

A data analysis pipeline integrating ion mobility and high-resolution mass spectrometry for non-target screening in environmental studies

**Julien Sade, Sabrina Guérin, Stéphane Mottelet,
Vincent Rocher, Régis Moilleron and Julien Le Roux**

Laboratoire Eau Environnement et Systèmes Urbains (LEESU)
Univ Paris Est Créteil, Créteil, France

julien.sade@u-pec.fr

SETAC Europe 2024, 5-9 May, Seville, Spain



What are the various HRMS strategies for characterising contaminants?

Technique	Targeted	Suspect	Non-target
Question	Are compounds x, y, & z present in this sample?	Which compounds of a defined list are present in this sample?	Which compounds are present in this sample?
Type	Known-knowns	Known-unknowns	Known-unknowns & unknown-unknowns

Confidence levels in contaminants annotation

Technique	Targeted	Suspect	Non-target
Question	Are compounds x, y, & z present in this sample?	Which compounds of a defined list are present in this sample?	Which compounds are present in this sample?
Type	Known-knowns	Known-unknowns	Known-unknowns & unknown-unknowns



1

Level
data requirements

MS, MS/MS, RT,
Reference standard

Confidence levels in contaminants annotation

Technique	Targeted	Suspect	Non-target
Question	Are compounds x, y, & z present in this sample?	Which compounds of a defined list are present in this sample?	Which compounds are present in this sample?
Type	Known-knowns	Known-unknowns	Known-unknowns & unknown-unknowns
Level	1		2 & 3
data requirements	MS, MS/MS, RT, Reference standard		MS, MS/MS, Library MS/MS MS, isotope/adduct

Confidence levels in contaminants annotation

Technique	Targeted	Suspect	Non-target
Question	Are compounds x, y, & z present in this sample?	Which compounds of a defined list are present in this sample?	Which compounds are present in this sample?
Type	Known-knowns	Known-unknowns	Known-unknowns & unknown-unknowns

Level

1

2 & 3

4 & 5

data requirements

MS, MS/MS, RT,
Reference standard

MS, MS/MS, Library MS/MS

MS, isotope/adduct

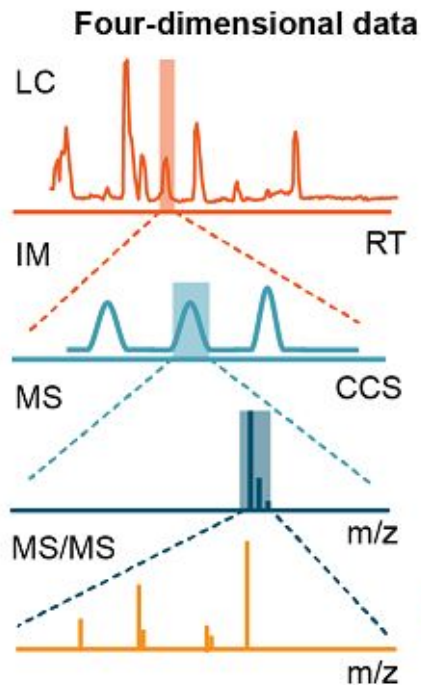
Potential false positives

Confidence levels in contaminants annotation: impact of CCS data

Vion[®] IMS QTof



Waters
THE SCIENCE OF WHAT'S POSSIBLE.™

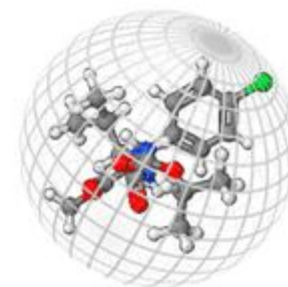
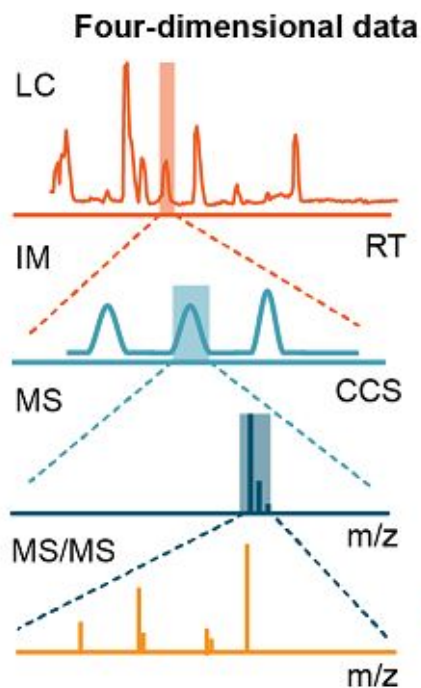


