2 , but it is expanded by the vertices of L_2^2 and the edges of L_2^2 that remain unchanged under rule p_2 .

An example of a full rule composition is shown in Fig. 3. The two rules in the example, which in this case are also chemical reactions, are part of the Formose grammar. The Formose grammar consists of two pairs of rules. The first pair of rules, (from now on denoted as p_0 and p_1), implements both directions of the keto-enol tautomerism. One direction, p_1 , is visualized in Fig. 3. The second pair, (from now on denoted as p_2 , p_3) is the aldol-addition and its reverse respectively. The reverse (p_3) is also visualized in Fig. 3. We see that the left side of p_1 is isomorphic to a subgraph of one of the components of the right side of p_3 . Composing the two rules by subgraph matching yields a third rule, $p_1 \circ p_3$.

In general, we require here that the connected components of R_1 and L_2 satisfy either $e_2(L_2^i) \subseteq$



Fig. 3: Composition of two rules from the Formose reaction $(p_1 \text{ and } p_3)$; the following rule names will be used: p_0 : forward keto-enol tautomerism (which corresponds to the reverse of p_1), p_1 : backward keto-enol tautomerism, p_2 : forward addol addition (which corresponds to the reverse of p_3), and p_3 : backward addol addition. The atom mapping and matching morphism are implicitly given in these drawings by corresponding positions of the atoms. The context K thus consists of all the atoms as well as the chemical bonds (edges) shown in black in both the left and the right graph of each rule.

 $e_1(R_1^j)$ or $e_1(R_1^j) \cap e_2(L_2^i) = \emptyset$. We furthermore exclude the trivial case of parallel rules in which only the second alternative is realized. In other extreme, if all components L_2^i satisfy $e_1(L_2^i) \subseteq e_2(R_1)$, the partial composition becomes a full composition. Formally, these alternatives are described by different dependency graphs E and/or different morphisms e_1 and e_2 . Pragmatically we can understand this as a matching μ of L_2 and R_1 as in Fig. 4. Specifying μ of course removes the ambiguity from the definition of the rule composition; hence we write $p_2 \circ_{\mu} p_1$ to emphasize the matching μ .

Constructing Rule Compositions

Given two rules, p_1 and p_2 , it is not only interesting to know if a partial composition is defined, but also to create the set of all possible compositions

$$\{p_2 \circ_{\mu_1} p_1, p_2 \circ_{\mu_2} p_1, \dots, p_2 \circ_{\mu_k} p_1\}$$

explicitly. This set in particular contains also all full compositions. The following describes an algorithm for enumerating all partial compositions.

Enumerating the Matchings μ

The key to finding all compositions is the enumeration of all matchings μ that respect out restrictions on overlaps between connected components. We thus start from the sets $\left\{R_1^1, R_1^2, \ldots, R_1^{\#R_1}\right\}$ and $\left\{L_2^1, L_2^2, \ldots, L_2^{\#L_2}\right\}$ of connected components of R_1 and L_2 , resp. In the first set we find all subgraph matches $L_2^i \subseteq R_1^j$ (represented as the corresponding matchings μ_{ij}) and arrange the result in a matrix of lists of subgraph matches, Fig. 5a.

The matching matrix is extended by a virtual column to account for the possibility that L_2^i is not matched with any component of R_1 . Every partial (and full) composition is now defined by a selection of one submatch from each row of the matrix, see Supplemental Material for an example. The converse is not true, however: Not every selection of matches correspond to a partial composition. In particular, we exclude the case that only entries from the virtual column are selected. In addition, the sub-matches must be disjoint to ensure that the combined match is injective. The latter conditions needs to be checked only when more than one submatch is



Fig. 4: The (partial) composition of two rules is mediated by the dependency graph E and the two matching morphisms e_1 and e_2 . Since these are subgraph isomorphisms in our case, E is simply the union $e_1(R_1) \cup e_2(L_2)$. The (partial) match $e_1(R_1) \cap e_2(L_2)$ can be understood as a matching μ between R_1 and L_2 , i.e., as a 1-1 relation of the matching nodes and edges. Whenever an edge is matched, then so are its incident vertices.



