

UNIVERSITY OF SOUTHERN DENMARK

DEPARTMENT OF MATHEMATICS
AND COMPUTER SCIENCE

MASTER THESIS – 31/12-2017

Computational Synthesis Planning Using Big Data

Computational Syntese planlægning ved hjælp af big data

Author:
Henrik Schulz

Supervisors:
Daniel Merkle



Abstract

Resumé

Contents

1	Acknowledgements	3
2	Introduction	3
2.1	Overview	3
3	Preliminaries	3
4	Finding The K-Best Synthetic Plans	4
4.1	Yen's Algorithm	4
4.2	Yen's Algorithm On Hypergraphs	4
5	Shortest Path	8
5.1	Dynamic Approach	8
	5.1.1 Approach	8
	5.1.2 Testing	9
	5.1.3 Problems	9
5.2	Nielsens Algorithm	9
	5.2.1 Approach	9
	5.2.2 Testing	9
	5.2.3 Problems	9
	5.2.4 Optimizing	9
6	Work With Beilstein Data	9
6.1	The Database	9
6.2	Data Assessment	10
6.3	Results	16
	6.3.1 Strychnine	16
	6.3.2 Colchicine	16
	6.3.3 Dysidiolide	16
	6.3.4 Asteriscanolide	16
	6.3.5 Lepadiformine	16
7	Conclusion	16
8	Appendix	18

1 Acknowledgements

2 Introduction

kkk [2]

2.1 Overview

Hvad indeholder hver sektion??

3 Preliminaries

This section contains definitions that will be used throughout this paper. It is assumed that the reader have a basic understanding of graph theory.

Hypergraphs

A directed hypergraph h is a set V of vertices and a set E of hyperedges, where each hyperedge $e = (T(e), H(e))$ is an ordered pair of non-empty multi-sets of vertices. The set $T(e)$ is denoted as the tail of the hyperedge and $H(e)$ is the head. If $|H(e)| = 1$ then the hyperedge is denoted as a B-hyperedge. If all edges in the hypergraph is B-hyperedges, then the graph is denoted as a B-hypergraph. This paper will only consider hypergraphs that are B-hypergraphs. A hypergraph $H' = (V', E')$ is a subhypergraph of $H = (V, E)$ if $V' \subseteq V$ and $E' \subseteq E$. [3]

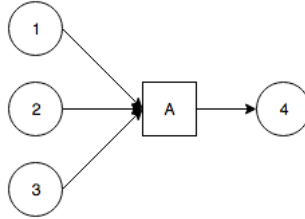


Fig. 1: Example of hyperedge A . $T(A) = \{1, 2, 3\}$, $H(A) = \{4\}$

Hyperpaths

A path P_{st} from s to t in a B-hypergraph is a sequence $P_{st} = \langle e_1, e_2, e_3, \dots, e_q \rangle$ of B-hyperedges such that $s \in T(e_1)$ and $t = H(e_q)$ and $H(e_i) \in T(e_{i+1})$ for $i = 1..q - 1$. Its length $|P_{st}|$ is the number q of hyperedges. If $t \in T(e_1)$, then P_{st} is a cycle. A hypergraph is acyclic if it does not contain any cycles. [3]

A hyperpath $\pi_{st} = (V_\pi, E_\pi)$ from a source vertex s to a target vertex t in a B-hypergraph H is a subhypergraph of H with the following properties: If $t = s$, the $V_\pi = \{s\}$ and $E_\pi = \emptyset$. Otherwise, E_π can be ordered in a sequence $\langle e_1, e_2, \dots, e_q \rangle$ such that

1. $T(e_i)\{s\} \cup \{H(e_1), H(e_2), \dots, H(e_{i-1})\}$ for all i
2. $t = H(e_q)$

3. Every $v \in V_\pi \setminus \{t\}$ has at least one outgoing hyperarc in E_π , and t has zero.
 4. Every $v \in V_\pi \setminus \{s\}$ has at least one ingoing hyperarc in E_π , and s has zero.
- [3]

4 Finding The K-Best Synthetic Plans

This section describes how to find the K-Best paths of a hypergraph modifying an algorithm made by Jin Y. Yen. The section starts out by describing the work flow of the algorithm, and then proceeds to deal with the alterations that are needed to run the algorithm on a hypergraph.

4.1 Yen's Algorithm

Yen's algorithm is an algorithm that computes the K-shortest paths for a graph with non-negative edges. It was published in 1971 and uses any shortest path algorithm to find the best path and then proceeds to find the $K - 1$ deviation of the best path. [4]

It starts out by finding the best path using a shortest path algorithm. Once the best path has been found it uses the path to find all the potential second best paths by fixing and removing edges in the graph.

