

A Software Package for Chemically Inspired Graph Transformation

Jakob L. Andersen^{1,9}(✉), Christoph Flamm^{2,8},
Daniel Merkle¹(✉), and Peter F. Stadler^{2,3,4,5,6,7}

¹ Department of Mathematics and Computer Science,
University of Southern Denmark, 5230 Odense, Denmark
{jlandersen,daniel}@imada.sdu.dk

² Institute for Theoretical Chemistry, University of Vienna, 1090 Wien, Austria
xtof@tbi.univie.ac.at

³ Bioinformatics Group, Department of Computer Science,
and Interdisciplinary Center for Bioinformatics,
University of Leipzig, 04107 Leipzig, Germany

⁴ Max Planck Institute for Mathematics in the Sciences, 04103 Leipzig, Germany

⁵ Fraunhofer Institute for Cell Therapy and Immunology, 04103 Leipzig, Germany

⁶ Center for Non-coding RNA in Technology and Health,
University of Copenhagen, 1870 Frederiksberg, Denmark
stadler@bioinf.uni-leipzig.de

⁷ Santa Fe Institute, 1399 Hyde Park Rd, Santa Fe, NM 87501, USA

⁸ Research Network Chemistry Meets Microbiology,
University of Vienna, 1090 Wien, Austria

⁹ Earth-Life Science Institute, Tokyo Institute of Technology, Tokyo 152-8550, Japan
jlandersen@elsi.jp

Abstract. Chemical reaction networks can be automatically generated from graph grammar descriptions, where transformation rules model reaction patterns. Because a molecule graph is connected and reactions in general involve multiple molecules, the transformation must be performed on multisets of graphs. We present a general software package for this type of graph transformation system, which can be used for modelling chemical systems. The package contains a C++ library with algorithms for working with transformation rules in the Double Pushout formalism, e.g., composition of rules and a domain specific language for programming graph language generation. A Python interface makes these features easily accessible. The package also has extensive procedures for automatically visualising not only graphs and transformation rules, but also Double Pushout diagrams and graph languages in form of directed hypergraphs. The software is available as an open source package, and interactive examples can be found on the accompanying webpage.

Keywords: Double Pushout · Chemical graph transformation system · Graph grammar · Rule composition · Strategy framework

1 Introduction

It has been common practice in chemistry for more than a century to represent molecules as labelled graphs, with vertices representing atoms, and edges representing the chemical bonds between them [23]. It is natural, therefore, to formalize chemical reactions as graph transformations [6, 11, 13, 20]. Many computational tools for graph transformation have been developed; some of them are either specific to chemistry [21] or at least provide special features for chemical systems [18]. General graph transformation tools, such as AGG [24], have also been used to modelling chemical systems [11].

Chemical graph transformation, however, differs in one crucial aspect from the usual setup in the graph transformation literature, where a single (usually connected) graph is rewritten, thus yielding a graph language. Chemical reactions in general involve multiple molecules. Chemical graph transformations therefore operate on *multisets* of graphs to produce a chemical “space” or “universe”. A similar viewpoint was presented in [17], but here we let the basic graphs remain connected, and multisets of them are therefore dynamically constructed and taken apart in direct derivations.

Graph languages can be infinite. This is of course also true for chemical universes (which in general contain classical graph languages as subsets). In the case of chemistry, the best known infinite universes comprise polymers. The combinatorics of graphs makes it impossible in most cases to explore graph languages or chemical universes by means of a simple breadth-first search. This limitation can be overcome at least in part with the help of strategy languages that guide the rule applications. One such language has been developed for rewriting port graphs [12], implemented in the PORGY tool [5]. We have in previous work presented a similar strategy language [4] for transformation of multisets of graphs, which is based on partial application of transformation rules [2].

Here, we present the first part of the software package MedØIDatschgerl (in short: MØD), that contains a chemically inspired graph transformation system, based on the Double Pushout formalism [10]. It includes generic algorithms for composing transformation rules [2]. This feature can be used, e.g., to abstract reaction mechanisms, or whole pathways, into overall rules [3]. MØD also implements the strategy language [4] mentioned above. It facilitates the efficient generation of vast reaction networks under global constraints on the system. The underlying transformation system is not constrained to chemical systems. The package contains specialized functionalities for applications in chemistry, such as the capability to load graphs from SMILES strings [26]. This first version of MØD thus provides the main features of a chemical graph transformation system as described in [27].

The core of the package is a C++ 11 library that in turn makes use of the Boost Graph Library [22] to implement standard graph algorithms. Easy access to the library is provided by means of extensive Python bindings. In the following we use these to demonstrate the functionality of the package. The Python module provides additional features, such as embedded domain-specific languages for rule composition, and for exploration strategies. The package also provides

comprehensive functionality for automatically visualising graphs, rules, Double Pushout diagrams, and hypergraphs of graph derivations, i.e., reaction networks. A \LaTeX package is additionally included that provides an easy mechanism for including visualisations directly in documents.

In Sect. 2 we first describe formal background for transforming multisets of graphs. Section 3 gives examples of how graph and rule objects can be used, e.g., to find morphisms with the help of the VF2 algorithms [8,9]. Sections 4 and 5 describes the interfaces for respectively rule composition and the strategy language. Section 6, finally, gives examples of the customisable figure generation functionality of the package, including the \LaTeX package.

The source code of `MedØIDatschgerl` as well as additional usage examples can be found at <http://mod.imada.sdu.dk>. A live version of the software can be accessed at <http://mod.imada.sdu.dk/playground.html>. This site also provides access to the large collection of examples.