Fig. 5: Example of a match matrix and the same matrix with its virtual extension. The top row specifies 1 possibility for $L_2^1 \subseteq R_1^1$ and 2 for $L_2^1 \subseteq R_1^3$. The extended matrix further specifies that L_2^1 can be unmatched. The bottom rows can be interpreted similarly. We display the number of matchings instead of a representation of the matchings themselves.

selected from the same column.

Composing the Rules

The construction of the composition $p_2 \circ_{\mu} p_1$ of two rules p_1 and p_2 does not explicitly depend on the component structure of R_2 and L_1 because it is uniquely defined by the matching μ and the bijections of the nodes of L_i , K_i , and R_i for each of the two rules. We obtain L by extending L_1 with unmatched components of L_2 and R by extending R_2 by the unmatched components of R_1 . The corresponding extension of μ to a bijection $\hat{\mu}$ of the vertex sets of L and R is uniquely defined. The context K of the composite rule simply consists the common vertex set of L and R and all edges (x, y) of L for which $(\hat{\mu}(x), \hat{\mu}(y))$ is an edge in R. We note in passing that $\hat{\mu}$ defines the atom mapping of the composite transformation. The explicit construction of (R, K, L) is summarized as Algorithm 0.1.

The implementation of the algorithm naturally depends heavily on the representation of transformation rules, which in our implementation is the representation from the Graph Grammar Library (GGL) [15]. The representation is a single graph, with attached vertex and edge properties defining membership of L, K and R, as well as the needed labels.

Not all matchings define valid rule composition. For instance, consider an edge (u, v) that is present in R_1 and R_2 but not in L_2 and both u and v are in L_2 . This would amount to creating the edge by means of rule p_2 which was already introduced by p_1 . Since we do not allow parallel edges and thus regard such inconsistencies as undefined cases and reject the matching. Note that a parallel edge does not correspond to a "double bond" (which essentially is only an edge with a specific type).

Graph Binding

The composition of transformation rules, and thereby chemical reactions, makes it possible to create abstract meta-rules in a way that is similar to the combination of multiple functions into more abstract functions in functional programming. A related concept from (functional) programming that seems useful in the context of graph grammars is partial function application. Consider, for example, the binding of the number 2 to the exponentiation operator, yielding either the function $f(x) = 2^x$ or $f(x) = x^2$. In the framework of rule composition, we define graph binding as a special case.

Let G be a graph and $p_2 = (L_2, K_2, R_2)$ be a transformation rule. The binding of G to p_2 results in the transformation rule p = (L, K, R) which implements the partial application of p_2 on G. This is accomplished simply by regarding G as a rule $p_1 =$ $(\emptyset, \emptyset, G)$, and using partial composition; $p = p_2 \circ p_1$. Note that if $p_2 \circ p_1$ is a full composition, then p can be regarded as a graph H and $G \stackrel{p_2}{\Longrightarrow} H$ holds.

Graph binding allows a simplified representation

Algorithm 0.1: Composing p_1 and p_2 to p, by a given partial mapping

Input: $p_1 = (L_1, K_1, R_1)$ **Input**: $p_2 = (L_2, K_2, R_2)$ **Input**: μ , a partial matching between L_2 and R_1 **Output**: p = (L, K, R)1 $p \leftarrow \text{empty rule}$ **2** Copy vertices of p_1 to p**3 foreach** *vertex* $v \in p_2$ **do** 4 if v is not mapped by μ then Copy v to p $\mathbf{5}$ 6 else Change membership in L, K and R for vertex $\mu(v)$ 7 8 Copy edges of p_1 to pforeach $Edge \ e \in p_2$ do 9 if e is not mapped by μ then 10 Copy e to p11 else 12Change membership in L, K and R for edge $\mu(e)$ 13 14 Delete edges and vertices created by p_1 , but deleted by p_2 15 if matching condition not satisfied then abort 16 return p

of reactions. For instance, we can use this formal construction to omit uninteresting ubiquitously present molecules such as water by binding the graph of the water molecule to the transformation rule of a reaction that requires water. Similarly, graph unbinding can be defined as a transformation rule that destroys graphs. In a chemical application it can be used to avoid the explicit representation of uninteresting ubiquitous molecules such as the solvent.

