

# How simple can mass-spec files get?

Databases are a speedy, small, and simple solution for MS data storage and access

William Kumler and Anitra E. Ingalls

wkumler@uw.edu • [github.com/wkumler/mzsql](https://github.com/wkumler/mzsql)

University of Washington, School of Oceanography, Seattle, USA



Column-based tidy format is simple and intuitive for all MS types

Convert  $m/z$  and intensity tuples into database columns

Pair with other separation data

- Retention time (liquid chromatography)
- Drift time (ion mobility)
- X/Y coordinate (imaging MS)

Link with MS<sup>n</sup> data via scan number

Pair with filename to aggregate multifile

- Optimized *across* files
- Metadata tables saved alongside

Table: MS1

filename	scan_num	rt	mz	int
Source file	Scan number	Retention time	m/z ratio	Intensity

Smp_A	1	0.10	60.0452	6618
Smp_A	1	0.10	60.0532	2657

...millions of additional entries...

Smp_Z	1385	22.35	60.0456	158084
Smp_Z	1385	22.35	60.0531	4673

Table: MS2

filename	scan_num	prescan	rt	fragmz	prempz	int
Source file	Scan number	Precursor scan	Retention time	Fragment m/z	Precursor m/z	Intensity
Smp_A	2	1	0.12	51.0238	241.0894	36104
Smp_A	2	1	0.12	53.0394	241.0894	243165

Ion chromatogram extraction:

SELECT \* FROM MS1 WHERE mz BETWEEN min AND max

Retention time range subset:

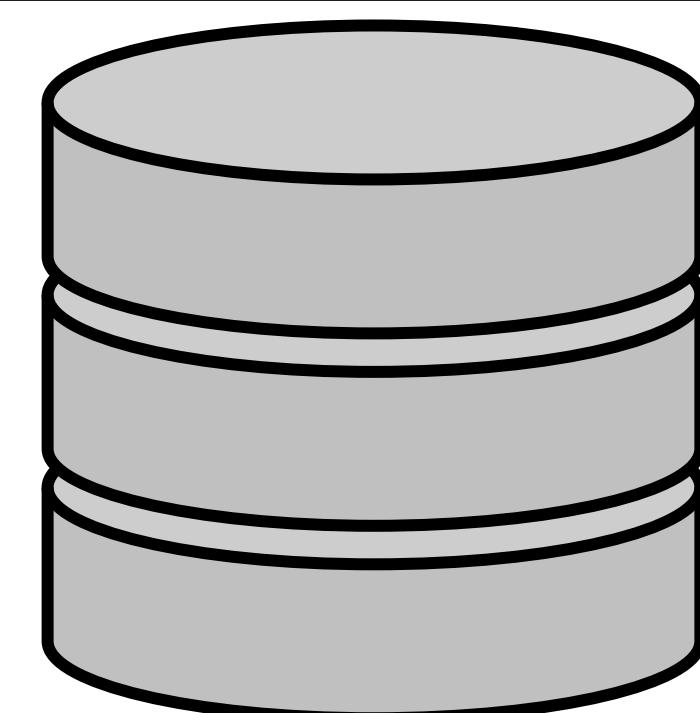
SELECT \* FROM MS1 WHERE rt BETWEEN min AND max

Fragment search:

SELECT \* FROM MS2 WHERE fragmz BETWEEN min AND max

Precursor search:

SELECT \* FROM MS2 WHERE prempz BETWEEN min AND max



Database  
.sqlite  
.duckdb  
.parquet

## Problem: very large in memory! Solution: simple databases

### Databases vs existing MS file options

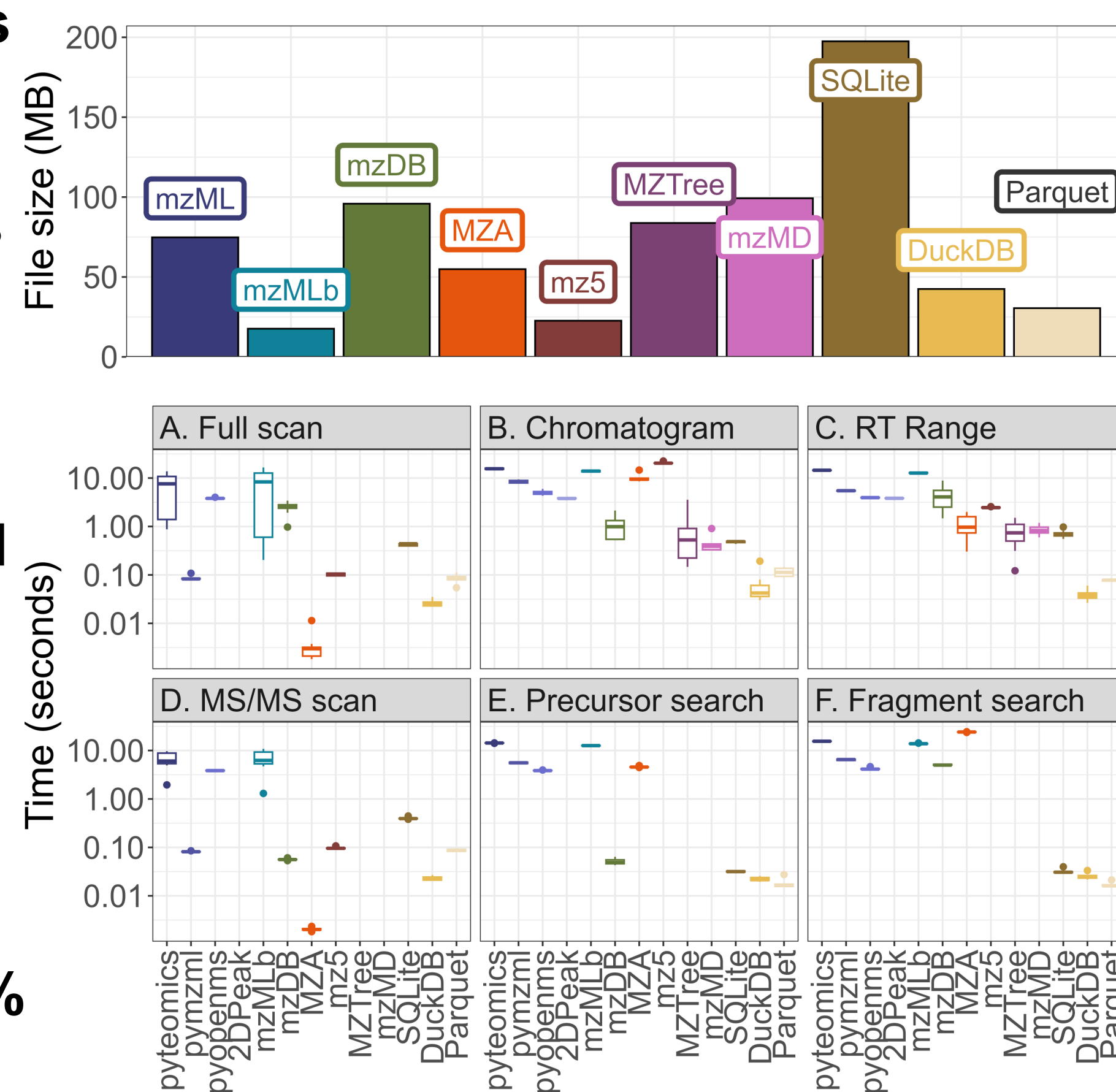
We compared 13 different mass-spectrometry file types in Python using six metrics

- Single MS<sup>1</sup> scan
- Ion chromatogram
- Consecutive MS<sup>1</sup> scans
- Single MS<sup>2</sup> scan
- Precursor  $m/z$  search
- Fragment  $m/z$  search