2 Transformation of Multisets of Graphs

The graph transformation formalism we use is a variant of the Double Pushout (DPO) approach (e.g., see [10] more details). Given a category of graphs \mathcal{C} , a DPO rule is defined as a span $p = (L \xleftarrow{l} K \xrightarrow{r} R)$, where we call the graphs L , K , and R respectively the *left side*, *context*, and *right side* of the rule. A rule can be applied to a graph G using a match morphism $m: L \rightarrow G$ when the *dangling condition* and the *identification condition* are satisfied [10]. This results in a new graph H , where the copy of L has been replaced with a copy of R . We write such a direct derivation as $G \xrightarrow{p,m} H$, or simply as $G \xrightarrow{R} H$ or $G \Rightarrow H$ when the match or rule is unimportant. The graph transformation thus works in a category \mathcal{C} of possibly disconnected graphs.

Let \mathcal{C}' be the subcategory of \mathcal{C} restricted to connected graphs. A graph $G \in \mathcal{C}$ will be identified with the multiset of its connected components. We use double curly brackets $\{\{\dots\}\}$ to denote the construction of multisets. Hence we write $G = \{\{g_1, g_2, \dots, g_k\}\}$ for an arbitrary graph $G \in \mathcal{C}$ with not necessarily distinct connected components $g_i \in \mathcal{C}'$. For a set $\mathcal{G} \subseteq \mathcal{C}'$ of connected graphs and a graph $G = \{\{g_1, g_2, \dots, g_k\}\} \in \mathcal{C}$ we write $G \in^* \mathcal{G}$ whenever $g_i \in \mathcal{G}$ for all $i = 1, \dots, k$.

We define a graph grammar $\Gamma(\mathcal{G}, \mathcal{P})$ by a set of connected starting graphs $\mathcal{G} \subseteq \mathcal{C}'$, and a set of DPO rules \mathcal{P} based on the category \mathcal{C} . The language of the grammar $L(\mathcal{G}, \mathcal{P})$ includes the starting graphs \mathcal{G} . Additional graphs in the language are constructed by iteratively finding direct derivations $G \xrightarrow{R} H$ with $p \in \mathcal{P}$ and $G, H \in \mathcal{C}$ such that $G \in^* L(\mathcal{G}, \mathcal{P})$. Each graph $h \in H$ is then defined to be in the language as well. A concise constructive definition of the language is thus $L(\mathcal{G}, \mathcal{P}) = \bigcup_{k=1}^{\infty} \mathcal{G}_k$ with $\mathcal{G}_1 = \mathcal{G}$ and

$$\mathcal{G}_{k+1} = \mathcal{G}_k \cup \bigcup_{p \in \mathcal{P}} \{h \in H \mid \exists G \in^* \mathcal{G}_k : G \xrightarrow{R} H\}$$

In `MØD` the objects of the category \mathcal{C} are all undirected graphs without parallel edges and loops, and labelled on vertices and edges with text strings.

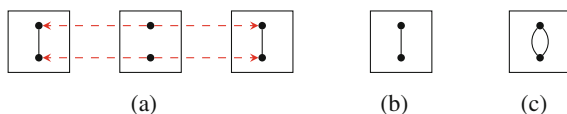


Fig. 1. The pushout object of (a) in the category of simple graphs is either not existing or is the graph depicted in (b) where the two edges are merged. For multigraphs the pushout object would be the graph depicted in (c).

The core algorithms can however be specialised for other label types. We also restrict the class of morphisms in \mathcal{C} to be injective, i.e., they are restricted to graph monomorphisms. Note that this restriction implies that the identification condition of rule application is always fulfilled.

The choice of disallowing parallel edges is motivated by the aim of modelling of chemistry, where bonds between atoms are single entities. While a “double bond” consists of twice the amount of electrons than a “single bond”, it does not in general behave as two single bonds. However, when parallel edges are disallowed a special situation arises when constructing pushouts. Consider the span in Fig. 1a. If parallel edges are allowed, the pushout object is the one shown in Fig. 1c. Without parallel edges we could identify the edges as shown in Fig. 1b. This approach was used in for example [7]. However, for chemistry this means that we must define how to add two bonds together, which is not meaningful. We therefore simply define that no pushout object exists for the span. A direct derivation with the Double Pushout approach thus additionally requires that the second pushout is defined.

The explicit use of multisets gives rise to a form of minimality of a derivation. If $\{\{g_a, g_b, g_b\} \xrightarrow{p,m} \{h_c, h_d\}\}$ is a valid derivation, for some rule p and match m , then the extended derivation $\{\{g_a, g_b, g_b, q\} \xrightarrow{p,m} \{h_c, h_d, q\}\}$ is also valid, even though q is not “used”. We therefore say that a derivation $G \xrightarrow{p,m} H$ with the left-hand side $G = \{g_1, g_2, \dots, g_n\}$ is *proper* if and only if

$$g_i \cap \text{img}(m) \neq \emptyset, \forall 1 \leq i \leq n$$

That is, if all connected components of G are hit by the match. The algorithms in MØD only enumerate proper derivations.

3 Graphs and Rules

Graphs and rules are available as classes in the library. A rule $(L \xleftarrow{l} K \xrightarrow{r} R)$ can be loaded from a description in GML [16] format. As both l and r are monomorphisms the rule is represented without redundant information in GML by three sets corresponding somewhat to the graph fragments $L \setminus K$, K , and $R \setminus K$ (see Fig. 2 for details).

Graphs can similarly be loaded from GML descriptions, and molecule graphs can also be loaded using the SMILES format [26] where most hydrogen atoms

are implicitly specified. A SMILES string is a pre-order recording of a depth-first traversal of the connected graph, where back-edges are replaced with pairs of integers.