By using the same first vertex as the original path but removing the first edge, it forces the shortest path algorithm to take another route through the graph and thereby creating a potential second best path. This is added to the list of potential paths and the algorithm can continue to derive other paths from the best path. By fixing the first edge in the previous best plan, Yen's algorithm forces the shortest path algorithm to take the first edge which it now shares with the best path. However, now the algorithm has removed the second edge from the original path and once again forces the shortest path algorithm to find alternative routes. This process is then repeated until we reach the next to last vertex in the best path.

By sorting the list of potential paths, it has the second best path at the start of the list and it can add it to the final list of best path. The algorithm then repeats on the second best path to find the third best path. This is done until all K-best paths have been found or there are no more paths to find.

4.2 Yen's Algorithm On Hypergraphs

We use the principles from Yen's algorithm to make our own algorithm that will work on hypergraphs. To handle the problem of generating all derived paths from our best path in our hypergraph, we use a method called Backwards-Branching. [3] [5] [6]

Algorithm 1: Backwards Branching for B-Hypergraph

```

1 function Back-Branch( $H, \pi$ )
2    $B = \emptyset$ 
3   for  $i = 1$  to  $q$  do
4     Let  $H^i$  be a new hypergraph
5      $H^i.V = H.V$ 

```

```

6      // Remove hyperarc from H
7       $H^i.E = H.E \setminus \{\pi.p(v_i)\}$ 
8      // Fix Back tree
9      for  $j = i+1$  to  $q$  do
10          $H^i.BS(v_j) = \setminus \{\pi.p(v_j)\}$ 
11          $B = B \cup \setminus \{H^i\}$ 
12     return B

```

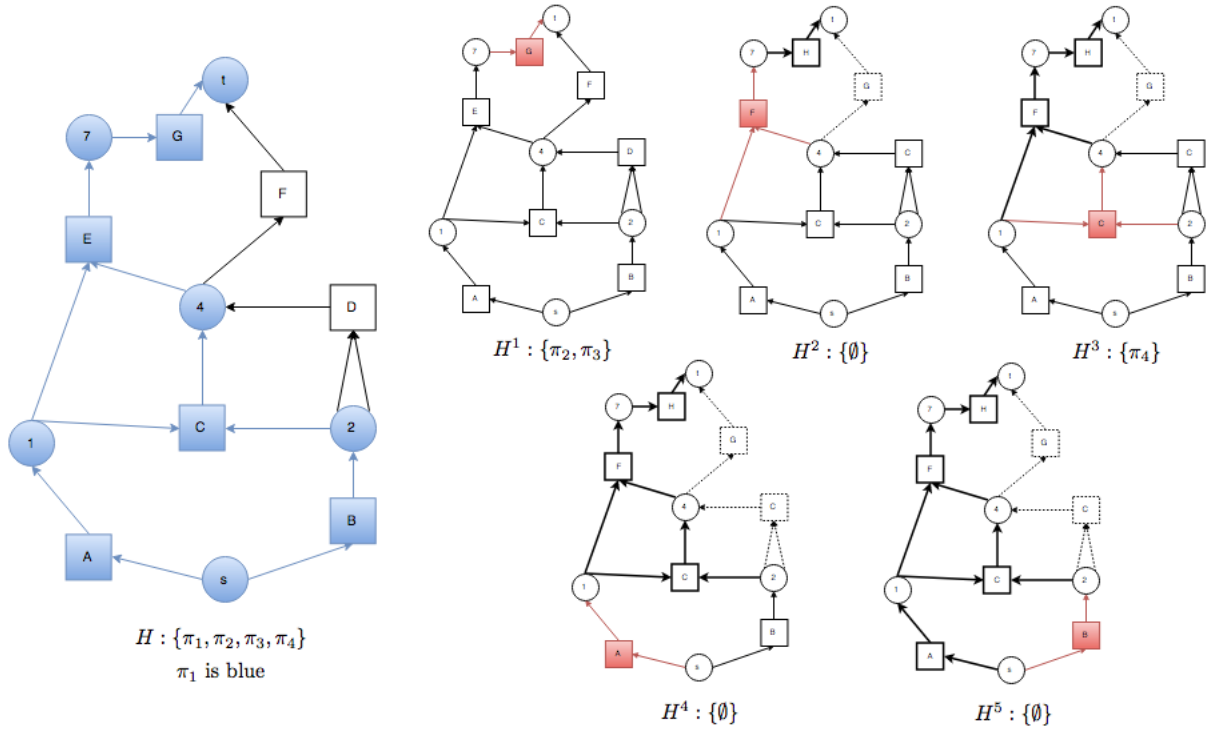


Fig. 2: An example of branching for hypergraphs. H is the original hypergraph. The vertices and hyperedges marked with blue are part of the best path. The rest of the figure illustrates the backwards branching. Each of the smaller figures shows a hypergraph H^i and how it is created from H . Dotted hyperedges and red hyperedges are not part of the hypergraph, but have been deleted due to branching. Hyperedges in bold are fixed hyperedges. When hyperedges are fixed it leads to other hyperedges being deleted (dotted). The caption beneath each hypergraph represents the possible paths that are available in the given hypergraph.