Ordering Rules

A wide variety of methods, including flux balance analysis, can be used to identify pathways or other subsets of reactions that are of interest. Adjacency of reactions in the original networks as well as their directionality can be used efficiently to prune the possible orders of rule compositions. The fact that multiple reactions are instantiations of the same transformation rule, as in the example discussed in detail in the next section, further reduces the search spaces.

Results and Discussion

We illustrate the use of transformation rule composition by deriving of meta-rules from the graph grammar consisting of the four rules necessary to represent the complete Formose reaction, see Fig. 3. The overall reaction pattern of the Formose cycle is $2g_0 + g_1 \rightarrow 2g_1$ with g_0 being formaldehyde and g_1 being glycolaldehyde. It amounts to the linear combination $\sum_{i=1}^{9} \rho_i$ of the eight reactions and the influx ρ_1 of g_0 listed in Fig. 3. It is important to notice that several of these reactions are instantiations of the same, well-known chemical transformations. We have forward keto-enol tautomerism ($\mathbf{p}_1: \rho_2, \rho_4, \rho_6$), backward keto-enol tautomerism ($\mathbf{p}_1: \rho_7, \rho_9$), forward aldol addition ($\mathbf{p}_2: \rho_3, \rho_5$), and backward aldol addition ($\mathbf{p}_3: \rho_8$). The composite rule models the complete autocatalytic cycle shown in Fig. 6 as a single meta-rule.

Throughout this section we will not explicitly distinguish between partial composition and full composition, and we interpret the composition operator \circ as right-associative to simplify the notation. Thus $p_i \circ p_j \circ p_k$ means $p_i \circ (p_j \circ p_k)$.

The rules are used in the autocatalytic cycle in the following order (starting with an keto-enol tautomerisation p_0):

$p_0, p_2, p_0, p_2, p_0, p_1, p_3, p_1\\$

As it is not possible to compose this sequence of rules directly, we start by binding glycolaldehyde g_1 to



Fig. 6: Above: Chemical reaction network for the Formose reaction; hyperedges are labeled with (i, \mathbf{p}_j) where *i* is the *i*-th reaction ρ_i in the rule composition. $\mathbf{p}_j, 0 \le j \le 3$, refers to a specific rule from the Formose reaction; Right: Resulting composed rule after the composition of the first *i* rules along the Formose cycle, context shown in black.

reaction p_0 , as the before-mentioned keto-enol tautomerisation is applied to molecule g_1 . The resulting rule is denoted as g_1 . The hyperedges in the chemical reaction network depicted in Fig. 6 are numbered according to the sequence that reflects in which order the Formose reaction takes place and consequently the order in which the rule composition subsequently is done. The first composition refers to the binding operation. This binding of glycolaldehyde results in a graph grammar rule, which is depicted in row 1 in the table depicted in Fig. 6, i.e., the rule $(\emptyset, \emptyset, g_1)$ (see "Graph Binding"). The numbers at the hyperedges $(2, 3, \ldots 9)$ refer to the second, third, \ldots , ninth reaction in the sequence of reactions given above. The graph grammar rule $p_i, 0 \le i \le 3$, used for the corresponding hyper-edge is given next to the sequence number. The rules inferred by a subsequent rule composition are given in rows 2 to 9 of the table.

The application of the final rule results in the composed meta-rule $p_1 \circ p_3 \circ \ldots \circ p_0 \circ g_1$. This rule precisely covers the reaction pattern of the Formose reaction, namely how two formaldehyde molecules and one (bound) glycolaldehyde are transformed to two glycolaldehyde molecules. However note, that the rule is general enough such that any pair of molecules with aldehyde groups can be used, i.e., the

inferred reaction pattern refers to a class of overall reactions and the product does not necessarily need to be glycolaldehyde.

The practical computation of these compositions takes less than a second in the current implementation. Even for substantially more general composition sequences the running time remains manageable. For instance, it takes less than 1 minute to compute all composition sequences with a length $k \leq 10$ of the form $p_{i_1} \circ p_{i_2} \circ \cdots \circ p_{i_k} \circ g_q$ with $i_j \in \{0, 1, 2, 3\}$, based on the binding of one of the influx molecules g_0 or g_1 . This results in 1875 different inferred composite rules.

