# Atomically resolve a metabolic reconstruction

Author(s): Hulda S. Haraldsdóttir and German A. Preciat Gonzalez, Systems Biochemistry Group, University of Luxembourg.

Reviewer(s): Catherine Clancy, Molecular Systems Physiology Group, University of Luxembourg.

Francisco J. Planes, Department of Biomedical Engineering and Sciences, Tecnun, University of Navarra.

## INTRODUCTION

Genome-scale metabolic network reconstructions have become a relevant tool in modern biology to study the metabolic pathways of biological systems *in silico*. However, a more detailed representation at the underlying level of atom mappings opens the possibility for a broader range of biological, biomedical and biotechnological applications than with stoichiometry alone.

A set of atom mappings represents the mechanism of each chemical reaction in a metabolic network, each of which relates an atom in a substrate metabolite to an atom of the same element in a product metabolite (Figure 1). To assign map reactions in a metabolic network reconstruction, one requires chemical structures in a data file format (SMILES, MDL MOL, InChIs), reaction stoichiometries, and an atom mapping algorithm.



Figure 1: Set of atom mappings for reaction L-Cysteine-L-homocysteine-Lyase (VMH ID: r0160).

Metabolites chemical structures can be obtained by different approaches such as draw them based on the literature using chemoinformatic software, or obtain them from metabolic databases either manually or using a computational software as suggested in [1]. Here we recommend downloading the metabolites structures in MDL MOL format for the latest human metabolic reconstruction Recon 3 [2] via the Virtual Metabolic Human database (VMH, http://vmh.life). Chemical structures and reaction stoichiometries from COBRA models are used to generate an MDL RXN file, which contains the information of a chemical reaction. Atom mapped reactions from Recon 3 can also be found in the VMH database in MDL RXN format. However, here we will atom map the chemical reactions using the Reaction Decoder Tool (RDT) algorithm [3], which was selected after comparing the performance of recently published algorithms [1]. However, despite its good performance (accuracy and availability) RDT algorithm does not map hydrogen atoms.

In this tutorial, we will identify the conserved metabolic moieties when mapping data for the dopamine synthesis network (DAS) extracted from Recon 3 [2] (Figure 2). Section 1 of the tutorial will cover obtaining and visualising an atom map of metabolic reactions, and section 2 of the tutorial covers the identification of conserved metabolic moieties.
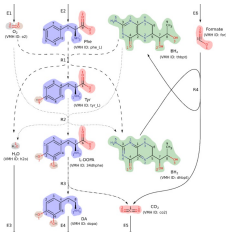
Figure 2: DAG: a small metabolic network consisting of reactions in the human dopamine synthesis pathway [1]. Atoms belonging to the same conserved moiety have identically coloured backgrounds.

## MATERIALS

To atom map reactions it is required to have Java version 8 and Linux. The atom mapping does not run on Windows at present.

On macOS, please make sure that you run the following commands in the Terminal before continuing with this tutorial:

```
$ /usr/bin/ruby -e "$(curl -fsSL https://raw.githubusercontent.com/Homebrew/install/master/install)"
$ brew install coreutils
```

On linux, please make sure that Java and ChemAxon directories are included. To do this, run the following commands:

```
$ export PATH=$PATH:/opt/opt/chemaxon/jchemsuite/bin/  (default location of JChem)
$ export PATH=$PATH:/usr/java/jre1.8.0_131/bin/  (default location of Java)
```

Also, in order to standardize the chemical reaction format it is required to have JChem downloaded from ChemAxon with its respective license.