However, this algorithm have a problem when working on a larger hypergraph. It demands that each time we make alterations on the hypergraph we have to make a copy, H^i , of the graph, H , with the exception of the hyperedges that is removed when fixing the back tree and removing $\pi.p(v_i)$.

This could easily work for smaller graphs, but if we use this on the hypergraph

that we generate from the beilstein database, we would have to copy a graph of multiple GigaBytes.

To handle this problem I came up with the idea of creating an overlay for the graph instead of copying it. The overlay would work as an transparent on top of the original graph, stating which edges still is accessible. This is done by creating a *vector<bool>* which has a length of R , where R is the number of reactions. Normally a reaction would contain at least 28 bytes of data:

- 3x ints of 4 bytes each
- 1x double of 8 bytes
- 1x *vector<int>* head of length one of at least 4 bytes
- 1x *vector<int>* tail of length N (number of educts) of at least 4 bytes

This can be reduced dramatically by using the *vector<bool>*, since c++ only uses 1 bit per boolean in the vector instead of the regular 1 byte per boolean.[7] This means if working on a hypergraph with 40 million reactions, we would be able to create an overlay using 5 MB of space per alternated graph, instead of copying a hypergraph were the reactions alone uses at least 1,12 GB per copy. As the figure below shows, we never change or remove anything on the hypergraph. We simply create the following overlay:

Reaction	A	B	C	D
Usable	true	true	true	false
Bit Representation	1	1	1	0

And then when trying to use an edge, we ask: "Does overlay at reaction A exist?". If yes, you can use it. If no, the edge have been "removed", and therefore cannot be used.

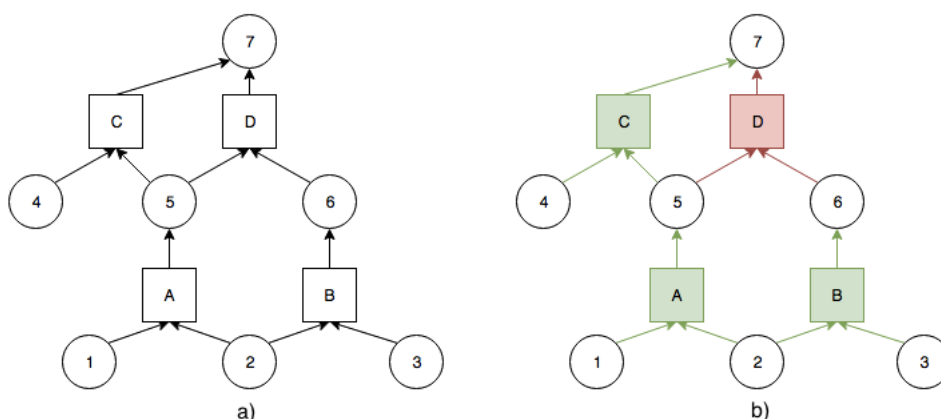


Fig. 3: a) The original hypergraph. b) The original hypergraph, but using the overlay. If the reaction is green it is still usable. If red then it have been "removed" from the hypergraph.

This of course means that the algorithm for back-branch have to be changed accordingly. Instead of the hypergraph as input we now give it an overlay. This overlay is changed so it fits with the new layout of the graph. Instead of deleting hyperedges in the copy, we now simply changes the boolean at the index of the hyperedge.id. True if we should "add" the hyperedge and false when we want to "remove" a hyperedge.

Algorithm 2: Backwards Branching for B-Hypergraph using overlay

```

1 function BackwardsBranching( $\pi$ , Overlay)
2   List B =  $\emptyset$ 
3   // q = Path length
4   for i = 1 to q do
5     //remove i'th hyperedge from Path in overlay
6     Set Overlay[ $\pi[i]$ ] to false
7     //fix the backtree
8     for j = i downto 1 do
9       vertex C  $\leftarrow$   $\pi[j]$ .head
10      for each hyperedge into C
11        Set Overlay[reaction.id] to false
12      Set Overlay[ $\pi[j]$ ] to true
13    endfor
14    B = B  $\cup$  {Overlay}
15  endfor
16  return B
17 endfunction

```

Now that we have the branching in place are we able to construct an algorithm that are similar to Yen's algorithm, but can run on hypergraphs. As input it takes a start vertex, a goal vertex, and an integer K , where K is the number of best paths we want to find.

It creates a heap, L , and a list, A , which will contain the K -best paths once the algorithm is finished. It then finds the best path using a shortest path algorithm and inserts it into the heap. Inside the loop it extracts the best path found from the heap and performs a backward branching, and finds all possible paths in the branches. If there is a path from s to t , then this path is added to the heap. The algorithm either terminates if the heaps is empty (No more paths available) or once it have found the K -best paths.