Confidence levels in contaminants annotation: impact of CCS data

Vion[®] IMS QTof



Waters
THE SCIENCE OF WHAT'S POSSIBLE.™

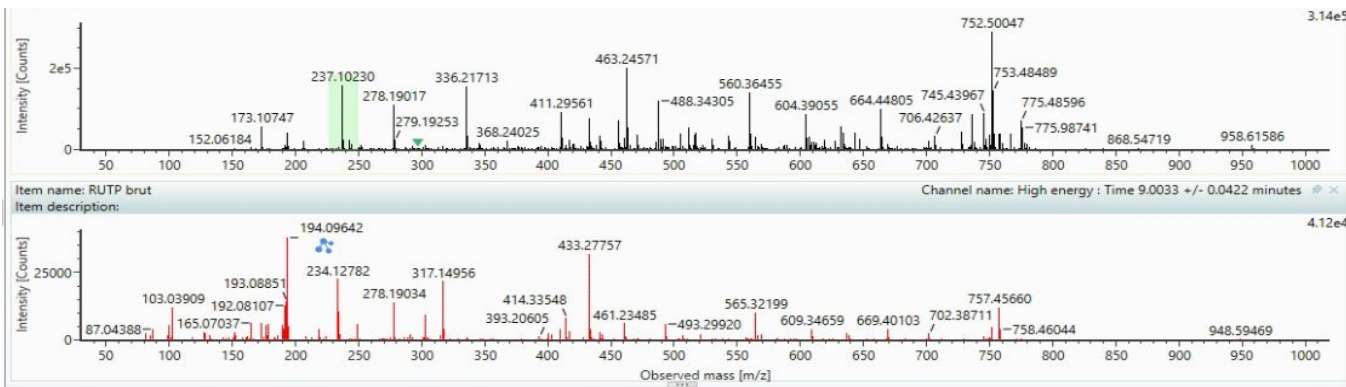


Why IMS?

- CCS : an additional identification parameter
- Spectral cleaning
- Reduction of false positive rate

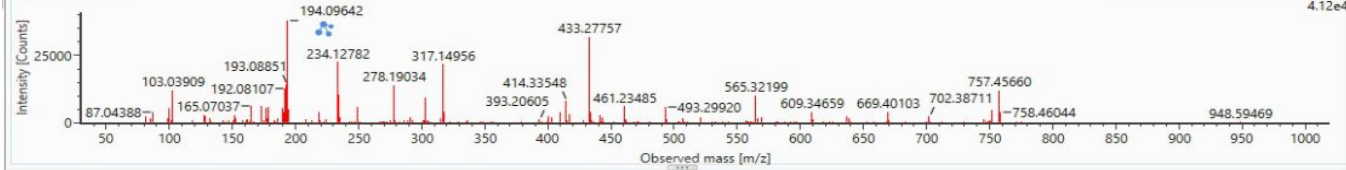
Confidence levels in contaminants annotation: impact of CCS data