Both input methods result in objects which internally stores the graph structure, where all labels are text strings. Figure 2 shows examples of graph and rule loading, using the Python interface of the software.

```

formaldehyde = graphGML("formaldehyde.gml")
caffeine = smiles("Cn1cnc2c1c(=O)n(c(=O)n2C)C")
ketoEnol = ruleGMLString("""rule [
  left [
    edge [ source 1 target 4 label "-" ]
    edge [ source 1 target 2 label "-" ]
    edge [ source 2 target 3 label "=" ]
    node [ id 3 label "O" ]
    node [ id 4 label "H" ]
  ]
  context [
    node [ id 1 label "C" ]
    node [ id 2 label "C" ]
  ]
  right [
    edge [ source 1 target 2 label "=" ]
    edge [ source 2 target 3 label "-" ]
    node [ id 3 label "O-" ]
    node [ id 4 label "H+" ]
  ]
] """)

```

Fig. 2. Creation of two graph objects and a transformation rule object in the Python interface. The (molecule) graph ‘formaldehyde’ is loaded from an external GML file, while the (molecule) graph ‘caffeine’ is loaded from a SMILES string [26], often used in cheminformatics. General labelled graphs can only be loaded from a GML description, and all graphs are internally stored simply as labelled adjacency lists. The DPO transformation rule ‘ketoEnol’ is loaded from an inline GML description. When the GML sections ‘left’, ‘context’, and ‘right’ are considered sets, they encode a rule ($L \leftarrow K \rightarrow R$) with $L = \text{‘left’} \cup \text{‘context’}$, $R = \text{‘right’} \cup \text{‘context’}$, and $K = \text{‘context’} \cup (\text{‘left’} \cap \text{‘right’})$. Vertices and edges that change label are thus specified in both ‘left’ and ‘right’. Note that in GML the endpoints of edges are described by ‘source’ and ‘target’, but for undirected graphs these tags have no particular meaning and may be exchanged. The graphs and rules are visualised in Fig. 5.

Graphs have methods for counting both monomorphisms and isomorphisms, e.g., for substructure search and for finding duplicate graphs. Counting the number of carbonyl groups in a molecule ‘mol’ can be done simply as

```

carbonyl = smiles("[C]=O") count = carbonyl.monomorphism(mol,
maxNumMatches=1337)

```

By default the ‘monomorphism’ method stops searching after the first morphism is found; alternative matches can be retrieved by setting the limit to a higher value.

Rule objects also have methods for counting monomorphisms and isomorphisms. A rule morphism $m: p_1 \rightarrow p_2$ on the rules $p_i = (L_i \xleftarrow{l_i} K_i \xrightarrow{r_i} R_i), i = 1, 2$ is a 3-tuple of graph morphisms $m_X: X_1 \rightarrow X_2, X \in \{L, K, R\}$ such that they commute with the morphisms in the rules. Finding an isomorphism between two rules can thus be used for detecting duplicate rules, while finding a monomorphism $m: p_1 \rightarrow p_2$ determines that p_1 is at least as general as p_2 .

4 Composition of Transformation Rules

In [2,3] the concept of rule composition is described, where two rules $p_1 = (L_1 \leftarrow K_1 \rightarrow R_1), p_2 = (L_2 \leftarrow K_2 \rightarrow R_2)$ are composed along a common subgraph given by the span $R_1 \leftarrow D \rightarrow L_2$. Different types of rule composition can be defined by restricting the common subgraph and its relation to the two rules. MØD implements enumeration algorithms for several special cases that are motivated and defined in [2,3]. The simplest case is to set D as the empty graph, denoted by the operator \bullet_\emptyset , to create a composed rule that implements the parallel application of two rules. In the most general case, denoted by \bullet_\cap , all common subgraphs of R_1 and L_2 are enumerated. In a more restricted setting R_1 is a subgraph of L_2 , denoted by \bullet_{\subseteq} , or, symmetrically, L_2 is a subgraph of R_1 , denoted by \bullet_{\supseteq} . When the subgraph requirement is relaxed to only hold for a subset of the connected components of the graphs we denoted it by \bullet_{\subseteq}^c and \bullet_{\supseteq}^c .

The Python interface contains a mini-language for computing the result of rule composition expressions with these operators. The grammar for this language of expressions is shown in Fig. 3.

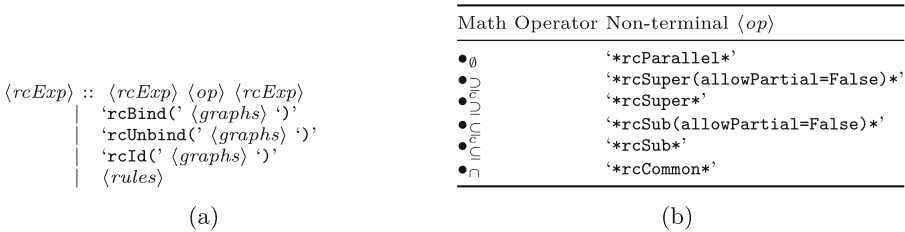


Fig. 3. Grammar for rule composition expressions in the Python interface, where $\langle graphs \rangle$ is a Python expression returning either a single graph or a collection of graphs. Similarly is $\langle rules \rangle$ a Python expression returning either a single rule or a collection of rules. The pseudo-operators $\langle op \rangle$ each correspond to a mathematical rule composition operator (see [2,3]). The three functions 'rcBind', 'rcUnbind', and 'rcId' refers to the construction of the respective rules $(\emptyset \leftarrow \emptyset \rightarrow G)$, $(G \leftarrow \emptyset \rightarrow \emptyset)$, and $(G \leftarrow G \rightarrow G)$ from a graph G .