Algorithm 3: K-Shortest Paths Algorithm in B-Hypergraph

```

1 function YenHyp(s, t, K)
2   L = new heap with elements (overlay,  $\pi$ )
3   A = List of shortest paths
4   //(Graph is default overlay (all true))
5    $\pi$  = shortestPath(Graph, s,t)
6   Insert (Graph,  $\pi$ ) into L
7   for k = 1 to K do
8     if L =  $\emptyset$ 
9       Break
10    endif
11    (Overlay',  $\pi'$ ) = L.pop

```



```

12   add  $\pi'$  to A
13   for all  $\text{Overlay}^i$  in  $\text{BackwardBranching}((\text{Overlay}', \pi'))$  do
14        $\pi^i = \text{shortestPath}(\text{Overlay}^i, s, t)$ 
15       if  $\pi^i$  is complete
16           Insert(  $H^i, \pi^i$ ) into L
17       endif
18   endfor
19 endfor
20 return A
21 endfunction

```

5 Shortest Path

This section describes two different approaches to the shortest path problem in a hypergraph. A dynamic approach proposed by Carsten Grønberg Lützen and Daniel Fentz Johansen [6] and a Dijkstra inspired approach proposed by Lars Relund Nielsen, Kim Allan Andersen and Daniele Pretolani [5].

5.1 Dynamic Approach

kkk [6]

kkkk [3]

5.1.1 Approach

Algorithm 4: Dynamic programming for finding the best path

```

1 function Min(V)
2   if (V) is starting material
3       return Cost of V
4   endif
5   mincost  $\leftarrow$  inf
6   for all  $e \in BS(V)$  do
7       cost  $\leftarrow$  cost of e
8       for all  $u \in \text{Tail}(e)$  do
9           cost  $\leftarrow$  cost + Min(u)
10      endfor
11      if mincost  $\leq$  cost
12          mincost  $\leftarrow$  cost
13          V.minedge  $\leftarrow$  e
14      endif
15  endfor
16  return mincost
17 endfunction

```

5.1.2 Testing**5.1.3 Problems****5.2 Nielsens Algorithm**

k [5]

5.2.1 Approach

Algorithm 5: Dynamic programming for finding the best path

```

1 Initialization: Set  $W(u) = \inf \forall u \in V, k_j = 0 \forall e_j \in E, Q = \{s\}$  and  $W(s) = 0$ 
2 function SBT-Dijkstra
3   while ( $Q \neq \emptyset$ ) do
4     select and remove  $u \in Q$  such that  $W(u) = \min\{W(x) | x \in Q\}$ 
5     for ( $e_j \in FS(u)$ ) do
6        $k_j < -k_j + 1$ 
7       if ( $k_j = |T(e_j)|$ )
8          $v \leftarrow h(e_j)$ 
9         if ( $W(v) > w(e_j) + F(e_j)$ )
10          if ( $v \notin Q$ )
11             $Q \leftarrow Q \cup \{v\}$ 
12          endif
13           $W(v) \leftarrow w(e_j) + F(e_j)$ 
14           $p(v) \leftarrow e_j$ 
15        endif
16      endif
17    endfor
18  endwhile
19 endfunction

```

5.2.2 Testing**5.2.3 Problems****5.2.4 Optimizing****6 Work With Beilstein Data**

This section HOFFMANN, REAXYS

6.1 The Database

The Beilstein database is the largest database in the field of organic chemistry. Since 2009, the content has been maintained and distributed by Elsevier Information Systems in Frankfurt under the name "Reaxys". The content covers more than 200 years of chemistry and has been abstracted from several thousands of journal titles, books and patents. Today the data is drawn from selected journals (400 titles) and chemistry patents, and the extraction process for each reaction or substance data included needs to meet three conditions:

1. It has a chemical structure
2. It is supported by an experimental fact (property, preparation, reaction)
3. It has a credible citation

Journals covered include *Advanced Synthesis and Catalysis*, *Angewandte Chemie*, *Journal of American Chemical Society*, *Journal of Organometallic Chemistry*, *Synlett* and *Tetrahedron*. [8][9]

6.2 Data Assessment

During my work with the Beilstein database have I found several issues when it comes to using it to find shortest paths. First problem is that there are multiple instances of the same compound, each with a different Reaxys IDs. As seen in fig.4 have I found four different IDs for the compound Dysidiolide. These were found when I tried to reproduce the different synthesis plans of Dysidiolide from [1] using the referenced articles where each synthesis plan origins. E. J. Corey's version of Dysidiolide have the ID 8171938, Boukouvalas have two different with ID 7601810 and 7910427 and Danishefsky have the ID 7910428.

Since we have four different IDs for Dysidiolide we can't state a single goal compound to our program that would result in giving us these three synthesis plans. To handle this problem the program can take several goal compounds as an input, and by creating a dummy node t are we able to give the illusion of a single target.