Low energy

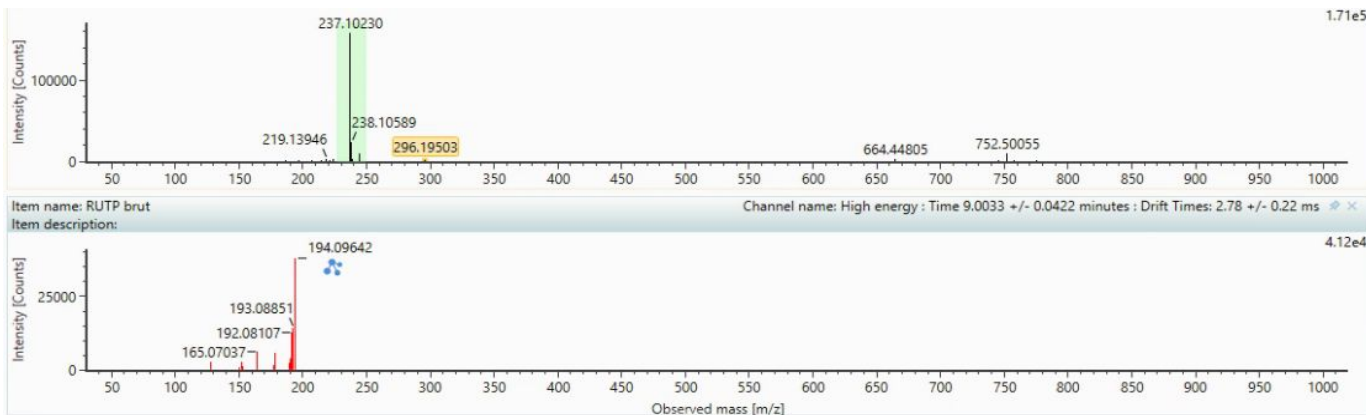


NO IMS

High energy



Low energy

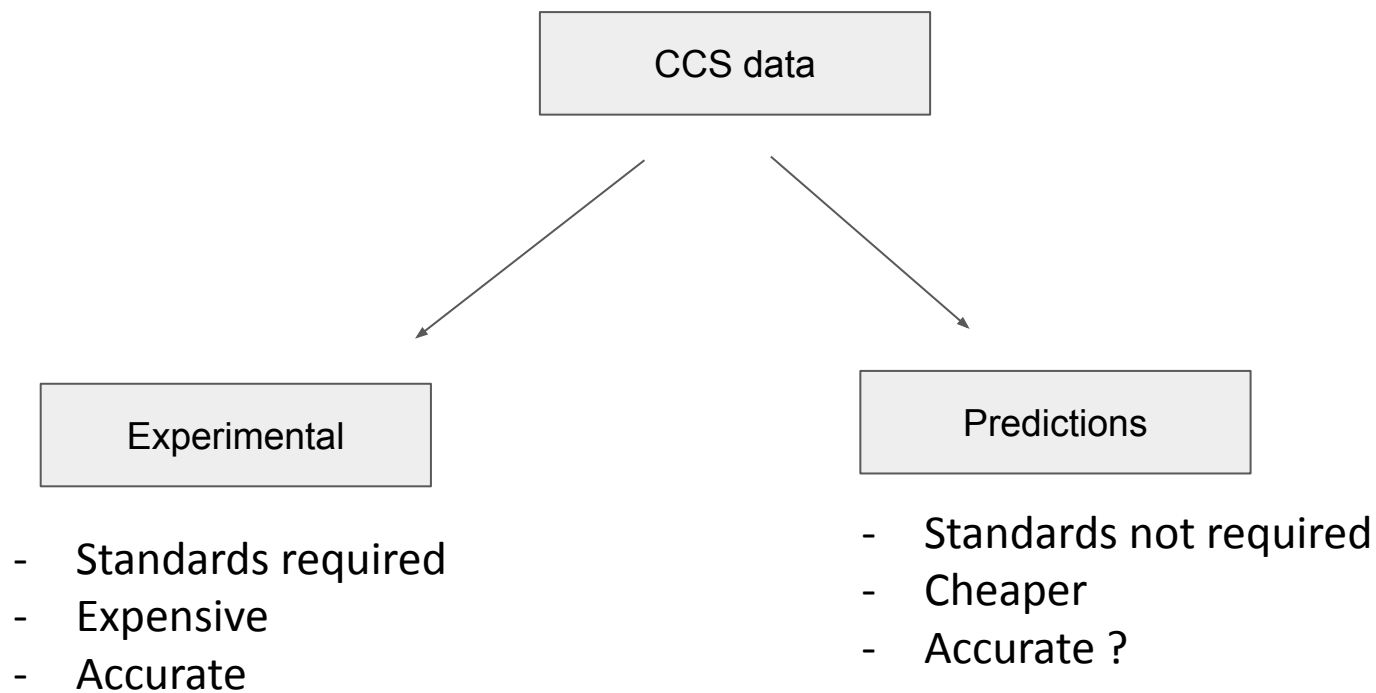


IMS

High energy

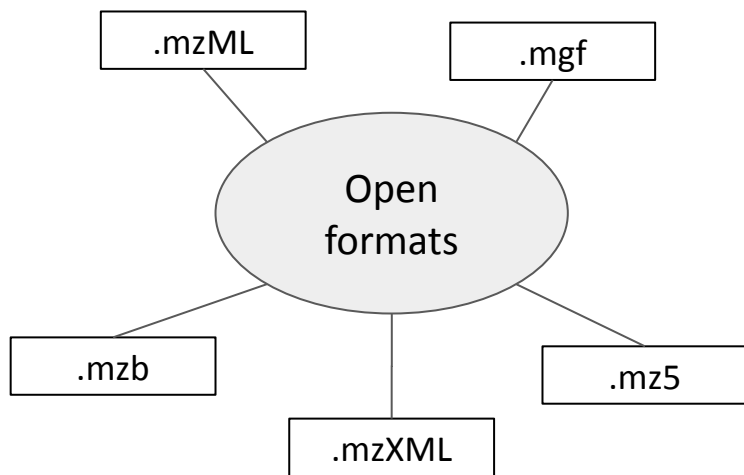
Better spectral matching with MS libraries!

CCS data for NTS



CCSbase

Challenges with formats and softwares including ion mobility (IMS)



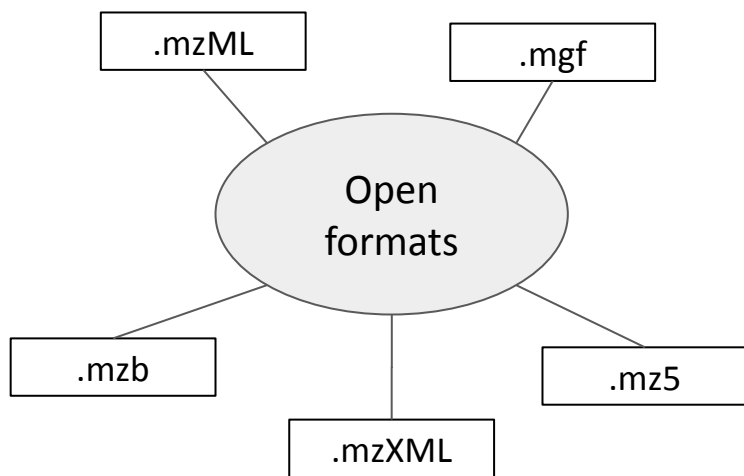
Issues with proprietary formats: lack of interoperability, closed ecosystems, unknown algorithms...