Databases outperformed on all metrics other than known scan extraction

- Often by 10-100x

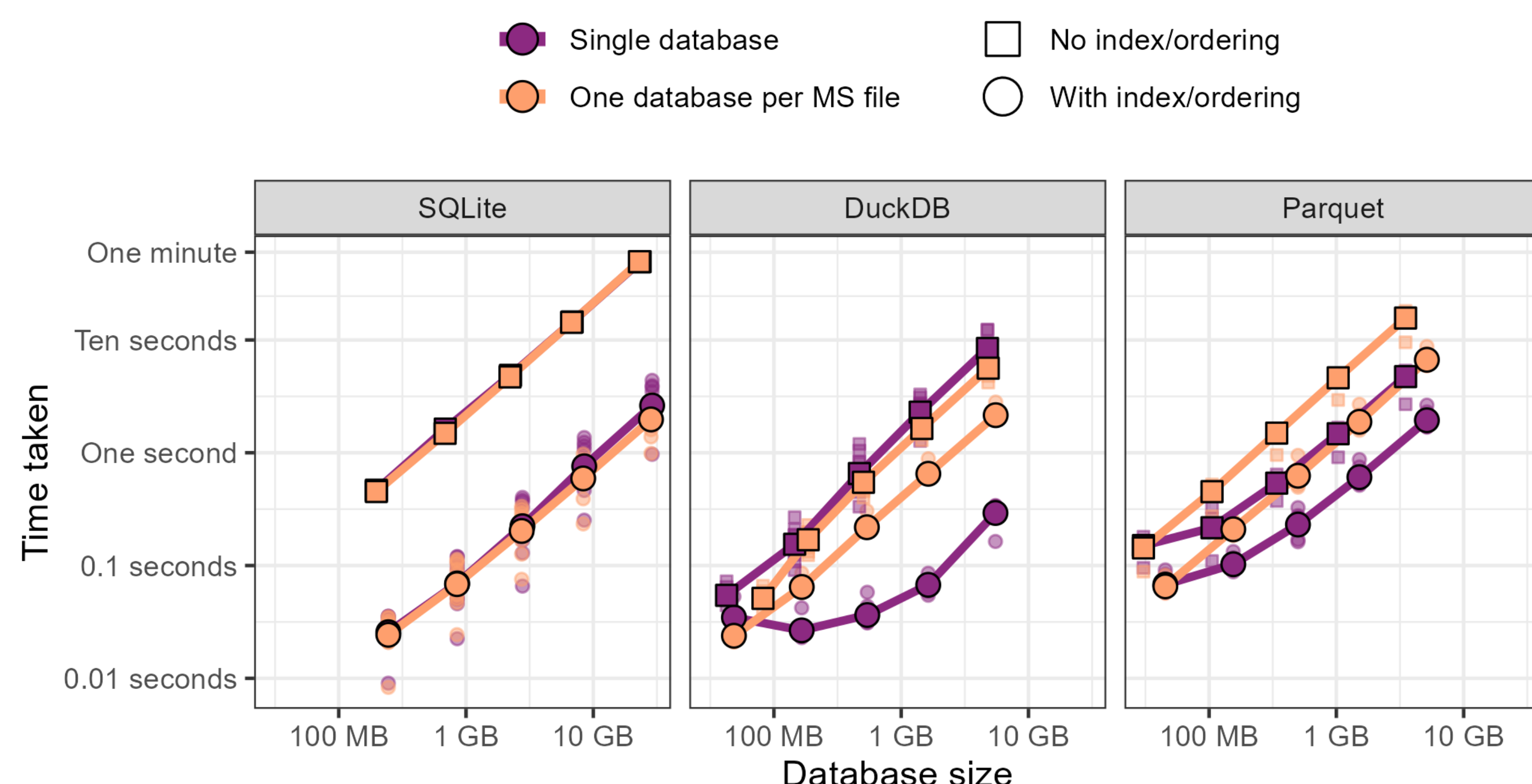
Parquet and DuckDB methods also reduced the on-disk file size ~50%



