

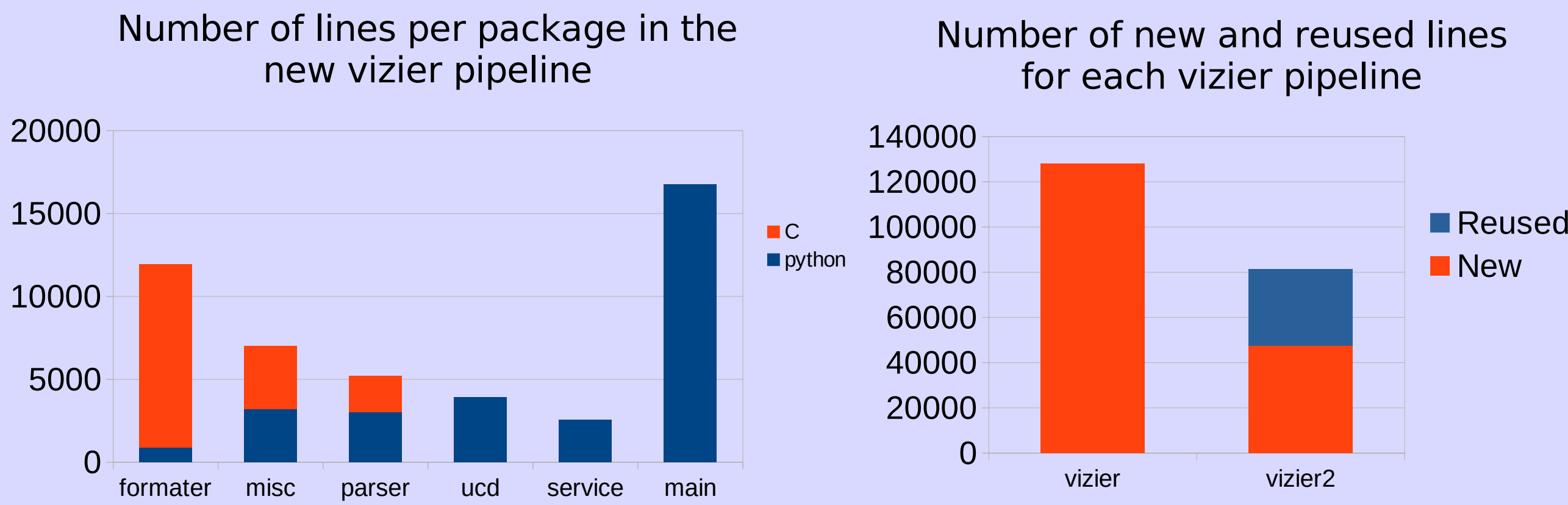
# Re-writing the VizieR catalogue ingestion pipeline