## SECTION 1 Atom mapping of reactions

Atom mappings for the internal reactions of a metabolic network reconstruction are performed by the function `obtainAtomMappingsRDT`. The main inputs are a COBRA model structure and a directory containing the molecular structures in MDL MOL format. For this tutorial, using the RDT algorithm, the atom mappings are generated based on the molecular structures contained in `cobratoolbox/tutorials/atomicallyResolveReconstruction/data/molFiles` (molFileDir) and the reconstructed DAS network without hydrogen atoms (model).

```
global CBTDIR
tutorialDir = fileparts(which('tutorial_atomicallyResolveReconstruction.mlx'));
model = readCbModel([tutorialDir filesep 'data' filesep 'subDas.mat']); % The subnetwork of the dopamine synthesis network

model =

        S: [30×28 double]
      mets: [30×1 cell]
        b: [30×1 double]
    csense: [30×1 char]
      rxns: [28×1 cell]
        lb: [28×1 double]
        ub: [28×1 double]
         c: [28×1 double]
     osense: -1
     genes: {6×1 cell}
     rules: {28×1 cell}
 metCharges: [30×1 double]
metFormulas: {30×1 cell}
metSmiles: {30×1 cell}
  modelID: 'subDas'
```

```
molFileDir = [tutorialDir filesep 'data' filesep 'molFiles']; % The chemical structures of metabolites
```

The function `obtainAtomMappingsRDT` generates 4 different directories containing:

- the atom mapped reactions in MDL RXN format (directory `atomMapped`),
- the images of the atom maps reactions (directory `images`),
- additional data for the atom mapped reaction (SMILES, and product and reactant indexes) (directory `rxnData`), and
- the unmapped MDL RXN files (directory `rxnFiles`).

The input variable `outputDir` indicates the directory where the folders will be generated (by default the function assigns the current directory).

```
outputDir = [pwd filesep 'output']
```

For some reactions, the RDT algorithm cannot compute the atom mappings (for a large reaction is generated an MDL RXN v3000 which is not compatible with the RDT algorithm). Therefore, it is necessary to assign a maximum time of processing `maxTime` (by default the function assign 30 minutes as a maximum time for computing an atom mapping for each reaction).

```
maxTime = 1800; % seconds
```

The function `obtainAtomMappingsRDT` generates atom mapped reactions in a standard canonical format but it is **REQUIRED** to have a ChemAxon license installed. However, the reactions can be atom mapped without being standardized. The variable `isChemaxonInstalled` contains a logical value defined by the user if the license is installed or not.

```
isChemaxonInstalled = false; % Change variable to "true" if Chemaxon is installed
```

We can obtain the atom mapping using `obtainAtomMappingsRDT`:

```
if ispc
    error('Error: atom mapping function should be run on Linux or MAC.')
else
    tic
    try
        standardisedRxns = obtainAtomMappingsRDT(model, molFileDir, outputDir, maxTime, isChemaxonInstalled);
    end
    toc
```

```
   end
```

Generating RXN files.

Computing atom mappings for 6 reactions.

```
4 reactions were atom mapped
0 reactions are not standardized
0 reactions were not mapped
```

RDT algorithm was developed by:
SA Rahman et al.: Reaction Decoder Tool (RDT): Extracting Features from Chemical
Reactions, Bioinformatics (2016), doi: 10.1093/bioinformatics/btw064
Elapsed time in 15.084866 seconds.

The output, **standardDadenBasis**, is a list of atom mapped mass balanced reactions.

### TIMING

The time to compute atom mappings for metabolic reactions depends on the size of the genome-scale model and the size of the molecules in the reactions. The above example took ∼1 min or less if **isChemAxonInstalled = false**.

### Visualising results

The images directory contains a graphical representation of the atom mapped reactions. They show the bijection between atoms and each of the metabolite pairs are coloured for an easy visualisation. Figure 3 shows the atom mapped reaction to produce dopamine and $CO_2$ from L-DOPA.



Figure 3: Reaction 3-hydroxy-L-tyrosine carboxy-lyase atom mapped (VMH ID: 3HLYTCL) here represented as R3. Images generated by RDT algorithm, also shows where a reaction centre occurs. The rxnfiles directory contains for all atom mapped reactions a corresponding MDL RXN file (Figure 4). Contained within these files are information of the chemical reaction, such as:

- the name of the reaction (on line 2 of the file),
- the chemical formula (on line 4 of the file),
- the number of substrates and products (on line 5 of the file), and
- specific information for each of the molecules (from line 6 onwards, with the identifier SMOL).

```
 1   $RXN
 2   RI
 3
 4   pew_llci + ohepr(n) + alci(-v  sym_llci + bhepr(n) + b2vlci
 5     0   1
 6   $MOL
 7   pew_llci
 8
 9       Vent677  1514177174120
10    12 12  0  0  0  0  0  0  0 999 V2000
11      -0.7146    2.0415    0.0000 C   0  0  0  0  0  0  0  0  0  0  0  0
12       0.0000    2.4750    0.0000 C   0  1  0  0  0  0  0  0  0  0  0  0
13      -0.7146    2.9420    0.0000 C   0  0  0  0  0  0  0  0  0  0  0  0
14      -0.7146    1.2275    0.0000 C   0  0  0  0  0  0  0  0  0  0  0  0
15       1.4289    2.4080    0.0000 O   0  0  0  0  0  0  0  0  0  0  0  0
16       1.4289   -0.0200    0.0000 O   0  0  0  0  0  0  0  0  0  0  0  0
17      -0.7146   -0.4520    0.0000 C   0  0  0  0  0  0  0  0  0  0  0  0
18       0.0000    0.0200    0.0000 C   0  0  0  0  0  0  0  0  0  0  0  0
19       0.0000    0.8290    0.0000 C   0  0  0  0  0  0  0  0  0  0  0  0
20       1.4289    0.3095    0.0000 O   0  0  0  0  0  0  0  0  0  0  0  0
21      -0.7146    2.7720    0.0000 O   0  0  0  0  0  0  0  0  0  0  0  0
22      -0.7146    3.7124    0.0000 O   0  0  0  0  0  0  0  0  0  0  0  0
23      1  1  1  0  0  0  0
24      2  3  1  0  0  0  0
25      3  4  1  0  0  0  0
26      4  5  1  0  0  0  0
27      4  7  2  0  0  0  0
28      5  6  1  0  0  0  0
29      6  9  1  0  0  0  0
30      2 11  1  0  0  0  0
31     10 12  2  0  0  0  0
32     10 11  1  0  0  0 -1
33     10  8  2  0  0  0  0
34     10  8  1  0  0  0 -1
35   M  END
36   $MOL
37   shepr(n)
```

Figure 4: A MDL RXN file stored in the rxnFake directory.