### Optimization: index / order / aggregate

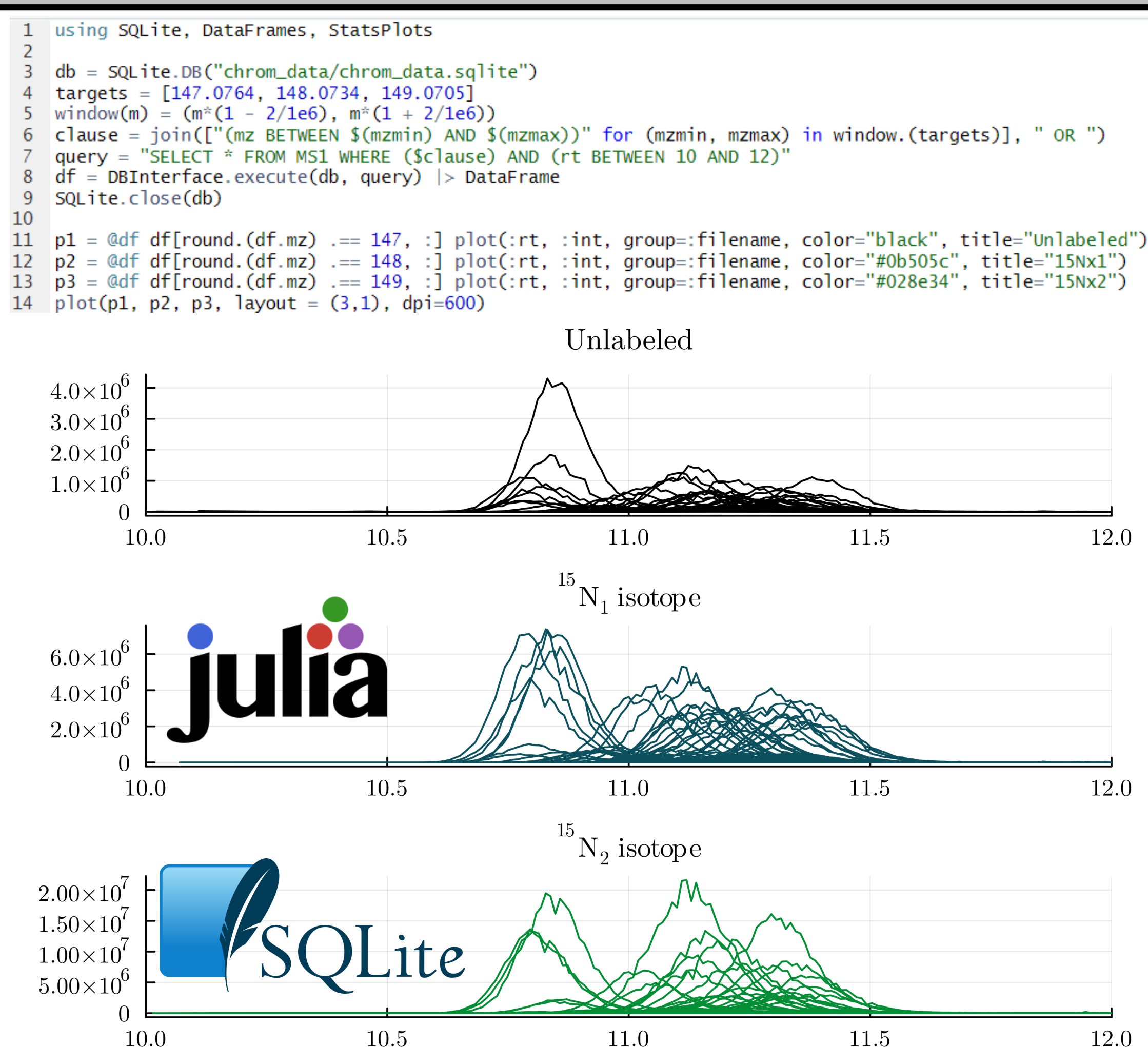
Indexing/ordering databases improved again by 10-100x