**Data conversion takes a long time
~ 20 minutes for 1 replicate**

File size ~ 2-10 GB per file

**Most tools do not read or process open-format data including IMS (yet)
Mzmine, MS-dial, XCMS, patRoan...**

Challenges with formats and softwares including ion mobility (IMS)



Issues with proprietary formats: lack of interoperability, closed ecosystems, unknown algorithms...

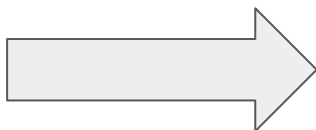
Data conversion takes a long time
~ 20 minutes for 1 replicate

File size ~ 2-10 GB per file

Most tools do not read or process open-format data including IMS (yet)
Mzmine, MS-dial, XCMS, patRoan...

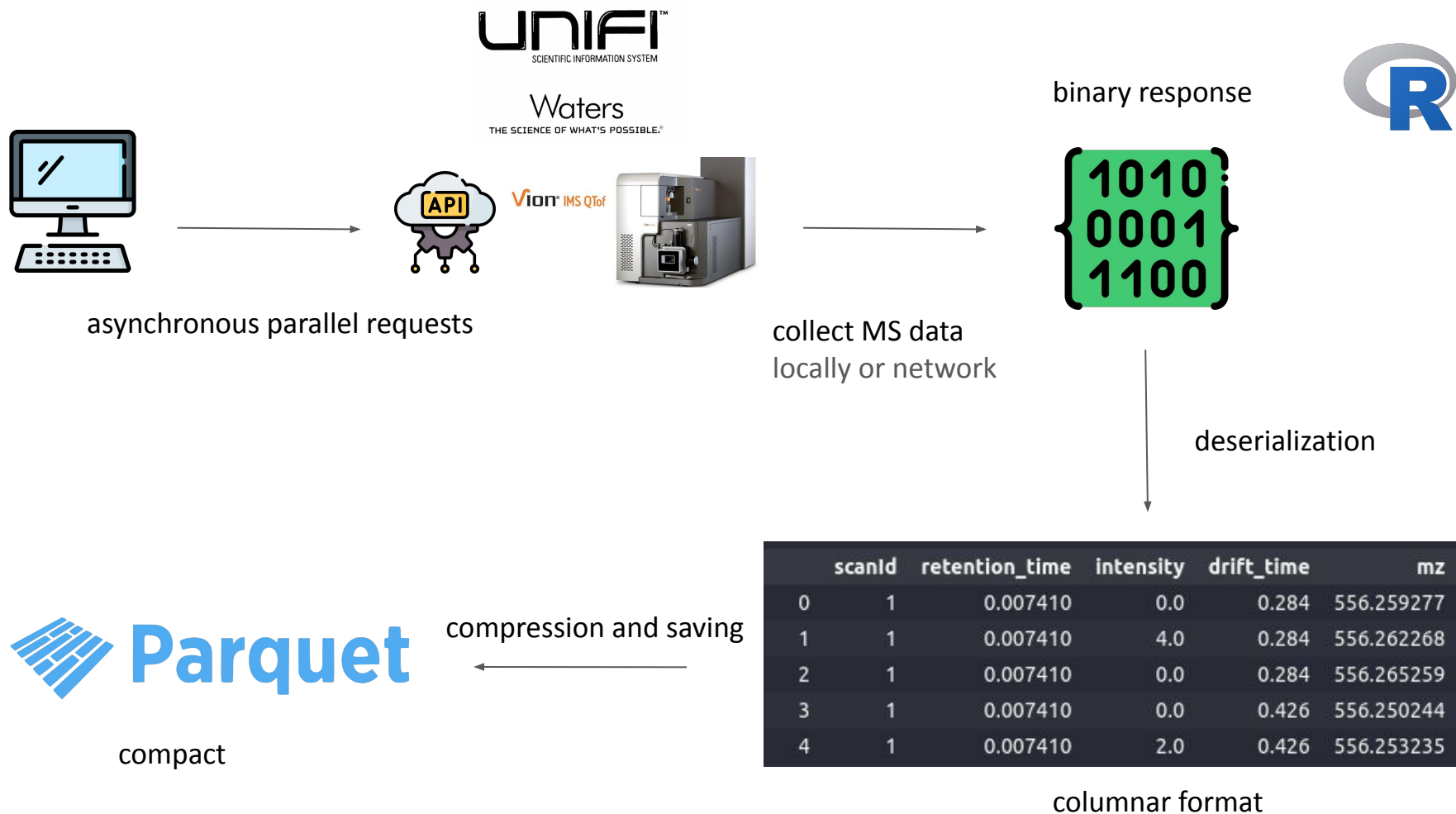
Our solution: arcMS

- R package
- Free, open-source
- Including IMS
- Under our control
- parquet file format