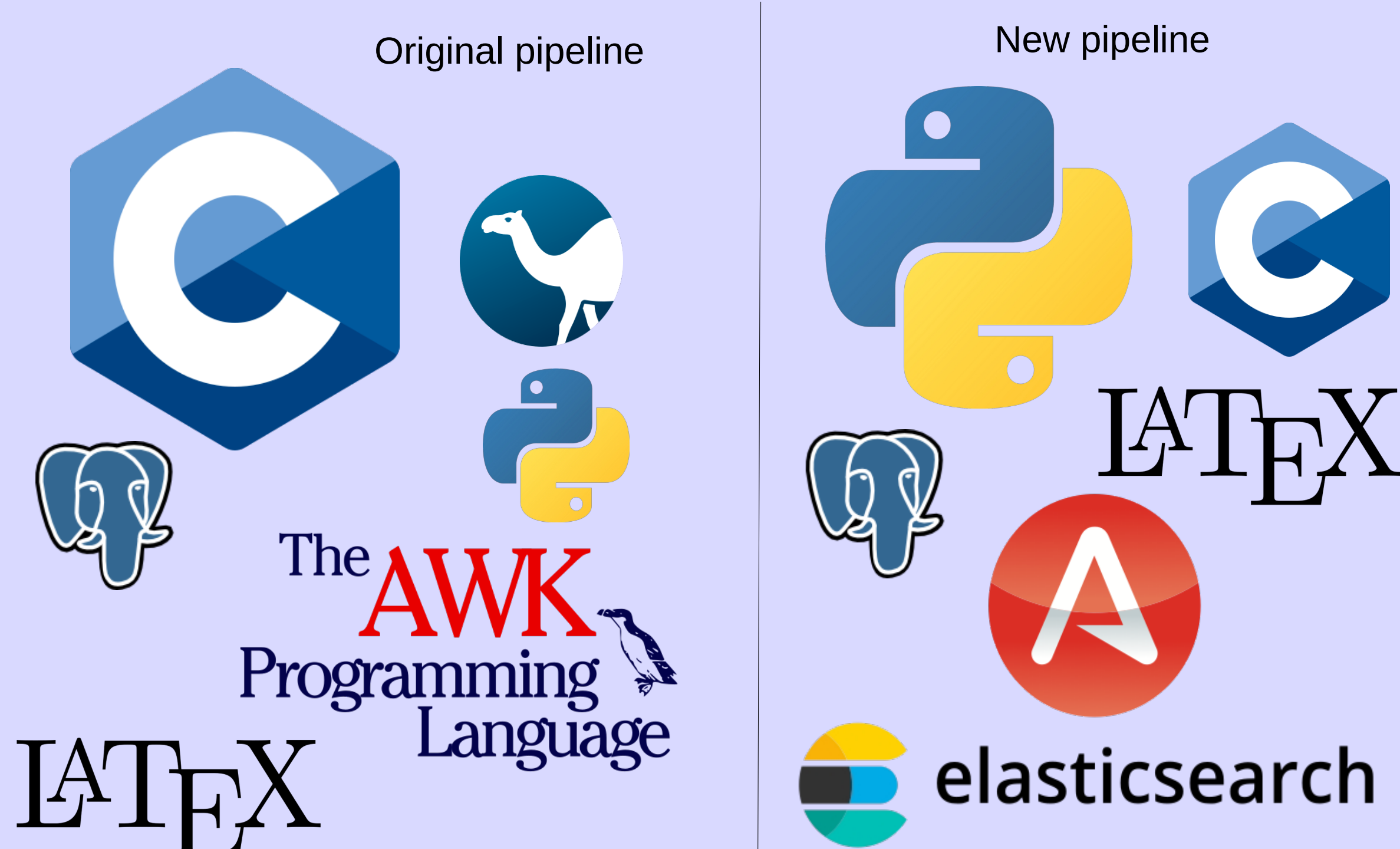
VizieR is a curated library of astronomical catalogues that provides enriched, verified data. Users have multiple ways of accessing the data, be it through the html pages, TAP queries, or interoperating with Topcat or astropy. On the internal side, VizieR works with a PostgreSQL database to store the data, as well as multiple files that describe and enrich catalogues, adding descriptions or links to other resources for instance. The internal workflow by which catalogues are ingested into VizieR is a semi-automated workflow : the description of catalogues and their enrichment with metadata are done by the documentalists at CDS, with the help of tools that streamline the process (for instance, the UCD builder, that gives proposition of UCDs for columns in the catalogue, that documentalists can verify and edit). The first part of the workflow builds the metadata around the catalogue, and relies heavily on the standardised description of the catalogue, while the second part, populating the database, needs no human intervention.