Its implementation is realised using a series of global objects with suitable overloading of the multiplication operator. A rule composition expression can be passed to an evaluator, which will carry out the composition and discard

duplicate results, as determined by checking isomorphism between rules. The result of each $\langle rcExp \rangle$ is coerced into a list of rules, and the operators consider all selections of rules from their arguments. That is, if ‘P1’ and ‘P2’ are two rule composition expressions, whose evaluation results in two corresponding lists of rules, P_1 and P_2 . Then, for example, the evaluation of ‘P1 *rcParallel* P2’ results in the following list of rules:

$$\bigcup_{p_1 \in P_1} \bigcup_{p_2 \in P_2} p_1 \bullet_{\emptyset} p_2$$

Each of these rules encodes the parallel application of a rule from P_1 and a rule from P_2 .

In the following Python code, for example, we compute the rules corresponding to the bottom span ($G \leftarrow D \rightarrow H$) of a DPO diagram, arising from applying the rule $p = (L \leftarrow K \rightarrow R)$ to the multiset of connected graphs $G = \{\{g_1, g_2\}\}$.

```
exp = rcId(g1) *rcParallel* rcId(g2) *rcSuper(allowPartial=False)* p
rc = rcEvaluator(ruleList) res = rc.eval(exp)
```

Here, the rule composition evaluator is given a list ‘ruleList’ of known rules that will be used for detecting isomorphic rules. Larger rule composition expressions, such as those found in [3], can similarly be directly written as Python code.

5 Exploration of Graph Languages Using Strategies

A breadth-first enumeration of the language of a graph grammar is not always desirable. For example, in chemical systems there are often constraints that can not be expressed easily in the underlying graph transformation rules. In [4] a strategy framework is introduced for the exploration of graph languages. It is a domain specific programming language that, like the rule composition expressions, is implemented in the Python interface, with the grammar shown in Fig. 4. The language computes on sets of graphs. Simplified, this means that each execution state is a set of connected graphs. An *addition strategy* adds further graphs to this state, and a *filter strategy* removes graphs from it. A *rule strategy* enumerates direct derivations based on the state, subject to acceptance by filters introduced by the *left-* and *right-predicate strategies*. Newly derived graphs are added to the state. Strategies can be sequentially composed with the ‘>>’ operator, which can be extended to k -fold composition with the *repetition strategy*. A *parallel strategy* executes multiple strategies with the same input, and merges their output. During the execution of a program the discovered direct derivations are recorded as an annotated directed multi-hypergraph, which for chemical systems is a reaction network. For a full definition of the language see [4] or the MØD documentation.

A strategy expression must, similarly to a rule composition expression, be given to an evaluator which ensures that isomorphic graphs are represented by the same C++/Python object. After execution the evaluator contains the

```

⟨strat⟩ :: ⟨strats⟩ | ⟨strat⟩ '⟨>' ⟨strat⟩ | ⟨rule⟩
| 'addSubset(' ⟨graphs⟩ ') ' | 'addUniverse(' ⟨graphs⟩ ') '
| 'filterSubset(' ⟨filterPred⟩ ') ' | 'filterUniverse(' ⟨filterPred⟩ ') '
| 'leftPredicate[' ⟨derivationPred⟩ ']' (' ⟨strat⟩ ') '
| 'rightPredicate[' ⟨derivationPred⟩ ']' (' ⟨strat⟩ ') '
| 'repeat' [ '[' ⟨int⟩ ']' ] (' ⟨strat⟩ ') '
| 'revive(' ⟨strat⟩ ') '

```

Fig. 4. Grammar for the domain specific language for guiding graph transformation, embedded in the Python interface of the software package. The non-terminal $\langle strats \rangle$ must be a collection of strategies, that becomes a **parallel** strategy from [4]. The production $\langle strat \rangle \langle \rangle \langle strat \rangle$ results in a sequence strategy.

generated derivation graph, which can be visualised or programmatically used for subsequent analysis.

The strategy language can for example be used for the simple breadth-first exploration of a grammar with a set of graphs `startingGraphs` and a set of rules `ruleSet`, where exploration does not result in graphs above a certain size (42 vertices):

```

strat = (
    addSubset(startingGraphs)
    >> rightPredicate[
        lambda derivation: all(g.numVertices <= 42 for g in derivation.right)
    ](
        repeat(ruleSet)
    )
)
dg = dgRuleComp(startingGraphs, strat)
dg.calc()

```

The `dg` object is the evaluator which afterwards contains the derivation graph. More examples can be found in [1, 4] where complex chemical behaviour is incorporated into strategies. An abstract example can also be found in [4] where the puzzle game Catalan [14] is solved using exploration strategies.

6 Figure Generation

The software package includes elaborate functionality for automatically visualising graph, rules, derivation graphs, and derivations. The final rendering of figures is done using the TikZ [25] package for L^AT_EX, while the layouts for graphs are computed using Graphviz [15]. However, for molecule graphs it is possible to use the cheminformatics library Open Babel [19] for laying out molecules and reaction patterns in a more chemically familiar manner.

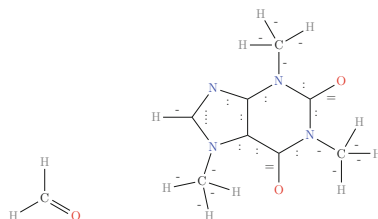
Visualisation starts by calling a `print` method on the object in question. This generates files with L^AT_EX code and a graph description in Graphviz format. Special post-processing commands are additionally inserted into another file. Invoking the post-processor will then generate coordinates and compile the final layout. In addition, an aggregate summary document is compiled that includes all figures for easy overview. Figure 5 shows an example, where the wrapper script `mod` provided by the package is used to automatically execute both a Python


```

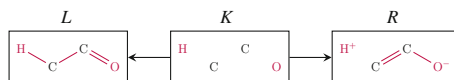
p = GraphPrinter()
p.setMolDefault()
p.collapseHydrogens = False
formaldehyde.print(p)
p.edgesAsBonds = False
caffeine.print(p)
p.setReactionDefault()
ketoEnol.print(p)

```

(a) Additional Python code to Fig. 2, for generating figures.