arcMS workflow

A UNIFI to .parquet conversion package



arcMS performance

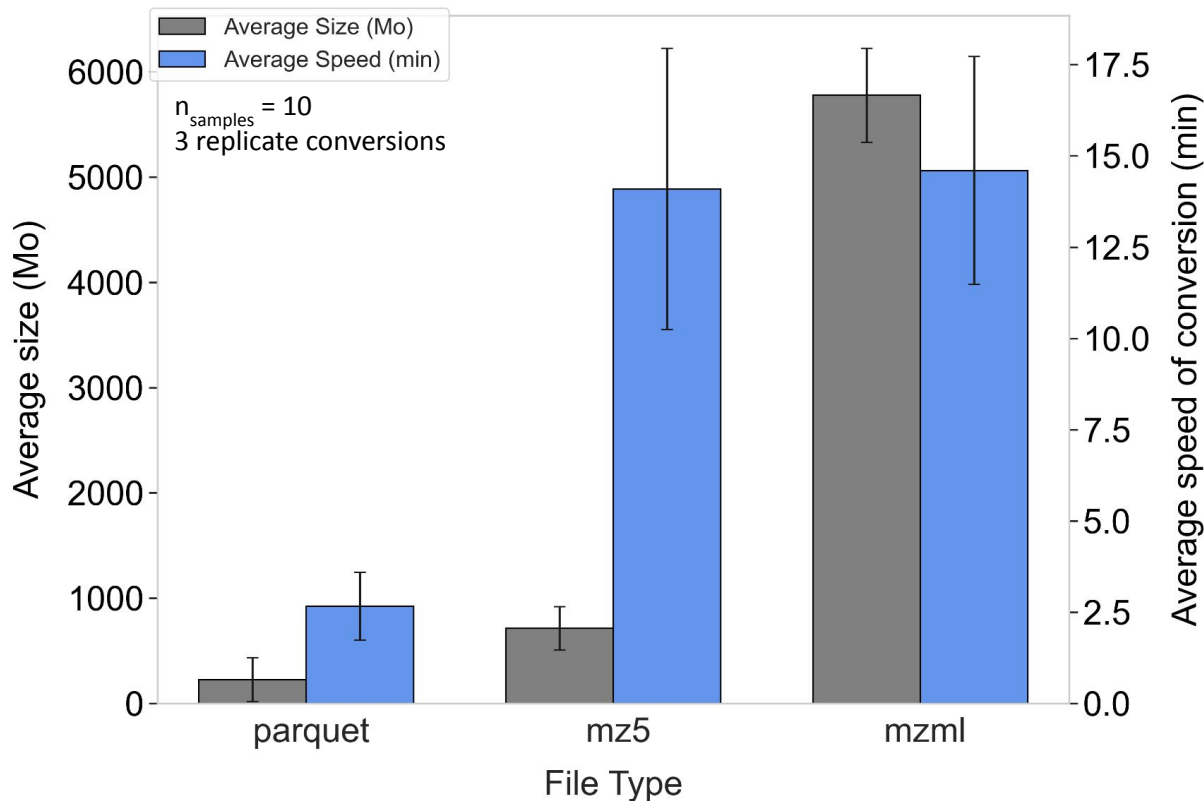
Data collection and conversion speed

File size

arcMS vs



arcMS vs mzML and mz5



- More compact
- Faster collection

arcMS philosophy

Opening and visualization of raw data

Columnar format

scanId	retention_time	intensity	drift_time	mz	
0	1	0.007410	0.0	0.284	556.259277
1	1	0.007410	4.0	0.284	556.262268
2	1	0.007410	0.0	0.284	556.265259
3	1	0.007410	0.0	0.426	556.250244
4	1	0.007410	2.0	0.426	556.253235



- Open source
- Modular, extensible, easy to manipulate, fast to read
- Compatible with several programming languages (R, python, Java, C++, C#,etc...)
- Compatible with the most commonly used libraries
 - For visualization : plotly, matplotlib, ggplot...
 - For data analysis : pandas, data.table, scikit-learn, numpy...