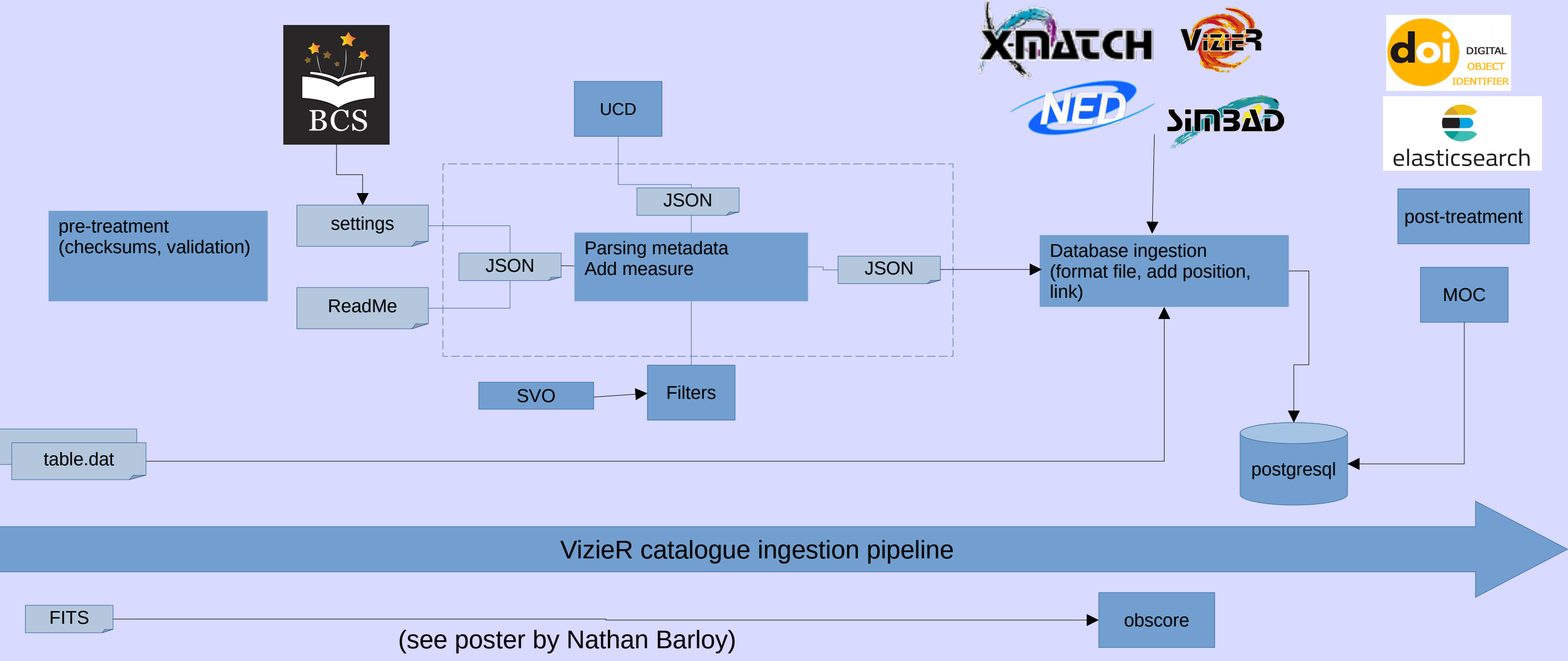
## Program length comparison



The original ingestion main program is written in C and is 17000 lines long. As such, maintaining the program and adding new features gets harder and harder. That is why we re-wrote the program : compartmentalising different parts of the workflow in different packages, which allowed us to reduce the size of the main program, making, in turn, the reading and understanding of the architecture easier for future developments. The new pipeline is in no way a simple transformation from C to python.



## Vizier catalogue ingestion pipeline



The new ingestion pipeline is written in C and Python, and reuses most of the ideas behind the original pipeline. The new pipeline is subdivided in different packages for different tools. During the building of the metadata, the standardised ReadMe file is read and combined with a settings file. They are transformed into a json file containing all metadata relevant to the catalogue. This file is then used to transform the data files : adding links to other data, adding columns such as position columns, merging files... Afterwards, data are ingested in the database.