(b) Automatically compiled figure of the two graphs loaded in Fig. 2.



(c) Automatically compiled figure of the DPO rule loaded in Fig. 2.

Fig. 5. Example of automatic visualisation of graphs and rules, using the post-processor. The Python code is an extension of the code from Fig. 2, and can be executed using the provided ‘mod’ script that invokes both the Python interpreter ‘python3’ and the post-processor, ‘mod_post’. Edges with special labels are as default rendered in a special chemical manner, as illustrated with the left graph of (b) (formaldehyde). In the right graph of (b) (caffeine) the edge labels are shown explicitly. Both graphs uses chemical colouring. The colouring of the transformation rule, (c) denote the differences between L , K , and R . (Color figure online)

script and subsequently the post-processor. The example also shows part of the functionality for chemical rendering options, such as atom-specific colouring, charges rendered in superscript, and collapsing of hydrogen vertices into their neighbours.

Derivation graphs can also be visualised automatically, where each vertex is depicted with a rendering of the graph it represents. The overall depiction can be customised to a high degree, e.g., by annotation or colouring of vertices and hyperedges using user-defined callback functions. Figure 6 illustrates part of this functionality.

Individual derivations of a derivation graph can be visualised in form of Double Pushout diagrams. The rendering of these diagrams can be customised similar to how rules and graph depictions can, e.g., to make the graphs have a more chemical feel. An example of derivation printing is illustrated in Fig. 7.

Composition of transformation rules is a core operation in the software, and for better understanding the operation we provide a mechanism for visualising individual compositions. An example of such a visualisation is shown in Fig. 8, where only the left and right graphs of two argument rules and the result rule are shown. The composition relation is shown as red dashed lines between the left graph of the first rule and the right graph of the second rule.

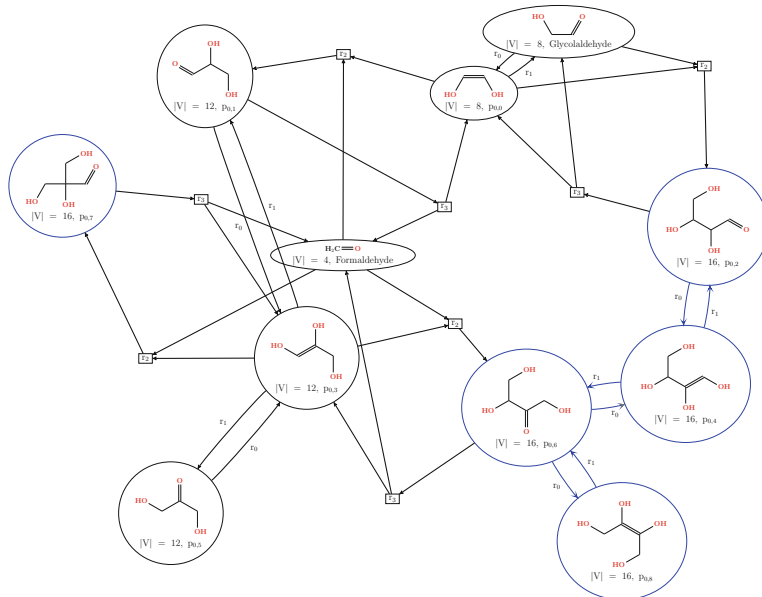
Including Figures in \LaTeX Documents. To make it easier to use illustrations of graphs and rules we have included a \LaTeX package in the software. It provides macros for automatically generating Python scripts that subsequently

```

p = DGPrinter()
p.pushVertexLabel(lambda g, dg: "|V| = %d" % g.numVertices)
p.pushVertexColour(lambda g, dg: "blue" if g.numVertices >= 16 else "")
dg.print(p)

```

(a) Python code for customised visualisation of a derivation graph ‘dg’.



(b) Example of automatically laid out and rendered derivation graph with custom labelling and colour.

Fig. 6. Example of derivation graph printing. Each vertex is as default labelled with the name of the graph it represents, and a figure of the graph is embedded. Each hyperedge is as default labelled with the name of the rule used in the derivation the hyperedge represents. A general hyperedge is represented by a box, but for hyperedges with only 1 head and 1 tail the box is omitted, and a single labelled arc is rendered. (Color figure online)

generate figures and \LaTeX code for inclusion into the original document. For example, the depictions in Fig. 5 are inserted with the following code.

```

\graphGML[collapse hydrogens=false][scale=0.4]{formaldehyde.gml}
\smiles[collapse hydrogens=false, edges as bonds=false][scale=0.4]
{Cn1cnc2c1c(=O)n(c(=O)n2C)C}
\ruleGML{ketoEnol.gml}{\dpoRule[scale=0.4]}

```

Each ‘\graphGML’ and ‘\smiles’ macro expands into an ‘\includegraphics’ for a specific PDF file, and a Python script is generated which can be executed to compile the needed files. The ‘\ruleGML’ macro expands into

```

\dpoRule[scale=0.4]{fileL.pdf}{fileK.pdf}{fileR.pdf}