arcMS philosophy

Opening and visualization of raw data

Columnar format

	scanId	retention_time	intensity	drift_time	mz
0	1	0.007410	0.0	0.284	556.259277
1	1	0.007410	4.0	0.284	556.262268
2	1	0.007410	0.0	0.284	556.265259
3	1	0.007410	0.0	0.426	556.250244
4	1	0.007410	2.0	0.426	556.253235



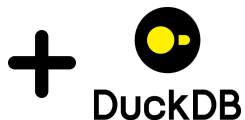
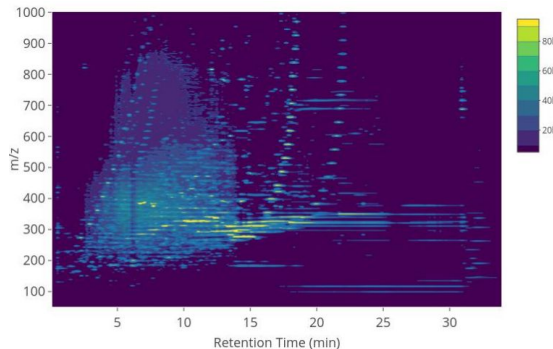
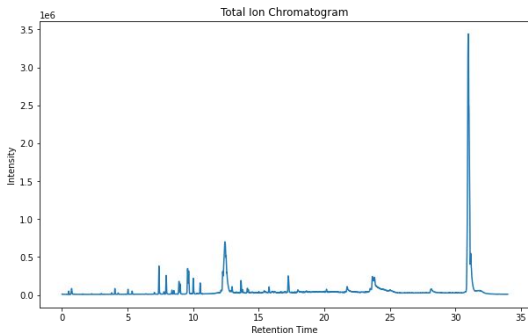
matplotlib