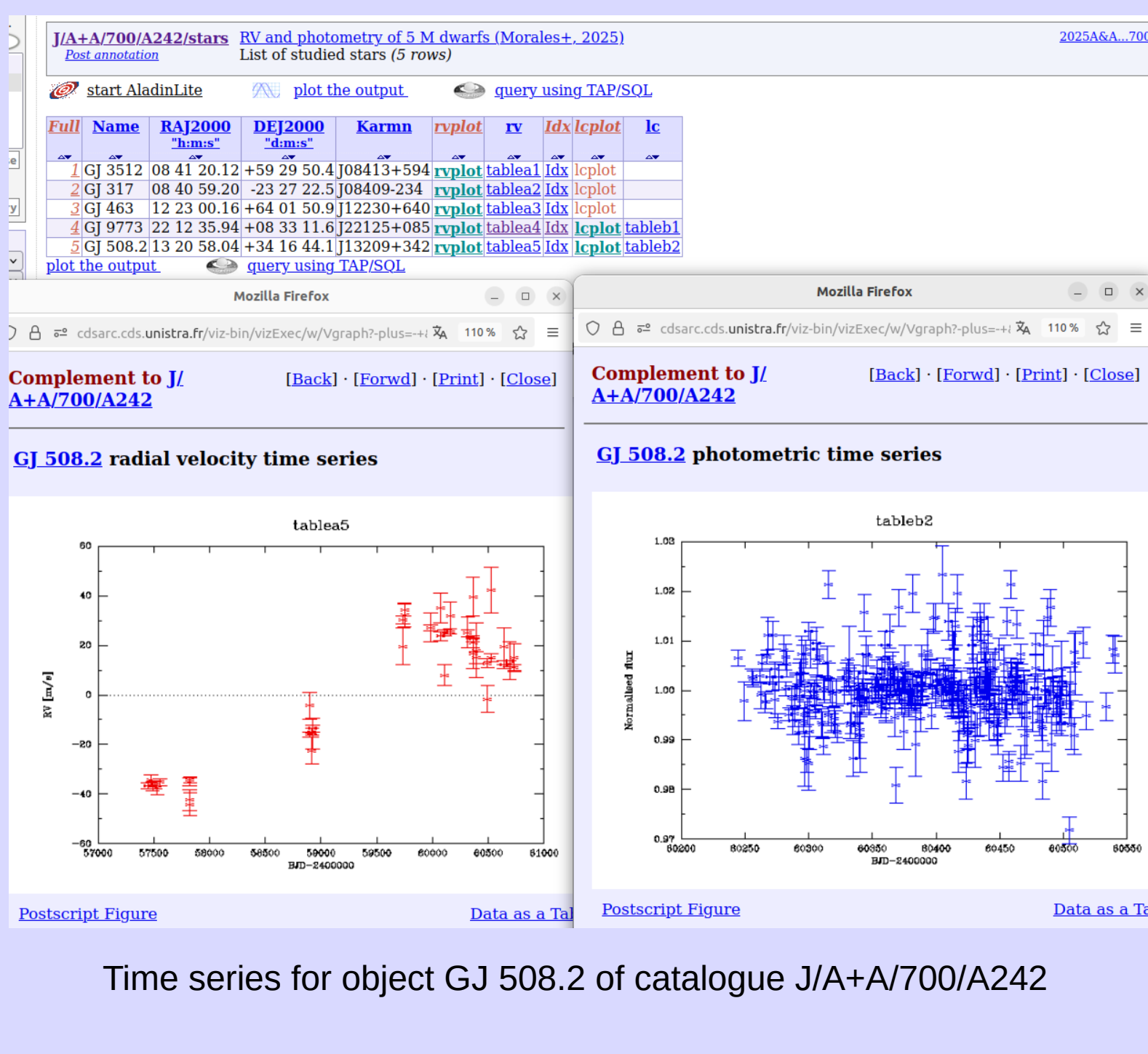
J110	Occultation Double-Star Observations (Evans 1983)
Catalogue of Occultation Double-Star Observations	
Evans D.S.	
-240 Colloquium No. 62, Current Techniques in Double and Multiple Star Research, Lowell Obs. Bull. 167, p. 73 (1983);	
-1983, ucd# 9.1, 736 1983, ucd# 9.1, 736 (Bull. 167, Volume 9)	
ADC Keywords: Stars, double and multiple ; occultations	
Description:	
This catalog contains data on 224 double stars observed photoelectrically during lunar occultations. The author cites the advantages of this method of double star detection as increased resolution, accuracy of the vector separations of roughly one-half arcsecond or better, and the opportunity to make photometric determinations of the magnitude differences between components. The vector separation is the true separation projected along a line perpendicular to the actual lunar limb. The catalog is a compilation of twelve years of observations from the literature (through roughly 1980).	
It is divided into three files. The first file, data1.dat, contains information on stars brighter than visual magnitude 6.7. The second, data2.dat, lists 540 catalog stars fainter than magnitude 6.7. The third file, data3.dat, contains data on faint stars with no SAO number. For these stars, data on their magnitudes or spectral types may be absent. In many cases there are multiple records per star, reflecting separate observations. The records are arranged by SAO number or other identifier, and contain visual magnitudes, spectral type, observing run number, a subjective grade of the probability of being double, the vector separation with computed error, position angle, and the lunar limb slope and its error. It also includes the magnitude difference between the components in (somewhat arbitrarily assigned) blue and red band passes. In the case of a triple star, the run number is repeated and the data for the triple given with magnitude differences from the brightest star.	
File Summary:	
Filename	Level
data1.dat	100
data2.dat	100
data3.dat	100