Polymerization can also be viewed as a pathway in a chemical reaction network, albeit one of potentially infinite size. The same methods applied to the automatic inference of the overall reaction pattern of the Formose cycle can be directly applied to detecting composition rules for polymerization reactions. Importantly, even if a chemical reaction network is not given, the approaches presented in this paper can be used to automatically find sequences of reactions that will lead to polymerization. This can be realized by a straight-forward post-processing step: all that needs to be done is to check whether an inferred composite rule exhibits a replicated functional unit. Such polymerization meta-rules also enable the analvsis of chemical systems with highly complex carbon skeletons such as the natural compound classes of the terpenes or the polyketides.

Conclusions

Graph grammars provide a convenient framework for modeling chemistries on different levels of abstraction. A chemically valid approach is to see any chemical reaction as a bi-molecular reaction. This requires graph grammar rules that cover changes of molecules in an rather explicit and detailed way. Understanding chemical reaction patterns usually requires spanning the chemical reaction networks based on such rules. Obviously, this approach suffers the inherent potential of an immense combinatorial explosion. In this paper we introduced the automatic inference of such higher-level chemical reaction pattern based on a formal approach for graph grammar rule combination. We analyzed the autocatalytic cycle of the Formose reaction and inferred its overall reaction pattern as a rule composition of nine rules. Rule composition is also naturally applicable to inferring patterns of polymerization reactions. Future work will include e.g. the analysis of terpene-based and hydrogen cyanide-based polymerization chemistry. Many of the enzyme reactions collected in metabolism databases such as KEGG [16] or MetaCyc [17] are in fact overall reactions of multistep mechanisms. The enzyme D-alanine transaminase (EC 2.6.1.21), for instance, achieves its chemical transformation in 12 elementary steps. Generating chemically correct atom-mappings of such overall reactions, a very important step e.g. in the interpretation of isotope tracer experiments, is infeasible with the currently available methods. In contrast a composite rule constructed from the individual enzymatic steps, as found for instance in the MACiE database [18], is guaranteed to yield the chemically correct atom mapping for the overal enzyme reaction.

Authors contributions

J.L.A. implemented the rule composition system. All authors contributed to the theory, the writing of the manuscript and approved the submitted manuscript.

Acknowledgments

This work was supported in part by the Volkswagen Stiftung proj. no. I/82719, the COST-Action CM0703 "Systems Chemistry", and the Danish Council for Independent Research, Natural Sciences.

References

- Klamt S, Haus UU, Theis F: Hypergraphs and Cellular Networks. PLoS Comput Biol 2009, 5(5):e1000385.
- Dittrich P, Ziegler J, Banzhaf W: Artificial chemistries
 a review. Artificial life 2001, 7(3):225–275.
- Benkö G, Flamm C, Stadler PF: A graph-based toy model of chemistry. J Chem Inf Comput Sci 2003, 43(4):1085–1093.
- Aittokallio T, Schwikowski B: Graph-based methods for analysing networks in cell biology. Brief Bioinform 2006, 7(3):243-255.
- Orth JD, Thiele I, Palsson BØ: What is flux balance analysis? Nature Biotech. 2010, 28:245–248.
- Schuster S, Fell DA, Dandekar T: A general definition of metabolic pathways useful for systematic organization and analysis of complex metabolic networks. *Nature Biotech.* 2000, 18:326–332.

