

COMPREHENSIVE SHARING AND MAPPING OF RP-LC RETENTION TIME INFORMATION

Jan Stanstrup, Urška Vrhovšek

Metabolomics technological platform, Department of Food Quality and Nutrition, Research and Innovation Centre, Fondazione Edmund Mach, San Michele all'Adige (TN), Italy. Email: jan.stanstrup@fmach.it

Objectives

- ❑ To build a user-driven database of retention times (RTs) of compounds.
- ❑ To map the RT of compounds between different chromatographic systems.
- ❑ To use these mappings to predict the RT of compounds not experimentally determined in the system of interest.

Background

Databases of experimental LC-MS data have been developed with great success with regards to compound fragmentation and these databases have recently been used by automated tools to assist compound identification.

But utilizing *only* the fragmentation means disregarding one, equally important, half of the information in LC-MS.

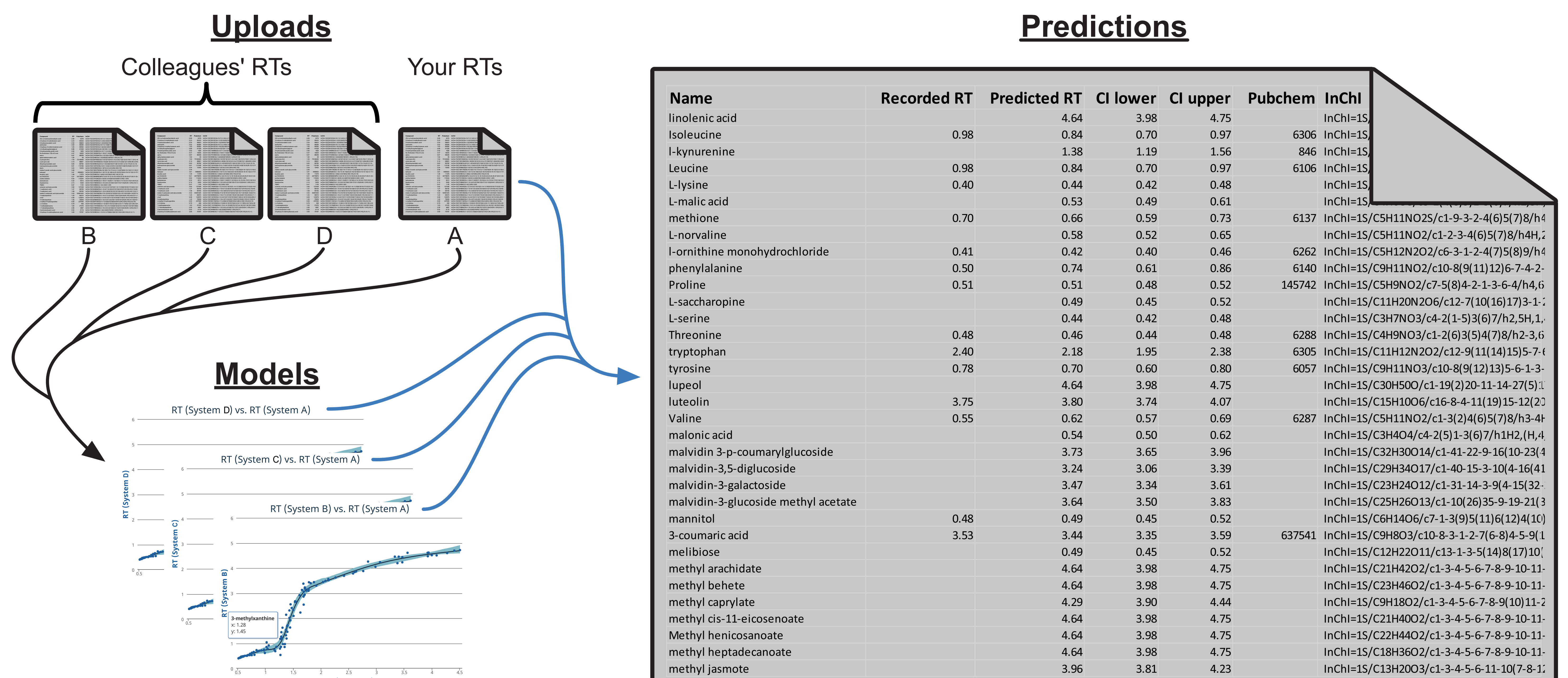
We therefore build a database of compounds' RTs and use this database to predict the RT of compounds in systems where they have not been experimentally determined.