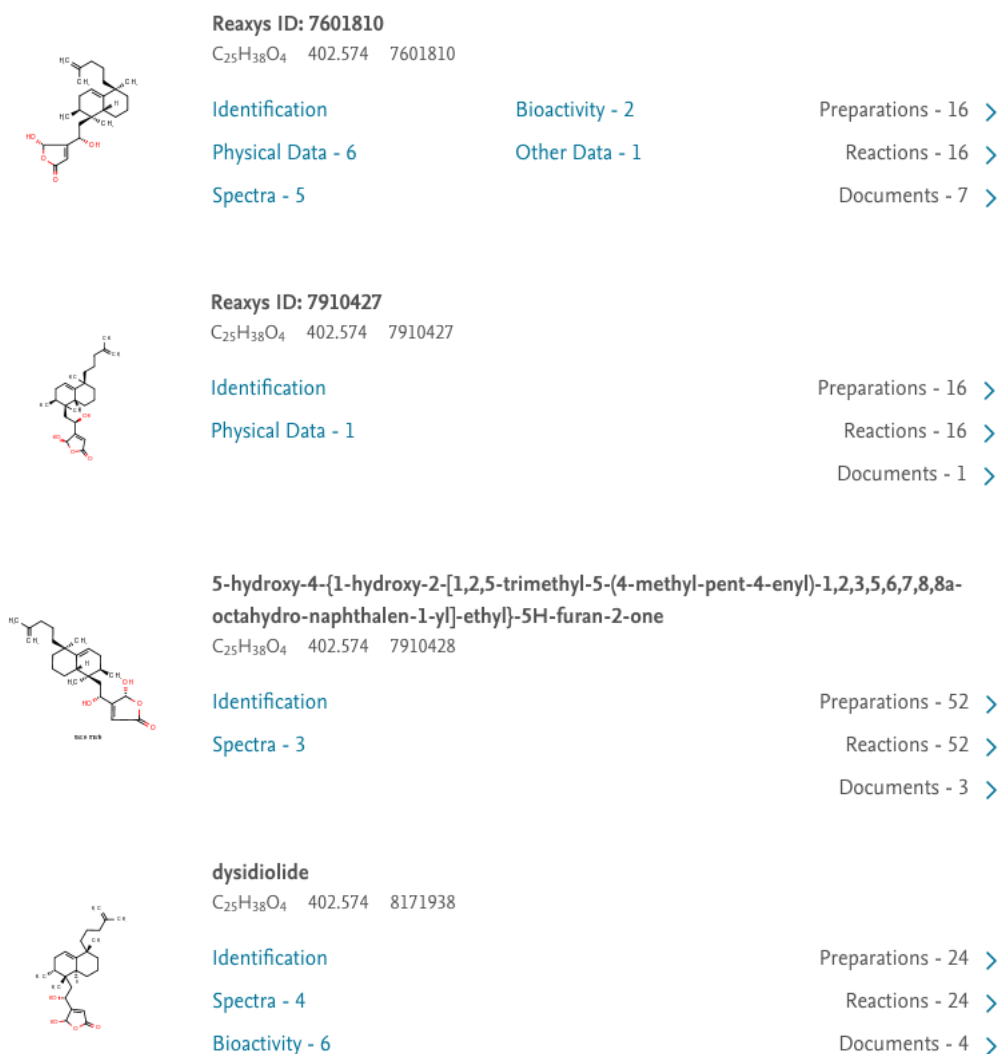


Fig. 4: Four different instances of Dysidiolide in Reaxys.

So what is causing this multiple ID issue? If we study fig. 4 can we see that the molecular weight and molecular formula are exactly the same, however if we look at the close up of the molecular structure in fig. 5 the compounds are not structured in the same way, even though the compounds are the same. The main differences is:

1. Which way some of the substructures are facing. Example: The lower $C_4O_3H_4$ is rotated differently in each instance or that we write H_2C instead of CH_2 .
2. How the bond between two chemical elements are notated. Example: The bond to the OH in the bottom of the structure is either a "single" (d), "single down" (a) (c) or "single up or down" (b).