plotly

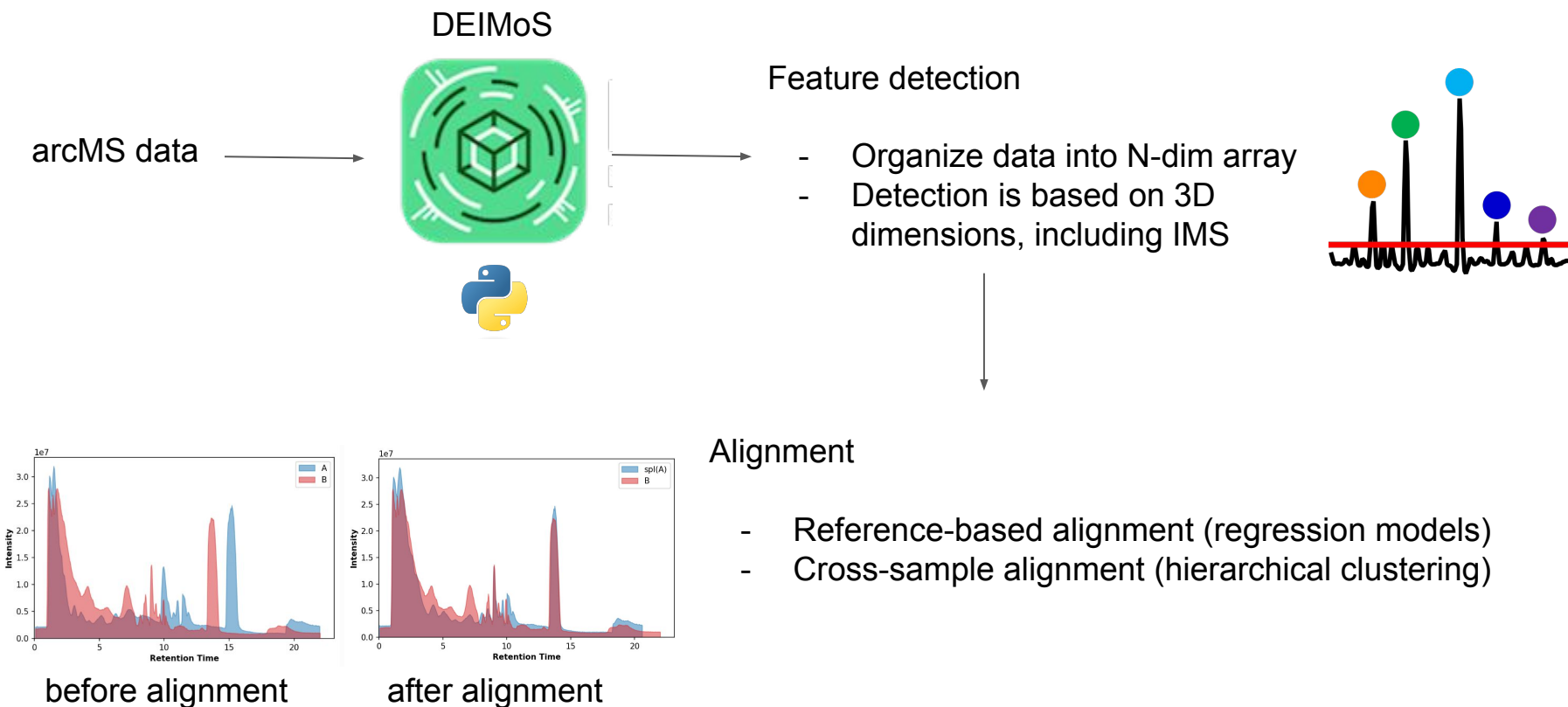
bokeh

- Open source
- Modular, extensible, easy to manipulate, fast to read
- Compatible with several programming languages (R, python, Java, C++, C#,etc...)
- Compatible with the most commonly used libraries
 - For visualization : plotly, matplotlib, ggplot...
 - For data analysis : pandas, data.table, scikit-learn, numpy...



for on-disk operations! (low RAM usage)

arcMS integration with DEIMoS: a package to process LC-IMS-MS data



Applications of the data pipeline

- **Observatory of Paris wastewater (Greater Paris sanitation authority):**

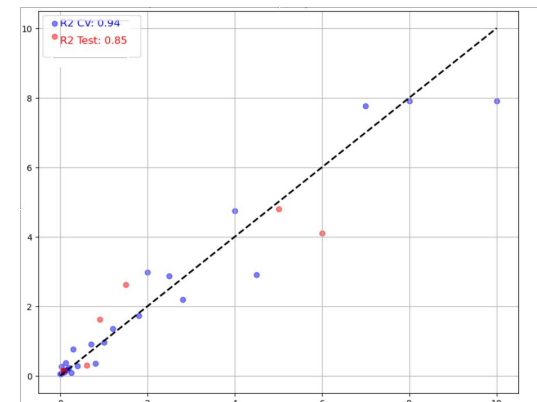
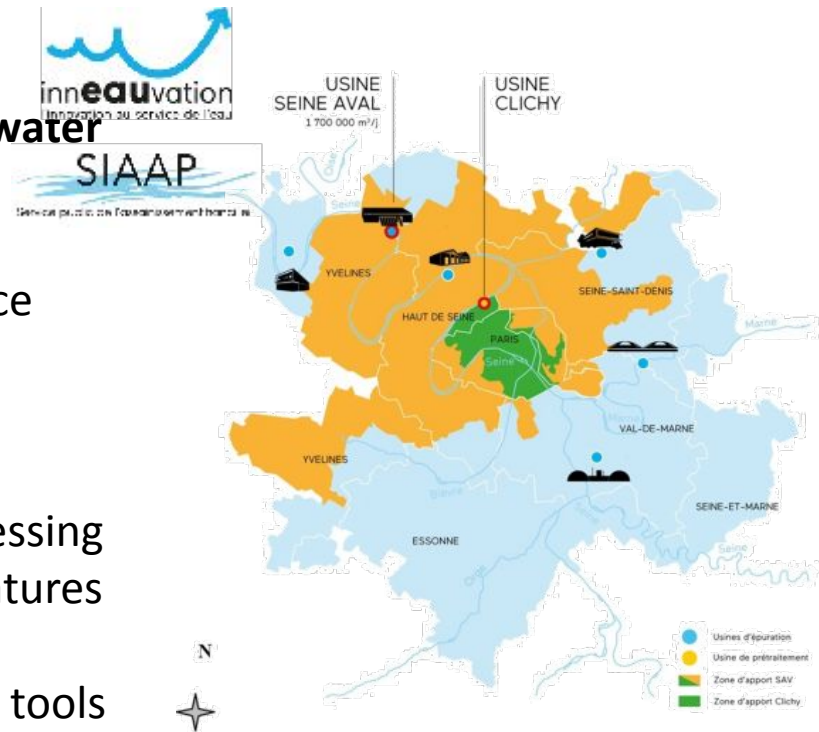
- Goal: 10-year analyses
- achieving a 15x reduction in storage space
- Cost savings

- **Data processing and statistics:**

- Development of an application for processing spectral data including IMS features (peak picking, alignment...)
- Visualization tools (chromatograms, contour plot...)
- Statistical analysis (descriptive statistics, unsupervised models (PCA), and supervised models)

- **Contamination model:**

- Development of predictive models for contamination of surface water samples



Annotation with CCS

urban wastewater discharges



retention_time	mz	intensity	drift_time	CCS
18.215646	705.409058	717.0	0.071	63.983055
17.262668	485.280487	9355.0	0.071	66.071616
17.580382	485.277679	1835.0	0.071	66.071648
18.198939	705.409058	946.0	0.213	66.197677
17.287745	509.290314	1073.0	0.142	66.925016
...