- Price ND, Reed JL, Papin JA, Wiback SJ, Palsson BØ: Network-based analysis of metabolic regulation in the human red blood cell. J. Theor. Biol. 2003, 225:185–194.
- Larhlimi A, Bockmayr A: A new constraintbased description of the steady-state flux cone of metabolic networks. *Discr. Appl. Math.* 2009, 157:2257–2266.
- Félix L, Rosselló F, Valiente G: Efficient Reconstruction of Metabolic Pathways by Bidirectional Chemical Search. Bull. Math. Biol. 2009, 71:750–769.
- Ehrig H, Ehrig K, Prange U, Taenthzer G: Fundamentals of Algebraic Graph Transformation. Berlin, D: Springer-Verlag 2006.
- Rosselló F, Valiente G: Graph grammars and molecular biology. Lect. Notes Comp. Sci. 2005, 3393:116– 133.
- Danos V, Feret J, Fontana W, Harmer R, Hayman J, Krivine J, Thompson-Walsh C, Winskel G: Graphs, Rewriting and Pathway Reconstruction for Rule-Based Models. Leibniz International Proceedings in Informatics 2012. in press.
- Golas U: Analysis and Correctness of Algebraic Graph and Model Transformations. Wiesbaden, D: Vieweg+Teubner 2010.
- Ehrig H, Habel A, Kreowski HJ, Parisi-Presicce F: Parallelism and Concurrency in High-Level Replacement Systems. Math. Struct. Comp. Science 1991, 1:361–404.
- Flamm C, Ullrich A, Ekker H, Mann M, Hogerl D, Rohrschneider M, Sauer S, Scheuermann G, Klemm K, Hofacker I, Stadler PF: Evolution of metabolic networks: A computational frame-work. *Journal of* Systems Chemistry 2010, 1(4).
- Kanehisa M, Goto S, Sato Y, Furumichi M, Tanabe M: KEGG for integration and interpretation of largescale molecular data sets. *Nucleic Acids Res.* 2012, 40:D109–D114.
- 17. Caspi R, Altman T, Dreher K, Fulcher CA, Subhraveti P, Keseler IM, Kothari A, Krummenacker M, Latendresse M, Mueller LA, Ong Q, Paley S, Pujar A, Shearer AG, Travers M, Weerasinghe D, Zhang P, Karp PD: The MetaCyc database of metabolic pathways and enzymes and the BioCyc collection of pathway/genome databases. Nucleic Acids Res. 2012, 40:D742-D753.
- Holliday GL, Andreini C, Fischer JD, Rahman SA, Almonacid DE, Williams ST, Pearson WR: MACiE: exploring the diversity of biochemical reactions. *Nucleic Acids Research* 2012, 40:D783–D789,.
- Ehrig H, Heckel R, Korff M, Löwe M, Ribeiro L, Wagner A, Corradini A: Algebraic Approaches to Graph Transformation Part II: Single Pushout Approach and Comparison with Double Pushout Approach. In Handbook of Graph Grammars and Computing by Graph Transformations. Edited by Rozenberg G, Singapore: World Scientific 1997:247-312.
- Löwe M: Algebraic approach to single-pushout graph transformation. Theor. Comp. Sci. 1993, 109:181–224.

 H E: Introduction to the Algebraic Theory of Graph Grammars. Lect. Notes Comp. Sci. 1979, 13:169.

Appendix A: SPO and DPO in Artificial Chemistry Models

In the double pushout (DPO) framework a production $L \stackrel{l}{\leftarrow} K \stackrel{r}{\rightarrow} R$ is defined by three graphs (the left graph L, the right graph R, and the context graph K) and two graph morphisms $l: K \to L$ and $r: K \to R$. In the single pushout (SPO) formalism, $L \stackrel{p}{\rightarrow} R$ is specified by only two graphs L and R and a partial graph morphism p. In other words, the production is specified by a total graph morphism $dom(p) \to R$, where dom(p) is a subgraph of L. While DPO lives on the category of graphs and graph morphisms, SPO is built upon the category of graphs and partial graph morphisms, see [19] for a detailed comparison of the two approaches.

There is no real difference between the rules themselves since DPO and SPO productions can be translated into each other: $L \xrightarrow{p} R$ translates to $L \xleftarrow{l} dom(p) \xrightarrow{r} R$ where l is the inclusion of dom(p)in L and r is the domain restriction of p to dom(p). Conversely, $L \xleftarrow{l} K \xrightarrow{r} R$ translates to $L \xrightarrow{p} R$ where the partial morphism $p = r \circ l^{-1}$ is well-defined provided l is injective and dom(p) = l(K) [20].

The main difference between the two frameworks lies in the application of the productions, i.e., in the resulting (direct) graph derivations. SPO derivations are complete, that is, for each production $L \xrightarrow{p} R$ and each match $m : L \xrightarrow{m} G$ there is a derivation $G \xrightarrow{p,m} H$. In contrast, in DPO the corresponding derivation exists if and only if the *gluing condition* is satisfied:

- (I) There are not distinct elements x, y of L with m(x) = m(y) and $y \notin l(K)$.
- (D) No edge e of $G \setminus m(L)$ is incident to a node in $m(L \setminus l(K))$.