```

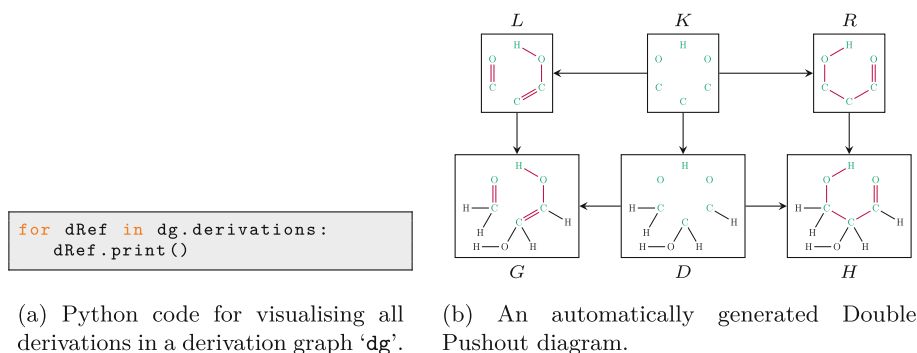


Fig. 7. Example of visualisation of derivations. Each derivation from a derivation graph can be printed, with the same customisation options as for graphs and rules. Additional colouring is used to highlight the image of the rule into the lower span. (Color figure online)

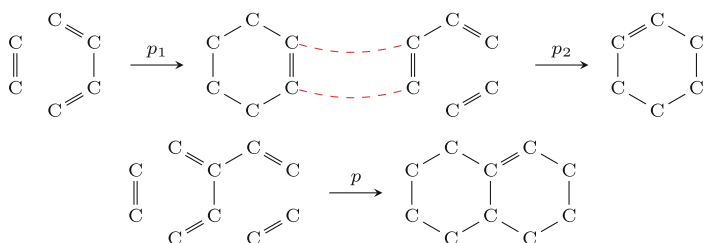


Fig. 8. Visualisation of the composition of two rules $p_i = (L_i \leftarrow K_i \rightarrow R_i, i = 1, 2$, along the a common subgraph of R_1 and L_2 , indicated by the dashed red lines. Only the left and right graphs of both rules, and the resulting rule, are shown. The rendering can be customised in the same manner as the rendering for graphs and rules can. (Color figure online)

where the three PDF files depict the left side, context, and right side of the rule. The `\dpoRule` macro then expands into the final rule diagram with the PDF files included.

7 Summary

MedØDatschgerl is a comprehensive software package for DPO graph transformation on multisets of undirected, labelled graphs. It can be used for generic, abstract graph models. By providing many features for handling chemical data it is particularly well-suited for modelling generative chemical systems. The package includes an elaborate system for automatically producing high-quality visualisations of graphs, rules, and DPO diagrams of direct derivations.

The first public version of MØD described here is intended as the foundation for a larger integrated package for graph-based cheminformatics. Future versions

will for example also include functionalities for pathway analysis in reaction networks produced by the generative transformation methods described here. The graph transformation system, on the other hand, will be extended to cover more complicated chemical properties such as radicals, charges, and stereochemistry.

Acknowledgements. This work is supported by the Danish Council for Independent Research, Natural Sciences, the COST Action CM1304 “Emergence and Evolution of Complex Chemical Systems”, and the ELSI Origins Network (EON), which is supported by a grant from the John Templeton Foundation. The opinions expressed in this publication are those of the authors and do not necessarily reflect the views of the John Templeton Foundation.

A Examples

The following is a short list of examples that show how MedØIDatschgerl can be used via the Python interface. They are all available as modifiable script in the live version of the software, accessible at <http://mod.imada.sdu.dk/playground.html>.

A.1 Graph Interface

Graph objects have a full interface to access individual vertices and edges. The labels of vertices and edges can be accessed both in their raw string form, and as their chemical counterpart (if they have one).

```
g = graphDFS("[R]{x}C([O-])CC=O")
print("|V| =", g.numVertices)
print("|E| =", g.numEdges)
for v in g.vertices:
    print("v%d: label='%s'" % (v.id, v.stringLabel), end="")
    print("\tas molecule: atomId=%d, charge=%d" % (v.atomId, v.charge), end="")
    print("\tis oxygen?", v.atomId == AtomIds.Oxygen)
    print("\td(v) =", v.degree)
    for e in v.incidentEdges: print("\tneighbour:", e.target.id)
for e in g.edges:
    print("(v%d, v%d): label='%s'" % (e.source.id, e.target.id, e.stringLabel), end="")
    try:
        bt = str(e.bondType)
    except LogicError:
        bt = "Invalid"
    print("\tas molecule: bondType=%s" % bt, end="")
    print("\tis double bond?", e.bondType == BondType.Double)
```

A.2 Graph Morphisms

Graph objects have methods for finding morphisms with the VF2 algorithms for isomorphism and monomorphism. We can therefore easily detect isomorphic graphs, count automorphisms, and search for substructures.

```
mol1 = smiles("CC(O)CO")
mol2 = smiles("C(C)CO")
# Check if there is just one isomorphism between the graphs:
isomorphic = mol1.isomorphism(mol2) == 1
print("Isomorphic?", isomorphic)
# Find the number of automorphisms in the graph,
# by explicitly enumerating all of them:
numAutomorphisms = mol1.isomorphism(mol1, maxNumMatches=1337)
print("|Aut(G)| =", numAutomorphisms)
```

```

# Let's count the number of methyl groups:
methyl = smiles("[CH3]")
# The symmetry of the group it self should not be counted,
# so find the size of the automorphism group of methyl.
numAutMethyl = methyl.isomorphism(methyl, maxNumMatches=1337)
print("|Aut(methyl)|", numAutMethyl)
# Now find the number of methyl matches,
numMono = methyl.monomorphism(moll, maxNumMatches=1337)
print("#monomorphisms =", numMono)
# and divide by the symmetries of methyl.
print("#methyl groups =", numMono / numAutMethyl)