- mz tolerance 5 ppm
- CCS tolerance 6%
- All IMS

Annotation with experimental CCS data (11K)

retention_time	mz	intensity	drift_time	CCS	Matching_Molecule
16.452719	399.251129	25720.756255	6.603	206.414514	tris(2-butoxyethyl) phosphate
4.790156	371.227448	25374.575866	5.609	181.595989	Nonoxynol-9_met115
10.366448	429.239502	24420.327584	6.603	205.519580	IRBESARTAN
4.840312	416.211487	18761.324646	6.035	191.086060	Famprofazone
4.028913	195.086945	16011.000000	3.550	139.806489	Caffeine
...



CCSbase

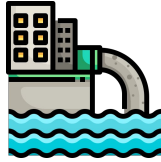
500 matches

Level 4

mz + CSS

Annotation with CCS

urban wastewater discharges



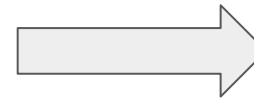
retention_time	mz	intensity	drift_time	CCS
18.215646	705.409058	717.0	0.071	63.983055
17.262668	485.280487	9355.0	0.071	66.071616
17.580382	485.277679	1835.0	0.071	66.071648
18.198939	705.409058	946.0	0.213	66.197677
17.287745	509.290314	1073.0	0.142	66.925016
...

- mz tolerance 5 ppm
- CCS tolerance 6%
- All IMS

Annotation with experimental CCS data (11K)

retention_time	mz	intensity	drift_time	CCS	Matching_Molecule
16.452719	399.251129	25720.756255	6.603	206.414514	tris(2-butoxyethyl) phosphate
4.790156	371.227448	25374.575866	5.609	181.595989	Nonoxynol-9_met115
10.366448	429.239502	24420.327584	6.603	205.519580	IRBESARTAN
4.840312	416.211487	18761.324646	6.035	191.086060	Famprofazone
4.028913	195.086945	16011.000000	3.550	139.806489	Caffeine
...

MS/MS confirmation



47 matches

Level 3



500 matches

Level 4



CCSbase

mz + CSS

Conclusions

- Fast automated pipeline for HRMS data including IMS
- Could be adapted to other file formats
- Versatility of Parquet format (space storage, reading and handling)

UNIFI to .parquet conversion package:

<https://github.com/leesulab/arcMS>

<https://leesulab.github.io/arcMS/>

leedulab / arcMS Public

<> Code Issues Pull requests Actions Projects Security Insights

main 4 Branches 5 Tags

Search Go to file

Code

Files

File/Folder	Description	Time
.github	removing branch pkgdown for github actions	3 months ago
R	changing Unifi to UNIFI in whole doc	3 months ago
inst	changing Unifi to UNIFI in whole doc	3 months ago
man	changing Unifi to UNIFI in whole doc	3 months ago
tests	parquetMS renamed to arcMS	3 months ago
vignettes	finalizing vignette api-configuration, adding link in READ...	3 months ago
.Rbuildignore	pkgdown init	3 months ago
.gitignore	starting vignette	3 months ago
DESCRIPTION	preparing v1.1.0	3 months ago
LICENSE	parquetMS renamed to arcMS	3 months ago
LICENSE.md	parquetMS renamed to arcMS	3 months ago
NAMESPACE	Connection params taken from object in environment, n...	4 months ago
NEWS.md	preparing v1.1.0	3 months ago
README.Rmd	finalizing vignette api-configuration, adding link in READ...	3 months ago
README.md	finalizing vignette api-configuration, adding link in READ...	3 months ago

About

Mass spectrometry data converter from UNIFI to Parquet and HDF5 formats

leedulab.github.io/arcMS/

converter r mass-spectrometry non-target

Readme

Unknown, MIT licenses found

Activity

Custom properties

1 star

1 watching

0 forks

Report repository

Releases

5 tags

Packages

No packages published

Contributors 2

arcMS 1.1.0 Reference Articles Changelog

arcMS

arcMS can convert HDMS⁵ data acquired with Waters UNIFI to tabular format for use in R or Python, with a small filesize when saved on disk. test

Two output data file formats can be obtained:

- the [Apache Parquet](#) format for minimal filesize and fast access. Two files are produced: one for MS data, one for metadata.
- the [HDF5](#) format with all data and metadata in one file, fast access but larger filesize.

arcMS stands for *accessible, rapid and compact*, and is also based on the french word *arc*, which means *bow*, to emphasize that it is compatible with the [Apache Arrow library](#).

Installation

You can install `arcMS` in R with the following command:

```
install.packages("pak")
pak::pkg_install("leedulab/arcMS")
```

To use the HDF5 format, the `rhdf5` package needs to be installed:

```
pak::pkg_install("rhdf5")
```

Usage

First load the package:

```
library("arcMS")
```

Links

[Browse source code](#)

License

[Full license](#)

[MIT](#) + file [LICENSE](#)

Citation

[Citing arcMS](#)

Developers

Julien Le Roux
Author, maintainer

Julien Sade
Author

Acknowledgements



¡ Muchas Gracias !