SPO is therefore more powerful than DPO in the sense that more general transformations can be implemented. On the other hand, the gluing condition ensures that there are no "side effects" such as dangling edges (which have to be eliminated by construction in SPO). These side effects make SPO transformation more difficult to understand and control in practical applications. A useful feature of DPO derivations is that productions are invertible [21]. As in chemical reactions, it suffices to exchange the roles of the left and the right graph, i.e., of products and educts. In contrast, the more general SPO derivations are not invertible in general.

In applications of graph transformation systems to modelling chemical reactions further restrictions are needed to account for the peculiarities of chemical transformation systems:

- (i) Reaction rules specify subgraphs. Therefore, the matching morphisms m and n are injective.
- (ii) Since atoms are conserved in chemical reactions, the restrictions l_V and r_V to the vertex set are bijections and determine the atom mapping. Edges model chemical bonds specified by electron pairs. These can only be moved around in the molecules but not collapsed onto the same bond with the same type. The morphisms l and r therefore must be injective.

Lemma 1. Conditions (i) and (ii) imply the gluing condition.

Proof. Since m is injective, i.e., m(x) = m(y) implies x = y, condition (I) is satisfied. Condition (ii) implies that $l_V(K)$ is the vertex set of L, i.e., $L \setminus l(K)$ contains no vertices. Thus $m(L \setminus l(K))$ is empty so that condition (D) is trivially satisfied. \Box

As far as models of chemistry are concerned, therefore, SPO and DPO graph transformations are equivalent.

We prefer to work with the DPO framework for several reasons. First, the explicit exposure of the context graph K provides a convenient starting point for considering transition states e.g. in terms of the pair of subgraphs $(L \setminus l(K), R \setminus r(K))$. In addition, in our experience it is helpful to explicitly construct K in the process of designing rule sets of particular types of chemistry such as Diels Alder reactions or aldol condensations appearing in this contribution. Maybe more importantly, the DPO framework appears more convenient when building analysis tools such as coarse graining operations into the rewriting system. Finally, the framework of (injective) graph morphisms, in our view, is a more convenient basis for mathematical investigations than the partial graph morphisms on which SPO is built.

Appendix B: Example of Enumeration of Compositions

In this Appendix we show the complete result of the composition of two (artificial) rules, p_1 and p_2 , including the selection of submatches from the match matrix. The two rules are depicted in Fig. 7 with the extended match matrix of the composition $p_2 \circ p_1$, that corresponds to the example of an extended match matrix as given in the paper. The rules in this section are all depicted with vertices that have an additional index. The numbering of the components is in increasing order wrt. to these indices, e.g., L_2^1 denotes the component connecting nodes A, 0 and B, 1 and L_2^2 denotes the component connecting nodes B, 2 and C, 3.



Fig. 7: The two rules p_1 and p_2 , and the extended match matrix of the composition $p_2 \circ p_1$. The components of both R_1 and L_2 are numbered in the same order as the vertex indices.

In the following we will enumerate all valid selections of submatches based on the extended match matrix and give the corresponding resulting rule composition. The chosen matches are depicted as • in the extended match matrix. If several matches can be found (in our example this is true for the component L_2^1 , that can be matched twice in R_1^3), the • has an index.

Composition 1



Result of composition 1

	R_1^1	R_1^2	R_1^3	R_1^{\emptyset}
L_2^1			\bullet_1	
L_{2}^{2}		•		



Result of composition 2

Composition 3





Result of composition 3

Composition 4





Result of composition 4

	R_1^1	R_1^2	R_1^3	R_1^{\emptyset}
L_2^1	•			
L_{2}^{2}			•	



Result of composition 5

Composition 6





Result of composition 6

Invalid Selection

	R_1^1	R_1^2	R_{1}^{3}	R_1^{\emptyset}
L_2^1			\bullet_2	
L_2^2			•	

This selection of submatches is invalid, as they are not disjoint (node B, 9 would be matched twice).



Result of composition $7\,$

Composition 8



Result of composition 8

Composition 9





Result of composition 9





Result of composition 10