```

A.3 Rule Loading

Rules must be specified in GML format.

```

# A rule (L <- K -> R) is specified by three graph fragments:
# left, context, and right
destroyVertex = ruleGMLString('rule [ left [ node [ id 1 label "A" ] ] ]')
createVertex = ruleGMLString('rule [ right [ node [ id 1 label "A" ] ] ]')
identity = ruleGMLString('rule [ context [ node [ id 1 label "A" ] ] ]')
# A vertex/edge can change label:
labelChange = ruleGMLString("""rule [
left [ node [ id 1 label"A" ] edge [ source 1 target 2 label"A" ] ]
# GML can have Python-style line comments too
context [ node [ id 2 label"Q" ] ]
right [ node [ id 1 label"B" ] edge [ source 1 target 2 label"B" ] ]
]""")
# A chemical rule should probably not destroy and create vertices:
ketoEnol = ruleGMLString("""rule [
left [
edge [ source 1 target 4 label"--" ] edge [ source 1 target 2 label"--" ]
edge [ source 2 target 3 label"=" ]
node [ id 3 label"O" ] node [ id 4 label"H" ]
]
context [
node [ id 1 label"C" ] node [ id 2 label"C" ]
]
right [
edge [ source 1 target 2 label"=" ] edge [ source 2 target 3 label"--" ]
node [ id 3 label"O-" ] node [ id 4 label"H+" ]
]
]""")
# Rules can be printed, but label changing edges are not visualised in K:
ketoEnol.print()
# Add with custom options, like graphs:
p1 = GraphPrinter()
p2 = GraphPrinter()
p1.disableAll()
p1.withTexttt = True
p1.withIndex = True
p2.setReactionDefault()
for p in inputRules: p.print(p1, p2)
# Be careful with printing options and non-existing implicit hydrogens:
p1.disableAll()
p1.edgesAsBonds = True
p2.setReactionDefault()
p2.simpleCarbons = True # !!
ketoEnol.print(p1, p2)

```

A.4 Rule Composition 1 — Unary Operators

Special rules can be constructed from graphs.

```

glycolaldehyde.print()
# A graph G can be used to construct special rules:
# (\emptyset <- \emptyset -> G)
bindExp = rcBind(glycolaldehyde)
# (G <- \emptyset -> \emptyset)
unbindExp = rcUnbind(glycolaldehyde)
# (G <- G -> G)
idExp = rcId(glycolaldehyde)
# These are really rule composition expressions that have to be evaluated:
rc = rcEvaluator(inputRules)
# Each expression results in a lists of rules:
bindRules = rc.eval(bindExp)
unbindRules = rc.eval(unbindExp)
idRules = rc.eval(idExp)
postSection("Bind Rules")

```

```

for p in bindRules: p.print()
postSection("Unbind Rules")
for p in unbindRules: p.print()
postSection("Id Rules")
for p in idRules: p.print()

```

A.5 Rule Composition 2 — Parallel Composition

A pair of rules can be merged to a new rule implementing the parallel transformation.

```

rc = rcEvaluator(inputRules)
# The special global object 'rcParallel' is used to make a pseudo-operator:
exp = rcId(formaldehyde) *rcParallel* rcUnbind(glycolaldehyde)
rules = rc.eval(exp)
for p in rules: p.print()

```

A.6 Rule Composition 3 — Supergraph Composition

A pair of rules can (maybe) be composed using a supergraph relation.

```

rc = rcEvaluator(inputRules)
exp = rcId(formaldehyde) *rcParallel* rcId(glycolaldehyde)
exp = exp *rcSuper* ketoEnol_F
rules = rc.eval(exp)
for p in rules: p.print()

```

A.7 Reaction Networks 1 — Rule Application

Transformation rules (reaction patterns) can be applied to graphs (molecules) to create new graphs (molecules). The transformations (reactions) implicitly form a directed (multi-)hypergraph (chemical reaction network).

```

# Reaction networks are expanded using a strategy:
strat = ( # A molecule can be active or passive during evaluation.
  addUniverse(formaldehyde) # passive
  >> addSubset(glycolaldehyde) # active
  # Each reaction must have a least 1 active educt.
  >> inputRules )
# We call a reaction network a 'derivation graph'.
dg = dgRuleComp(inputGraphs, strat)
dg.calc()
# They can also be visualised.
dg.print()

```

A.8 Reaction Networks 2 — Repetition

A sub-strategy can be repeated.

```

strat = ( addUniverse(formaldehyde)
  >> addSubset(glycolaldehyde)
  # Iterate the rule application 4 times.
  >> repeat[4](inputRules) )
dg = dgRuleComp(inputGraphs, strat)
dg.calc()
dg.print()