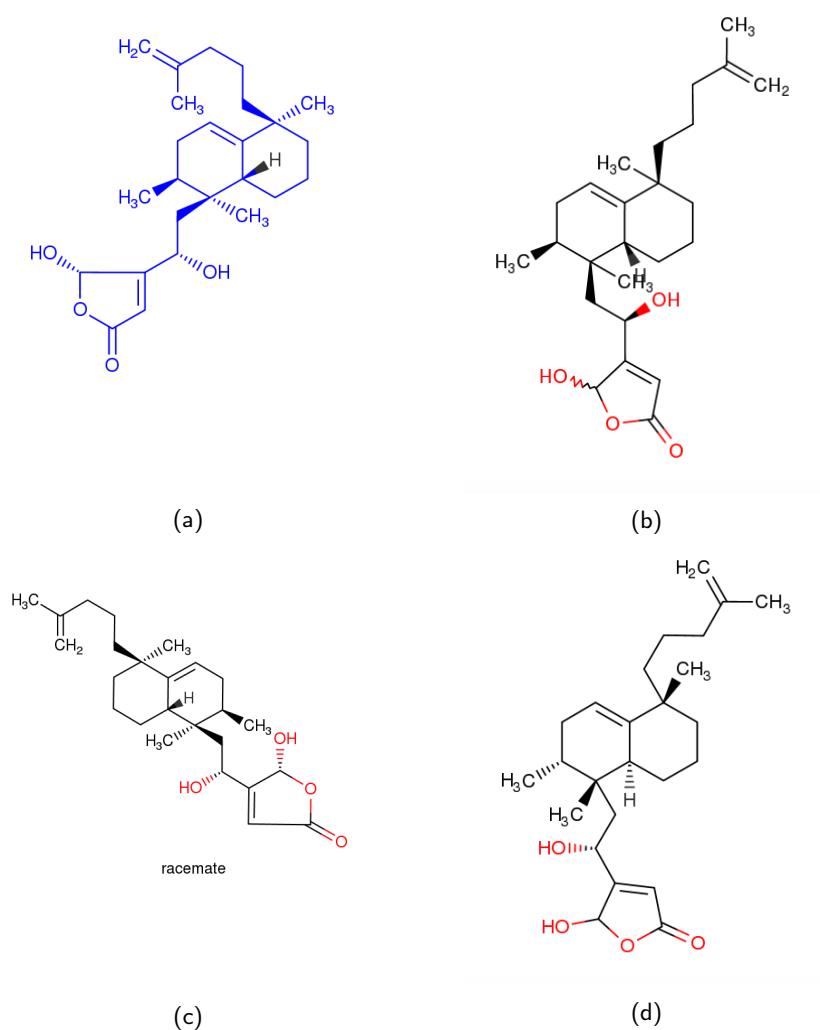


Fig. 5: Different versions of Dysidiolide: ID 7601810 (a), ID 7910427(b), ID 7910428(c), ID 8171938(d)

Second issue is the problem of a reaction not having an educt or a product(fig: 6). This leaves the reaction incomplete and makes it useless in the graph construction. If there is no educt the hyperedge created will have an indegree of 0, and thereby making it unreachable. If there is no product the hyperedge will have an outdegree of 0, making it a deadend.