Byte-by-byte Description of file: data1.dat	
Bytes	Format
1-6	16
10-17	A8
19-23	F5.2
24-33	A10
35-42	A8
44-45	I2
46-47	I2
48-49	I2
51	I1
52	A1
54-60	F7.1
62-65	F4.1
66	A1
68-70	A1
71	A1
73-75	A1
77-78	I2
80-83	F4.2
85-88	F4.2
89	A1
90-94	F5.2
96-99	F4.2
100	A1

Example of a Readme file (see 10.1051/aas:2000169)

## New feature

Using the new basis for the handling of the metadata in the VizieR pipeline, we are able to integrate new standards in the service. A long-standing issue for VizeR has been the accessibility of associated data (images, spectra...) of catalogues. Indeed, apart from the visual html interface, it is impossible to access linked resources. To address this problem, we are implementing the Datalink standard (https://www.ivoa.net/documents/DataLink/20231215/REC-DataLink-1.1.html). We can add a column to the result of the user query, containing links to datalink instances, which themselves link back to the associated resources, with a short description for each.



```
<?xml version="1.4"?>
<!--RESOURCE type="results"-->
<TABLE>
  <FIELD name="ID" datatype="char" arraysize="1" ucd="meta.id.meta.main"/>
  <FIELD name="access.url" datatype="char" arraysize="1" ucd="meta.ref.url"/>
  <FIELD name="service.def" datatype="char" arraysize="1" ucd="meta.ref"/>
  <FIELD name="error.message" datatype="char" arraysize="1" ucd="meta.code.error"/>
  <FIELD name="semantics" datatype="char" arraysize="1" ucd="meta.code"/>
  <FIELD name="description" datatype="char" arraysize="1" ucd="meta.note"/>
  <FIELD name="content.type" datatype="char" arraysize="1" ucd="meta.code.nise"/>
  <FIELD name="content.qualifier" datatype="char" arraysize="1" ucd="meta.code.azim"/>
</TABLE>
<!--DATA-->
<TABLEDATA>
  <TB>
    <TR>
      <TD>http://cds.uistra.fr/viz-bin/v12Exec/u/9graph?+plus=+A/700/A242/table4
    </TD>
    <TD>
      <!--auxiliary-->
      <TB>
        <TR>
          <TD>Show the Radial Velocity curves (timeserie) from J/A+A/700/A242
        </TD>
        <TD>
          <!--text/html-->
          <!--timeserie-->
        </TD>
      </TB>
    </TD>
  </TB>
  <TB>
    <TR>
      <TD>http://cds.uistra.fr/viz-bin/v12Exec/u/9graph?+plus=+A/700/A242/table1
    </TD>
    <TD>
      <!--auxiliary-->
      <TB>
        <TR>
          <TD>Show the light curves (timeserie) from J/A+A/700/A242
        </TD>
        <TD>
          <!--text/html-->
          <!--timeserie-->
        </TD>
      </TB>
    </TD>
  </TB>
</TABLEDATA>
</TABLE>
</RESOURCE>
</VOTABLE>
```

## Retrocompatibility

A large part of the work behind writing the new VizieR pipeline has been making sure that it is backwards compatible with the ~30 thousand catalogues saved in our database. Although it mainly means abiding by pre-existing rules, it also means scouring the database to find exceptions, outliers, and even unspoken rules.

## Links



The VizieR website



The VizieR database of astronomical catalogues article



Latest Datalink documentation



Ivan Brossard  
i.brossard@astro.unistra.fr  
With the help of Gilles Landais  
gilles.landais@astro.unistra.fr