```

References

1. Andersen, J.L., Andersen, T., Flamm, C., Hanczyc, M.M., Merkle, D., Stadler, P.F.: Navigating the chemical space of HCN polymerization and hydrolysis: guiding graph grammars by mass spectrometry data. *Entropy* **15**(10), 4066–4083 (2013)
2. Andersen, J.L., Flamm, C., Merkle, D., Stadler, P.F.: Inferring chemical reaction patterns using rule composition in graph grammars. *J. Syst. Chem.* **4**(1), 4 (2013)
3. Andersen, J.L., Flamm, C., Merkle, D., Stadler, P.F.: 50 shades of rule composition. In: Fages, F., Piazza, C. (eds.) *FMMB 2014. LNCS*, vol. 8738, pp. 117–135. Springer, Heidelberg (2014)
4. Andersen, J.L., Flamm, C., Merkle, D., Stadler, P.F.: Generic strategies for chemical space exploration. *Int. J. Comput. Biol. Drug Des.* **7**(2/3), 225–258 (2014). TR: <http://arxiv.org/abs/1302.4006>
5. Andrei, O., Fernández, M., Kirchner, H., Melançon, G., Namet, O., Pinaud, B.: PORGY: strategy driven interactive transformation of graphs. In: *Proceedings of the 6th International Workshop on Computing with Terms and Graphs (TERMGRAPH 2011)*. *Electronic Proceedings in Theoretical Computer Science*, vol. 48, pp. 54–68 (2011)
6. Benkő, G., Flamm, C., Stadler, P.F.: A graph-based toy model of chemistry. *J. Chem. Inf. Comput. Sci.* **43**(4), 1085–1093 (2003)
7. Braatz, B., Golas, U., Soboll, T.: How to delete categorically - two pushout complement constructions. *J. Symb. Comput.* **46**(3), 246–271 (2011). *Applied and Computational Category Theory*
8. Cordella, L., Foggia, P., Sansone, C., Vento, M.: A (sub) graph isomorphism algorithm for matching large graphs. *IEEE Trans. Pattern Anal. Mach. Intell.* **26**(10), 1367 (2004)
9. Cordella, L.P., Foggia, P., Sansone, C., Vento, M.: An improved algorithm for matching large graphs. In: *Proceedings of the 3rd IAPR-TC15 Workshop on Graph-based Representations in Pattern Recognition*, pp. 149–159 (2001)
10. Corradini, A., Montanari, U., Rossi, F., Ehrig, H., Heckel, R., Löwe, M.: Algebraic approaches to graph transformation - Part I: Basic concepts and double pushout approach. In: Rozenberg, G. (ed.) *Handbook of Graph Grammars and Computing by Graph Transformation*. Chapter 3, pp. 163–245. World Scientific, Singapore (1997)
11. Ehrig, K., Heckel, R., Lajios, G.: Molecular analysis of metabolic pathway with graph transformation. In: Corradini, A., Ehrig, H., Montanari, U., Ribeiro, L., Rozenberg, G. (eds.) *ICGT 2006. LNCS*, vol. 4178, pp. 107–121. Springer, Heidelberg (2006)
12. Fernández, M., Kirchner, H., Namet, O.: A strategy language for graph rewriting. In: Vidal, G. (ed.) *LOPSTR 2011. LNCS*, vol. 7225, pp. 173–188. Springer, Heidelberg (2012)
13. Flamm, C., Ullrich, A., Ekker, H., Mann, M., Högerl, D., Rohrschneider, M., Sauer, S., Scheuermann, G., Klemm, K., Hofacker, I.L., Stadler, P.F.: Evolution of metabolic networks: a computational framework. *J. Syst. Chem.* **1**(4), 4 (2010)
14. *Increpare games: Catalan* (2011). <http://www.increpare.com/2011/01/catalan/>
15. Gansner, E.R., North, S.C.: An open graph visualization system and its applications to software engineering. *Softw. Pract. Exp.* **30**(11), 1203–1233 (2000)
16. Himsolt, M.: GML: a portable graph file format. <http://www.fim.uni-passau.de/fileadmin/files/lehrstuhl/brandenburg/projekte/gml/gml-technical-report.pdf>

17. Kreowski, H.J., Kuske, S.: Graph multiset transformation: a new framework for massively parallel computation inspired by DNA computing. *Nat. Comput.* **10**(2), 961–986 (2011)
18. Mann, M., Ekker, H., Flamm, C.: The graph grammar library - a generic framework for chemical graph rewrite systems. In: Duddy, K., Kappel, G. (eds.) ICMB 2013. LNCS, vol. 7909, pp. 52–53. Springer, Heidelberg (2013)
19. O’Boyle, N.M., Banck, M., James, C.A., Morley, C., Vandermeersch, T., Hutchison, G.R.: Open Babel: an open chemical toolbox. *J. Cheminformatics* **3**, 33 (2011)
20. Rosselló, F., Valiente, G.: Analysis of metabolic pathways by graph transformation. In: Ehrig, H., Engels, G., Parisi-Presicce, F., Rozenberg, G. (eds.) ICGT 2004. LNCS, vol. 3256, pp. 70–82. Springer, Heidelberg (2004)
21. Rosselló, F., Valiente, G.: Chemical graphs, chemical reaction graphs, and chemical graph transformation. *Electron. Notes Theor. Comput. Sci.* **127**(1), 157–166 (2005). Proceedings of the International Workshop on Graph-Based Tools (GraBaTs 2004) Graph-Based Tools 2004
22. Siek, J.G., Lee, L.Q., Lumsdaine, A.: Boost Graph Library: The User Guide and Reference Manual. Pearson Education, Upper Saddle River (2001). <http://www.boost.org/libs/graph/>
23. Sylvester, J.J.: On an application of the new atomic theory to the graphical representation of the invariants and covariants of binary quantics, with three appendices. *Am. J. Math.* **1**(1), 64–128 (1878)
24. Taentzer, G.: AGG: A graph transformation environment for modeling and validation of software. In: Pfaltz, J.L., Nagl, M., Böhlen, B. (eds.) AGTIVE 2003. LNCS, vol. 3062, pp. 446–453. Springer, Heidelberg (2004)
25. Tantau, T.: The TikZ and PGF Packages (2013). <http://sourceforge.net/projects/pgf/>
26. Weininger, D.: SMILES, a chemical language and information system. 1. Introduction to methodology and encoding rules. *J. Chem. Inf. Comput. Sci.* **28**(1), 31–36 (1988)
27. Yadav, M.K., Kelley, B.P., Silverman, S.M.: The potential of a chemical graph transformation system. In: Ehrig, H., Engels, G., Parisi-Presicce, F., Rozenberg, G. (eds.) ICGT 2004. LNCS, vol. 3256, pp. 83–95. Springer, Heidelberg (2004)