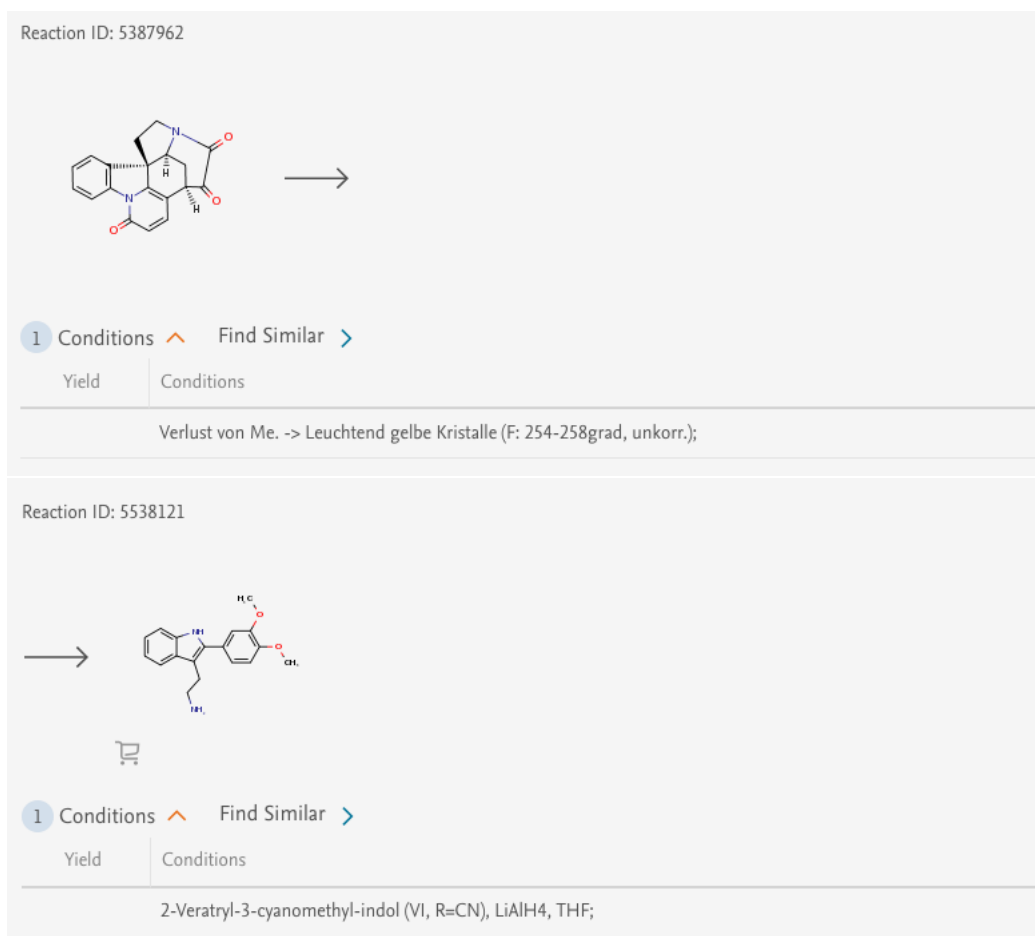


Fig. 6: Example of reactions that either is missing educts or products.

Third issue is that some reactions contains educts where a compound isn't linked to a compound in the database. This results that the name of the compound is written instead of the usual structure diagram (fig 7). This issue does not give problems when it comes to the construction of the hypergraph. The compound is simply just not added to the hypergraph. This however results in a slightly misleading result if the user does not look up the reaction in Reaxys where the name of the educt is stated. Example: If $A + B \rightarrow C$ but B is not given an compound ID the reaction would look like $A \rightarrow C$ in the hypergraph.



Fig. 7: Example of an reaction with a missing educt.

Fourth issue is multireactions. (fig 8) A multireaction is a reaction with its own ID, but it consists of an educt and a product where there are multiple reaction steps, s , between the two compounds. This means that instead of s different reactions with their own ID and yield we get a single reaction without a yield. The yield for each reaction in the multistep reaction is often stored as a part of the reaction text, but not easy extractable. The multistep reaction should however only consist of individual reactions that already are in the database. The solution to this problem have been to not include all reactions labelled as "Multi-step reaction" to the extraction from the database. Since the steps should be saved as individual reactions this would not cause any harm to possibility of finding the exact same path, but only using all s steps instead of one.

Reaction ID: 14038382

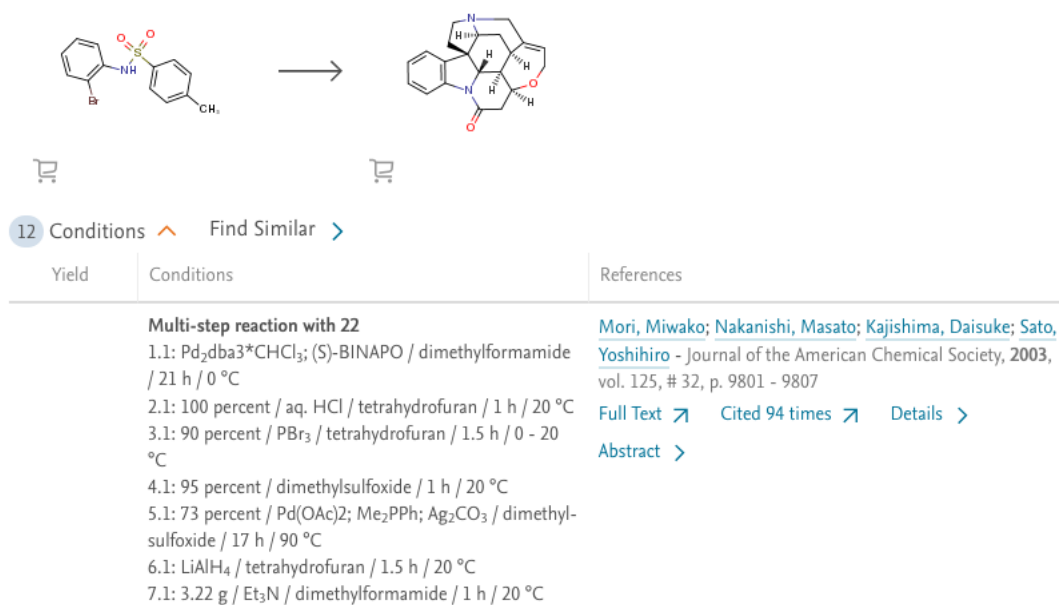


Fig. 8: Example of a multistep reaction.