Improvements magnified with multiple files in one database

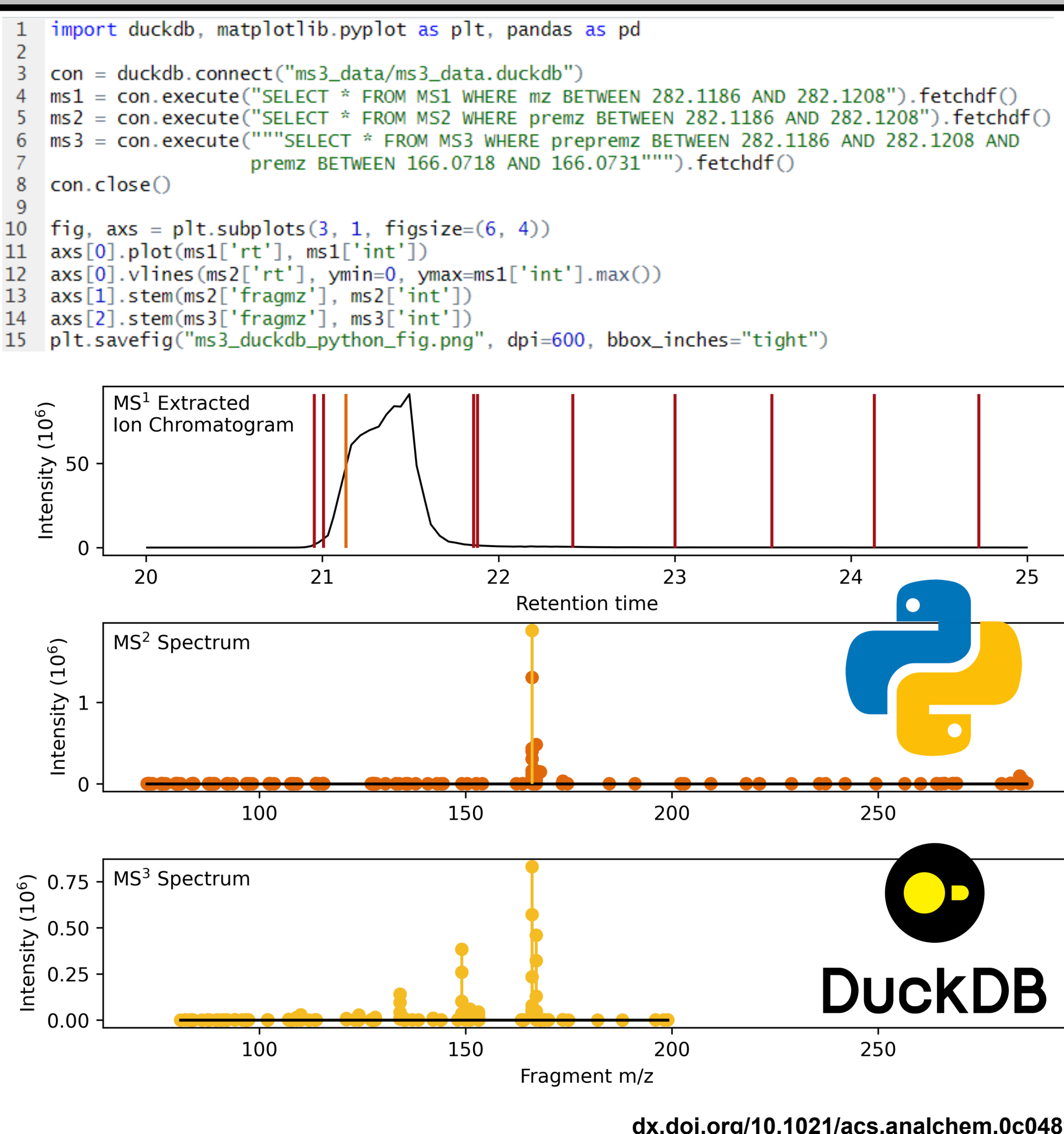


## Databases offer language agnostic data access with essentially zero memory overhead