Specific information for each of the molecules includes the name of the metabolite, its InChI key (if the metabolite does not contain an R group) and the number of atoms and bonds. Following this is the atom block, which contains detailed information on the coordinates, element, charge and atom mapping number for each of the atoms, and then finally, the bond block connects all the atoms in the metabolite.

```
regexp(filereadr(outputDir fileseep "atomMapped" fileseep "RXL.txt"), "\n", "\n", "split")$
```

```
RXN =
'$RXN'
''
''  EC-BLAST    RXI'
''
'  1   2'
'$MOL'
'3A60pDe[L]'
'  EC-BLAST  2112171832'
'IsC33XXry=4TDZGC1KGRGINC-1J2LTPRIZGA-N'
'  12 14  0  0  0  0  0  0  0  0999 V2000'
'    -0.7165    -1.7125    0.0000  O   0   0   0   0   0   0   0   0   0   1   0   0'
'     0.8688     1.3898    0.0000  N   0   0   0   0   0   0   0   0   0   2   0   0'
'    -0.7103    -1.7125    0.0000  C   0   0   0   0   0   0   0   0   0   3   0   0'
'     0.8888     2.0738    0.0000  C   0   0   0   0   0   0   0   0   0   4   0   0'
'    -0.7101    -1.0435    0.0000  C   0   0   0   0   0   0   0   0   0   5   0   0'
'    -0.7101     2.0625    0.0000  C   0   0   0   0   0   0   0   0   0   6   0   0'
'     0.7165     1.2825    0.0000  C   0   0   0   0   0   0   0   0   0   7   0   0'
'     1.1283     0.0238    0.0000  C   0   0   0   0   0   0   0   0   0   8   0   0'
'     1.1285     0.9038    0.0000  C   0   0   0   0   0   0   0   0   0   9   0   0'
'     0.7165    -0.0125    0.0000  C   0   0   0   0   0   0   0   0   0  10   0   0'
'     1.1163    -1.2725    0.0000  C   0   0   0   0   0   0   0   0   0  11   0   0'
'     0.8888     0.8098    0.0000  C   0   0   0   0   0   0   0   0   0  12   0   0'
'    -0.7101    -0.0125    0.0000  C   0   0   0   0   0   0   0   0   0  14   0   0'
'  2  1  2  0  0  0  0'
'  2  3  1  0  0  0  0'
'  4  2  1  0  0  0  0'
'  5  3  1  0  0  0  0'
'  6  1  1  0  0  0  0'
'  7  5  1  0  0  0  0'
'  7  8  2  0  0  0  0'
'  8  9  1  0  0  0  0'
'  9 10  2  0  0  0  0'
' 10 11  1  0  0  0  0'
' 10 12  2  0  0  0  0'
' 12 13  1  0  0  0  0'
'M  CHG  1    3  -1'
'M  CHG  1    5   1'
'M  END'
'$MOL'
'$SMOL[L]'
'  EC-BLAST  2112171832'
'IsC33XXry=5YF7TLLB3K2KU-DR7FFX2YGA-N'
'  11 11  0  0  0  0  0  0  0  0999 V2000'
'    -0.7103    -0.0125    0.0000  C   0   0   0   0   0   0   0   0   0  10   0   0'
'     0.8688     0.8038    0.0000  N   0   0   0   0   0   0   0   0   0   2   0   0'
'     0.8688     2.0738    0.0000  C   0   0   0   0   0   0   0   0   0   4   0   0'
'     0.7103     1.2375    0.0000  C   0   0   0   0   0   0   0   0   0   7   0   0'
'     1.1283     0.0238    0.0000  C   0   0   0   0   0   0   0   0   0   8   0   0'
'     1.1285     0.9038    0.0000  C   0   0   0   0   0   0   0   0   0   9   0   0'
'     0.7103    -1.0735    0.0000  C   0   0   0   0   0   0   0   0   0   5   0   0'
'     1.1283    -1.2725    0.0000  C   0   0   0   0   0   0   0   0   0  11   0   0'
'     0.7165    -0.0125    0.0000  C   0   0   0   0   0   0   0   0   0  18   0   0'
'  2  1  1  0  0  0  0'
'  2  3  2  0  0  0  0'
' 11  2  1  0  0  0  0'
'  5  1  1  0  0  0  0'
'  4  2  1  0  0  0  0'
'  7  5  1  0  0  0  0'
'  5  8  1  0  0  0  0'
'  9 11  1  0  0  0  0'
' 11 10  1  0  0  0  0'
'M  CHG  1    4   1'
'M  END'
'$MOL'
'1u2I[L]'
'  EC-BLAST  2112171832'
'IsC33Xry=IDRLTGRVJ-IL2I-DR7FFX2YGA-N'
'   5  2  0  0  0  0  0  0  0  0999 V2000'
'     1.4168    0.8898    0.0000  O   0   0   0   0   0   0   0   0   0   1   0   0'
'     1.4168    0.8898    0.0000  O   0   0   0   0   0   0   0   0   0   2   0   0'
'     1.3888    0.8898    0.0000  O   0   0   0   0   0   0   0   0   0   4   0   0'
'     1  2  2  0  0  0  0'
'     2  3  2  0  0  0  0'
'M  END'
```

The bdZabs directory contains the TXT information of the reaction including the SMILES format, which holds the standard canonical format of the reaction,

the reactant input atom index and the product input atom index.

```
regexp(filemod[outputDir fileSep 'txtData' fileSep 'R2.txt'][], 'is', 'split')'

ans =

   '/*'
   'SELECTED ASM MAPPING'
   '[S:1]=[C:2]([C:3]=[O:4])N[N+5][R+6][O-7]=[C:6]([N-9]([N-10])=[C:11][N+12][N-13]([N-14][C-15]([C-16]([C-17]=[C-18])[O-19]...
   '
   '//'
   'REACTANT INPUT ATOM INDEX:-->AAM IS'
   '[1>1, 2>2, 3>3, 4>4, 5>5, 6>6, 7>18, 8>11, 9>12, 10>19, 11>6, 12>2, 13>5, 14>1]'
   'PRODUCT INPUT ATOM INDEX:-->AAM IS'
   '[1>21, 2>20, 3>9, 4>5, 5>4, 6>3, 7>8, 8>8, 11>6, 12>12, 13>13, 14>16]'
```