Reaction ID: 9601319

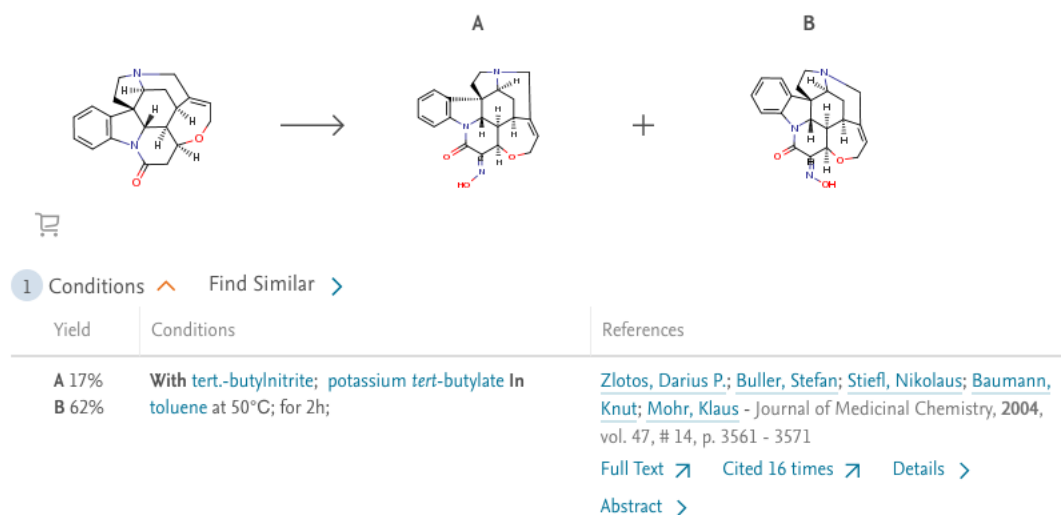


Fig. 9: Example of a reaction with multiple products

Fifth issue is reactions which have more than one product (fig. 9). If this becomes a part of our hypergraph it is no longer a B-hypergraph. This is a problem since both shortest paths algorithms only works on B-hypergraphs. This

problem is however solved after the graph have been created by the method *convertToBHypergraph()* in *Hypergraph.hpp*. The method iterates through the reaction list and if it encounters a non B-hyperedge, the hyperedge is added to a list of non B-hyperedges. For each of the non fixed B-hyperedges e , it creates $|H(e)|$ new hyperedges e_i , where $T(e_i) = T(e)$ and $H(e_i) = H(e)[i], \forall i \in \mathbb{N}, 1 \leq i \leq |H(e)|$ (fig. 10). Each of the new reactions contains an original-ID variable that points to the original hyperedge. This is only used to print the correct ID when the edge is used in a synthesis plan.

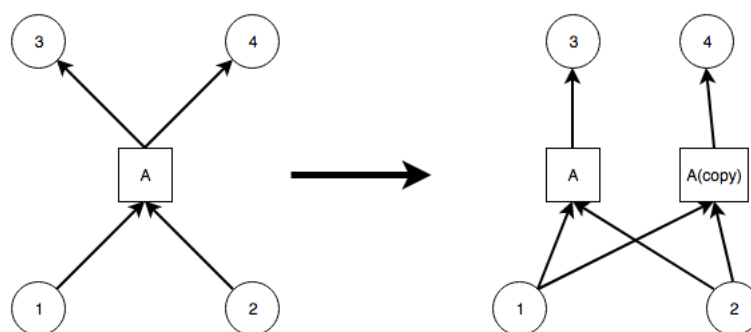


Fig. 10: Conversion to B-Hyperedge

6.3 Results

6.3.1 Strychnine

6.3.2 Colchicine

6.3.3 Dysidiolide

6.3.4 Asteriscanolide

6.3.5 Lepadiformine

kkkk [1]

7 Conclusion

Books

- [1] R. W. Hoffmann, *Elements of Synthesis Planning*. Springer Berlin Heidelberg, 2009.

Articles

- [2] S. Szymkuc, E. P. Gajewska, T. Klucznik, K. Molga, P. Dittwald, M. Startek, M. Bajczyk, and B. A. Grzybowski, “Computer-assisted synthetic planning: The end of the beginning”, *Angewandte Chemie International Edition*, no. 55, pp. 5904–5937, 2016.
- [3] R. Fagerberg, C. Flamm, R. Kianian, D. Merkle, and P. F. Stadler, “Finding the K best synthesis plans”, 2017, Unpublished Article.
- [4] J. Y. Yen, “Finding the K shortest loopless paths in a network”, *Management Science*, vol. 17, no. 11, pp. 712–716, Jul. 1971.
- [5] L. R. Nielsen, K. A. Andersen, and D. Pretolani, “Finding the K shortest hyperpaths: Algorithms and applications”, 2002.

Other

- [6] C. G. Lützen and D. F. Johansen, “A computational and mathematical approach to synthesis planning”, Master’s thesis, University of Southern Denmark, 2015.
- [7] Sep. 2017. [Online]. Available: http://en.cppreference.com/w/cpp/container/vector_bool.
- [8] Oct. 2017. [Online]. Available: <https://en.wikipedia.org/wiki/Reaxys>.
- [9] Oct. 2017. [Online]. Available: https://en.wikipedia.org/wiki/Beilstein_database.

8 Appendix