

An automated Metabolite Identification Pipeline using Mass Spectral Trees



Theo Reijmers¹, Julio Peironcely^{1,2}, Miquel Rojas-Chertó¹, Piotr Kasper¹, Leon Coulier², Albert Tas², Rob Vreeken¹, Thomas Hankemeier¹

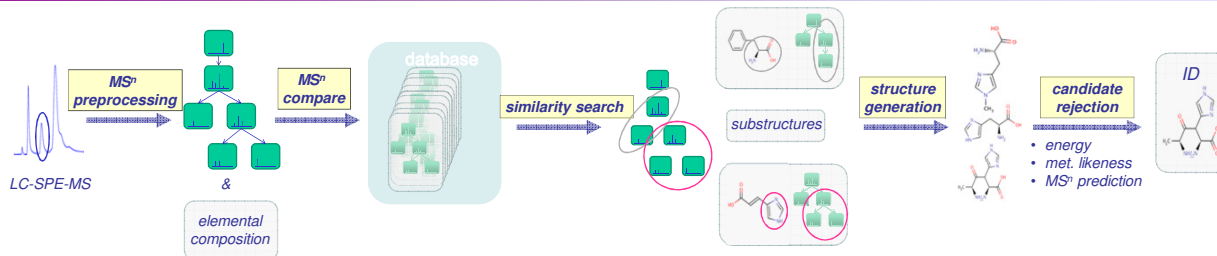
Netherlands Metabolomics Centre,
¹Analytical Biosciences, Leiden University, ²TNO Quality of Life, Zeist

t.reijmers@facdr.leidenuniv.nl
 Netherlands
 Metabolomics Centre

Introduction

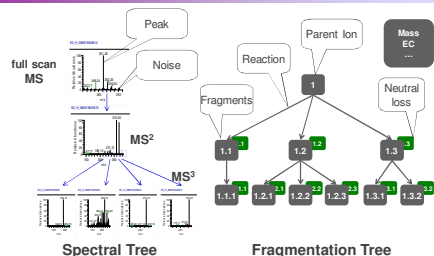
Identifying metabolites has been reported as one of the major bottlenecks in metabolomics. In part, this is due to the absence of good computational tools to automate metabolite identification. To address this issue, we have developed mathematical tools to process and compare multi-stage mass spectrometry data (MSⁿ), in order to extract as much as possible information from the fragmentation trees. In addition, candidate structures are generated computationally and filters are used to reject improbable chemical structures. This poster presents the integrated use of these tools in a pipeline fashion to identify metabolites in human urine.

Approach / Metabolite Identification Pipeline



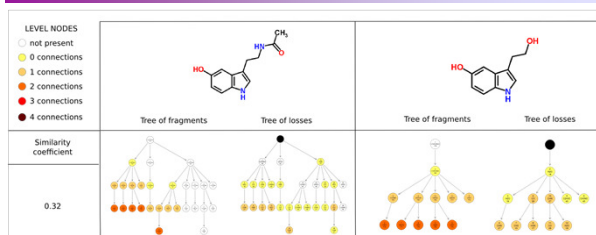
- MSⁿ spectral trees are processed into fragmentation trees using the MEF (1) tool. The nodes of the fragmentation and neutral loss tree are annotated with unique elemental compositions.
- Comparing the fragmentation tree with trees in a reference MSⁿ database (2) returns an identical tree (identity search) or multiple similar trees (similarity search). From the most similar trees maximum common substructures are extracted.
- The Open Structure Generator (OMG) (3) generates for the unknown metabolite all possible chemical structures for a given elemental composition and (multiple) substructures.
- This number of structures on the list of candidate structures is reduced after applying an internal energy and a metabolite likeness (4) filter. The fragmentation prediction tool MetFrag (5) further reduces the list by comparing the observed with predicted fragmentation spectra.

MSⁿ Preprocessing



From a pooled urine sample, mass spectral trees are acquired for 30 compounds with unknown identity. After preprocessing with the MEF tool, 30 fragmentation and neutral loss trees are obtained with unique elemental compositions assigned to the nodes.

MSⁿ Comparison



Similarity search in the NMC reference MSⁿ database of each unknown fragmentation tree yielded the following results: 9 trees are 100% similar to reference trees in database (identity search), 9 trees return similar trees from database (similarity search) and 12 trees do not show similar trees (similarity value < 10%).

Structure Generation

	Glycine	Phenylalanine	Malic acid	D-Cysteine	p-Cresol sulfate
Elemental Composition	C2H5NO2	C9H11NO2	C4H6O5	C3H7NO2S	C7H8O3S
# Output Molecules	84	277,810,163	8,070	3,838	10,203,389
1 Fragment	6	4,037,499	1,601	100	19,940
2 Fragments	93,137				948
3 Fragments	584				278

Generate
 Keep molecules if canonical augmentation
 All non-duplicated molecules
 CDK Nauty

Multiple similar trees are found for 6 of the 9 trees. Maximum common substructure (MCSS) of these similar hits is used as additional input for structure generator. Remaining 3 trees only return 1 similar hit in the database so no MCSS can be calculated.

Candidate Rejection

	146.08	181.06	227.08	151.08	170.04	183.05	196.06	262.04	185.09
Elemental Composition	C6H11 NO3	C8H8 N2O3	C8H10 N4O4	C9H10 O2	C7H7 N4O2	C6H6 N4O3	C9H9 NO4	C9H13 NO4P2	C8H12 N2O3
# hits	1	1	1	2	2	2	2	3	21
Candidate Structures	575,709	Billions	Billions	6.8M	150M	Billions	Billions	Billions	Billions
MCSS	NO	NO	NO						
Putative Identities MCSS				82	24K	4	8	9	Billions
Putative Identities Filtered				8	1K	4	5	7	

HMDB ZINC DB
 MDL Public Keys
 Random Forest Classifier
 Metabolite Likeness
 MetFrag

The number of generated candidate structures is further reduced after setting an internal energy threshold, applying a metabolite-likeness filter and predicting fragmentations using MetFrag.

Conclusions

- Presented pipeline, consisting of a number of separate tools, facilitates experts for (de-novo) identification of metabolites.
- Establishing a comprehensive MSⁿ database (for extracting substructures) is crucial.
- Visit the MetiTree (6) web-application, www.metitree.nl, for part of this pipeline.

1. Rojas-Chertó, M. et al. *Bioinformatics* 27, 2011, 2376-2383.
 2. Rojas-Chertó, M. et al. *Anal. Chem.* 84, 2012, 5524-5534.
 3. Peironcely, J.E. et al. *J. Cheminformatics*, 4, 2012, 21.

4. Peironcely, J.E. et al. *PLoS ONE*, 6, 2011, e28966.
 5. Wolf, S. et al. *BMC Bioinformatics*, 11, 2010, 148.
 6. Rojas-Chertó, M. et al. *Bioinformatics*, doi: 10.1093/bioinformatics/bts486.