## SECTION 2 Identifying conserved metabolic moieties

A conserved moiety is a group of atoms within compounds conserved by reactions, that follow identical paths through a metabolic network and therefore, its amount remains constant (Figure 5). Representative examples from energy metabolism include the AMP and NAD moieties. With the set of atom mappings for a metabolic network the set of linearly independent conserved moieties for the metabolic network can be identified, each of which corresponds to a particular identifiable molecular substructure[1].



Figure 5: A graphical representation of a conserved moiety.

In this section, we will identify conserved moieties in a subnetwork of the DAS network (Figure 3) by graph theoretical analysis of its atom transition network.
The method is described in[1]. This section consists of two parts:

Part 1 covers basics set up of the code.

Part 2 covers decomposition of a conserved moiety resulting from variable atom mappings between the recurring metabolite pair $O_i$ and $H_i O_i$.

### Part 1: Generate an atom transition network for DAS based on atom mappings for internal (mass and charge balanced) reactions.

**Step 1: Generate an atom transition network for DAS based on atom mappings for internal (mass and charge balanced) reactions.**

The atom transition network is generated based on the reconstructed DAS network [model] and atom mappings for internal reactions, obtained in the previous section and predicted with the RDT algorithm[1].

```
if ~isDemoInstalled
    copyfile[outputDir fileSep 'data' fileSep 'atomMapped'][], [outputDir fileSep 'atomMapped']]
end
atomMappedDir = [outputDir fileSep 'atomMapped'];
ATN = buildAtomTransitionNetwork(model, atomMappedDir);

Atom mappings found for 4 model reactions.
Generating atom transition network for reactions with atom mappings.
```

The output variable (atn) is a Matlab structure with several fields. atn.A is the incidence matrix of the directed graph representing the atom transition network. Each row represents a particular atom in one of the 11 DAS metabolites. atn.mets indicates which metabolite the each atom belongs to. To find rows of atn.A corresponding to atoms in $CO_2$, type:

```
icol = find(ismember(ATN.mets, 'co2[c]'))'
```

```
icol =

   39   40   60
```

The order of atoms in atn.A matches their order in MDL MOL files encoding metabolite structures (Figure 7), e.g., atn.A(60,:) is the row corresponding to the second oxygen atom (number 3 in Figure 6).



Figure 6: Rows for $CO_2$ atoms in atn.A are ordered as shown.

atn.elements contains the element symbols of atoms, e.g.,

```
ATN.elements(98)

ans =

  98
```

Each column of `ATN.s` represents a particular atom transition in one of the four internal reactions in DAS. Reaction identifiers for atom transitions are given in `ATN.runs`. To find all atom transitions that involve $CO_2$, run:

```
tco2 = find(any(ATN.B(ico2,:), 1))

tco2 =

   75   76   77   93   96   97

ATN.runs(tco2)

ans =

   'R3'   'R3'   'R3'   'R6'   'R6'   'R6'
```

i.e., three atom transitions in each of the reactions R3 and R6 involve atoms in $CO_2$. To find atoms connected to $CO_2$ atoms via these atom transitions, run:

```
cco2 = find(any(ATN.B(:, tco2) < 0,2));
ATN.mets(cco2)

ans =

   '34dhpheコ'   '34dhpheコ'   '34dhpheコ'   'forコ'   'forコ'   'forコ'
```

i.e., $CO_2$ atoms are connected to atoms in the metabolites L-DOPA (VMH ID: 34dhphe) and formate (VMH ID: for).