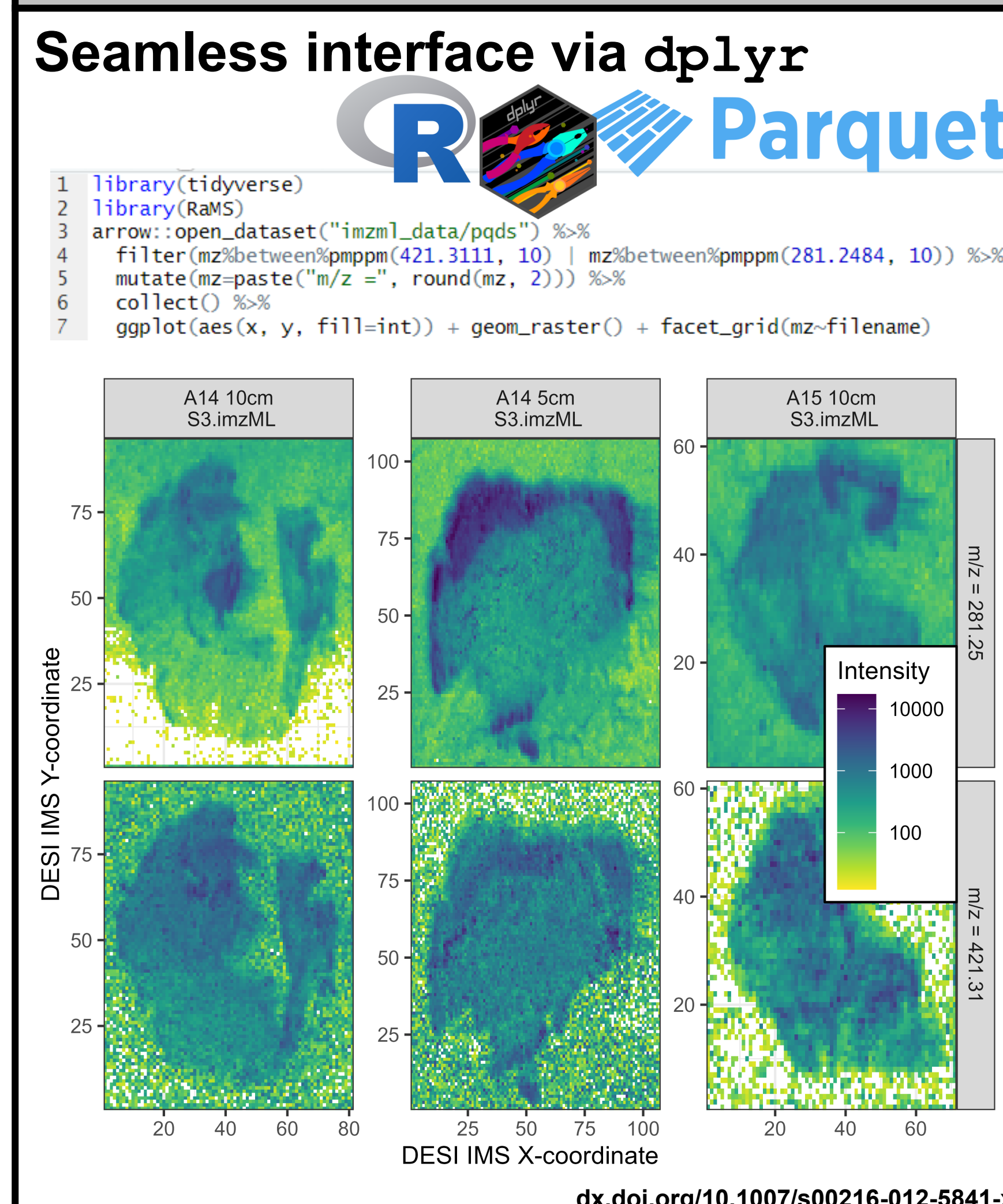
### EICs in Julia via SQLite



### MS<sup>3</sup> in Python via DuckDB



### Imaging in R via Parquet



# Simpler is often smaller *and* speedier