Conclusions

- ✓ The tool is available at www.predret.com. After uploading RTs of compounds in your own systems predicted RTs for other compounds become available.
- ✓ Identification efforts can be reserved for compounds that cannot be excluded by the RT predictions.
- ✓ Community support is required to expand the database.

Results

Building models between all chromatographic systems in the database allows us to predict the RT of compounds in systems where they have not been experimentally determined. The accuracy is dependent on the number of compounds measured in both systems used in the mapping step and on the slope of the curve at the point of prediction. With the current small database (< 1000 compounds) it was possible to predict up to 350 RTs with a median error between 0.04 and 0.19 min depending on the system.



Methods

The prediction tool is made available as a web application at www.predret.com. The user will upload a spreadsheet with RTs of compounds measured in their system along with molecular identifiers such as PubChem CIDs or InChIs.

The system will then map the RT of compounds between systems by building monotonically increasing smooth generalized additive models between RTs experimentally determined in two different systems. This model can then be used to predict the RT of a compound if the RT is known in one system, but not in the other. This mapping and prediction can be done between all systems added to the database.

COMPREHENSIVE SHARING AND MAPPING OF RPLC RETENTION TIME INFORMATION

Jan Stanstrup*, Urška Vrhovšek*

*Department of Food Quality and Nutrition, Research and Innovation Centre, Fondazione Edmund Mach (FEM), San Michele all'Adige (TN), Italy.

Data repositories have been developed with great success to the benefit of the scientific community with regards to the fragmentation of compounds^{1,2}. In addition, recently efforts have been made to use automated tools, often assisted by these databases of experimental data, to assist compound identification^{3–5}. However, these tools and databases only focus on one aspect of the experimental data: the fragmentation of the compounds formed in mass spectrometers. But utilizing the fragmentation is only using half the available information.

In GC analysis retention indexes are routinely used to make different systems comparable but for LC systems there are currently no coordinated focused efforts to share and exploit information regarding the retention time (RT) of compounds. The reason RT information has been neglected in LC systems is that the RT is specific to a specific chromatographic setup and there exists no agreed upon RT references.

We have therefore sought to rectify this by building a database of compounds' RTs. With this database we are able to map the RT of compounds between systems. RTs, experimentally determined in two different systems, of a number of compounds is used to build monotonically increasing smooth generalized additive models between the RTs in the two systems using the mgcv package⁶ for R. This model can then be used to predict the RT of a compound if the RT is known in one system, but not in the other. Building these models between all chromatographic systems in the database thus allows us to predict the RT of a high number of compounds in systems where they have not been experimentally determined.

The tool is completely web-based and available at www.predret.org. On the website it is possible to upload a spreadsheet containing RT information and subsequently download predicted RTs for other compounds based on the data available in the database.

We believe that this tool will greatly help the identification process since compounds that are not compatible with the observed RT can be disregarded. Confirmatory experiments can then be reserved for compounds that *could* have the observed RT. This will allow researchers to complete the feature annotation and compound identification process in a faster and more rational manner and thus save time and resources, both monetary and environmental.

References

1. Horai, H. *et al.* MassBank: a public repository for sharing mass spectral data for life sciences. *J. Mass Spectrom.* **45**, 703–714 (2010).
2. Wishart, D. S. *et al.* HMDB: a knowledgebase for the human metabolome. *Nucleic Acids Res.* **37**, D603–D610 (2009).
3. Gerlich, M. & Neumann, S. MetFusion: integration of compound identification strategies. *J. Mass Spectrom.* **48**, 291–298 (2013).
4. Wolf, S., Schmidt, S., Müller-Hannemann, M. & Neumann, S. In silico fragmentation for computer assisted identification of metabolite mass spectra. *BMC Bioinformatics* **11**, 148 (2010).
5. Peironcelly, J. E. *et al.* Automated Pipeline for De Novo Metabolite Identification Using Mass-Spectrometry-Based Metabolomics. *Anal. Chem.* **85**, 3576–3583 (2013).
6. Wood, S. N. Fast stable restricted maximum likelihood and marginal likelihood estimation of semiparametric generalized linear models. *J. R. Stat. Soc. Ser. B Stat. Methodol.* **73**, 3–36 (2011).