**Step 3: Identify conserved moieties in DAS by graph theoretical analysis of the atom transition network generated in Step 1.**

```
tic
[L,lambda,moietyFormulas,moietiesmets,moietiesvectors,atoms2moieties] = ...
    identifyConservedMoieties(model, ATN);
t = toc;
fprintf('Computation time: %.3e s\n', t); % Print computation time

Computation time: 3.2e-01 s
```

This function assigns the moiety matrix (L), the moiety supergraph (Lambda), the chemical formulas of moieties (moietyFormulas), and three vectors that map between the various inputs and outputs. The L×d moiety matrix L has a row for each metabolite and a column for each conserved moiety in DAS. Each column is a moiety vector, with elements corresponding to the number of instances of a conserved moiety in each metabolite. To find the number of instances of moiety 2 in L-DOPA, run

```
iLDOPA = find(ismember(model.mets, '34dhpheコ'))

iLDOPA = 7

full(L(iLDOPA, 2))

ans = 1
```

i.e., L-DOPA contains one instance of moiety 2.

The 18×17 moiety supergraph (Lambda) contains the graphs of each of seven conserved moieties in DAS (Figure 7).
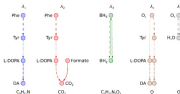


Figure 7: Graphs of the five conserved moieties in DAS. Each node represents an instance of a conserved moiety in a particular metabolite. Each directed edge represents conversion of a moiety between two metabolites. The chemical formula of each moiety is given below its graph.

Each row of Lambda represents a single instance of a conserved moiety in a particular metabolite. The vector moietiesvectors maps between the rows of Lambda and the columns of L. To obtain the incidence matrix of a particular moiety graph, e.g., A2 in Figure 7, run

```
l2 = find(moietiesvectors == 2);
c2 = find(any(L(xLambda(l2,:)));
lambda2 = full(Lambda(l2, c2))
```

```
lambda2 =
    -1     0     0     0
     1    -1     0     0
     0     1    -1     0
     0     0     1    -1
     0     0     0     1
```

The vector `moietiesInsets` maps the rows of Lambda to metabolite indices in the DAS reconstruction (`model1`). To find metabolites containing instances of moiety 2, run

```
m2 = moieties2mets(2);
mets2 = model1.mets(m2)
```

```
mets2 =
    'pam_L[c]'    'tyr_L[c]'    '34dhphe[c]'    'co2[c]'    'for[c]'
```

The chemical formula of moiety 2 is given by,

```
moietyFormulas(2)
```

```
ans = CO2
```

The vectors `atoms2moieties` maps each atom in the atom transition network for DAG to a particular instance of a conserved moiety. To find atoms in L-DOPA that belong to moiety 2, run

```
find(ismember(atoms2moieties, 12) & ismember(ATN.mets, '34dhphe[c]'))
```

```
ans =
    75    75    76
```

### Step 3: Classify moieties

```
types = classifyMoieties(L, model.S)
```

```
types =
    'Transitive'
    'Transitive'
    'Internal'
    'Transitive'
    'Transitive'
```

The internal moiety (λ3 in Figure 3) is conserved in both the open and closed DAG network, whereas the transitive and integrative moieties are only conserved in the closed network.[4]

## Part 2: Effects of variable atom mappings on recurring metabolite pairs

Here, we will again identify conserved moieties in DAG but with a slightly different set of atom mappings (Figure 6). The different atom mappings gives rise to a different atom transition network with a different set of conserved moieties. In particular, it contains a single composite moiety, λ6 in Figure 5, in place of the two moieties λ4 and λ5 in Figure 3. The composite moiety is the result of variable atom mappings between the recurring metabolite pair Q0 and H2O in reactions R1 and R2.
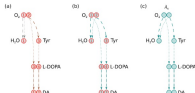


Figure 6: (a) Oxygen atom transitions used in Part 1. Oxygen atom 1 in O0 maps to the oxygen atom in H2O in both R1 and R2. These atom transitions contain two separate moieties, with two disconnected moiety graphs (λ4 and λ5 in Figure 7), and two linearly independent moiety vectors (L₄ and L₅ in Figure 3). (b) Oxygen atom transitions used in Part 2. A different oxygen atom maps from Q0 to H2O in R1 than in R2. These atom transitions contain only one composite moiety, λ6. The complete moiety graph arising from the oxygen atom transitions contains only one composite moiety, λ6. The complete moiety graph arising from the oxygen atom transitions is connected, with one composite moiety vector λ6.

### Step 1: Identify conserved moieties with the alternative set of atom mappings

```
% Create an alternative MDL RXN file
R2rxn = regexp(fileread([outputDir filesep 'atomMapped' filesep 'R2.rxn']), '\n', 'split')';
R2rxn{12} = 'alternativeR2';
R2rxn{165}(42:63) = '1$';
R2rxn{165}(62:63) = '18';
fid2 = fopen([outputDir filesep 'atomMapped' filesep 'alternativeR2.rxn'], 'w');
```

```
fprintf(fid2, 'main', B2run(1));
fclose(fid2);

% Create an alternative DAG model
alternativeModel = model;
alternativeModel.comps(2) = 'alternativeA2';

% Identify conserved moieties
ATN = buildAtomTransitionNetwork(alternativeModel, atomMappedDir);

Atom mappings found for 4 model reactions.
Generating atom transition network for reactions with atom mappings.

[L,Lambda,moietyFormulae,moieties,moietiesVectors,atoms2moieties] = ...
    identifyConservedMoieties(alternativeModel,ATN);
```

## Step 2: Decompose the composite moiety vector

First, extract the internal stoichiometric matrix for DAG, by running:

```
rbool = ismember(alternativeModel.rxns, ATN.rxns);
mbool = any(alternativeModel.S(:,rbool), 2);
N = alternativeModel.S(mbool, rbool);
```

To decompose the moiety vector computed in Step 1, run:

```
try
    changeCobraSolver('gurobi6', 'milp');
end
```

```
> Gurobi interface added to MATLAB path.
> gurobi (version TBD) is compatible and fully tested with MATLAB MIXING on your operating system.
```

```
D = decomposeMoietyVectors(L, N);
```

Note that you can use any MILP solver supported by the COBRA toolbox. The decomposed moiety matrix D is identical to the original moiety matrix computed in Part 1. Moiety vectors D(:,4) and D(:,5) are the linearly independent components of the composite moiety vector L(:,4) above.

```
full(D(:,[4 5]))'
```

```
ans =

    0    0    1    0    1    0    0    0    0
    0    0    1    1    0    0    0    0    0
```

One disadvantage of decomposing moiety vectors is that it is difficult to keep track of which atoms belong to the decomposed moieties. We can, however, estimate the chemical formulae of the decomposed moieties using the elemental matrix for DAG. The elemental matrix is a numerical representation of the chemical formulae of metabolites in DAG.

```
[E,elements] = constructElementalMatrix(alternativeModel.metFormulas, ...
    alternativeModel.metCharges);
decomposedMoietyFormulae = estimateMoietyFormulae(D, E, elements);
decomposedMoietyFormulae(4:5)'
```

```
ans =

    'O2'    'O'
```

i.e., each decomposed moiety contains an oxygen atom.

## References

1. Haraldsdóttir, H.S., Thiele, I., Fleming, R.M. Comparative evaluation of open source software for mapping between metabolite identifiers in metabolic network reconstructions: application to Recon 2. *J. Cheminform* 6(1), 2 (2014).
2. Elizabeth Brunk, et al. Recon 3D—A Three-Dimensional View of Human Metabolism and Disease. *Submitted*.
3. Rahman, S.A., et al. Reaction Decoder Tool (RDT): extracting features from chemical reactions. *Bioinformatics* 32(13), 2065–2066 (2016).
4. Preclat et al. Comparative evaluation of atom mapping algorithms for balanced metabolic networks: application to Recon 3D. *J. Cheminform*, In 38 (2017).
5. Hulda S. Haraldsdóttir and Ronan M. T. Fleming. Identification of conserved moieties in metabolic networks by graph theoretical analysis of atom transition networks. *PLOS Comput. Biol.* 12(11) (2016).
6. Imari Famili and B. Ø. Palsson. The convex basis of the left null space of the stoichiometric matrix leads to the definition of metabolically meaningful pools. *Biophys. J.* 85(1):16–26 (2003).