



university of  
 groningen

# Astro-WISE: Tracing and Using Lineage for Scientific Data Processing

Johnson Mwebaze, Danny Boxhoorn and Edwin Valentijn

*Kapteyn Astronomical Institute*

*University of Groningen*

*The Netherlands*

- Introduction
- Motivation
- Astro-WISE
- Framework
- Data Lineage and Data Processing
- Tracing Lineage in Astro-WISE
- Experiments [Case Studies]
- Related Work
- Conclusions and prospects

- Data growing exponentially
  - caused by successive generations of inexpensive sensors + exponentially faster computing
  - Old and New Projects in the pipeline; Creating accelerated acquisition of large volumes of observational data
- Changes the nature of scientific computing (caused by Technology, Sociology, Economics)
  - The Virtual Observatory, a collection of data centers each with unique collections of astronomical data.
  - how to believe, share and validate data as published.
  - how to selectively retrieve and process data from multiple archives interactively
- how to allow flexibility in e-science systems to enable researchers to work with and change methods on the fly while capturing all processing activities/events
- Scientific computing is revolving around data [Jim Gray] - Move analysis to the data.

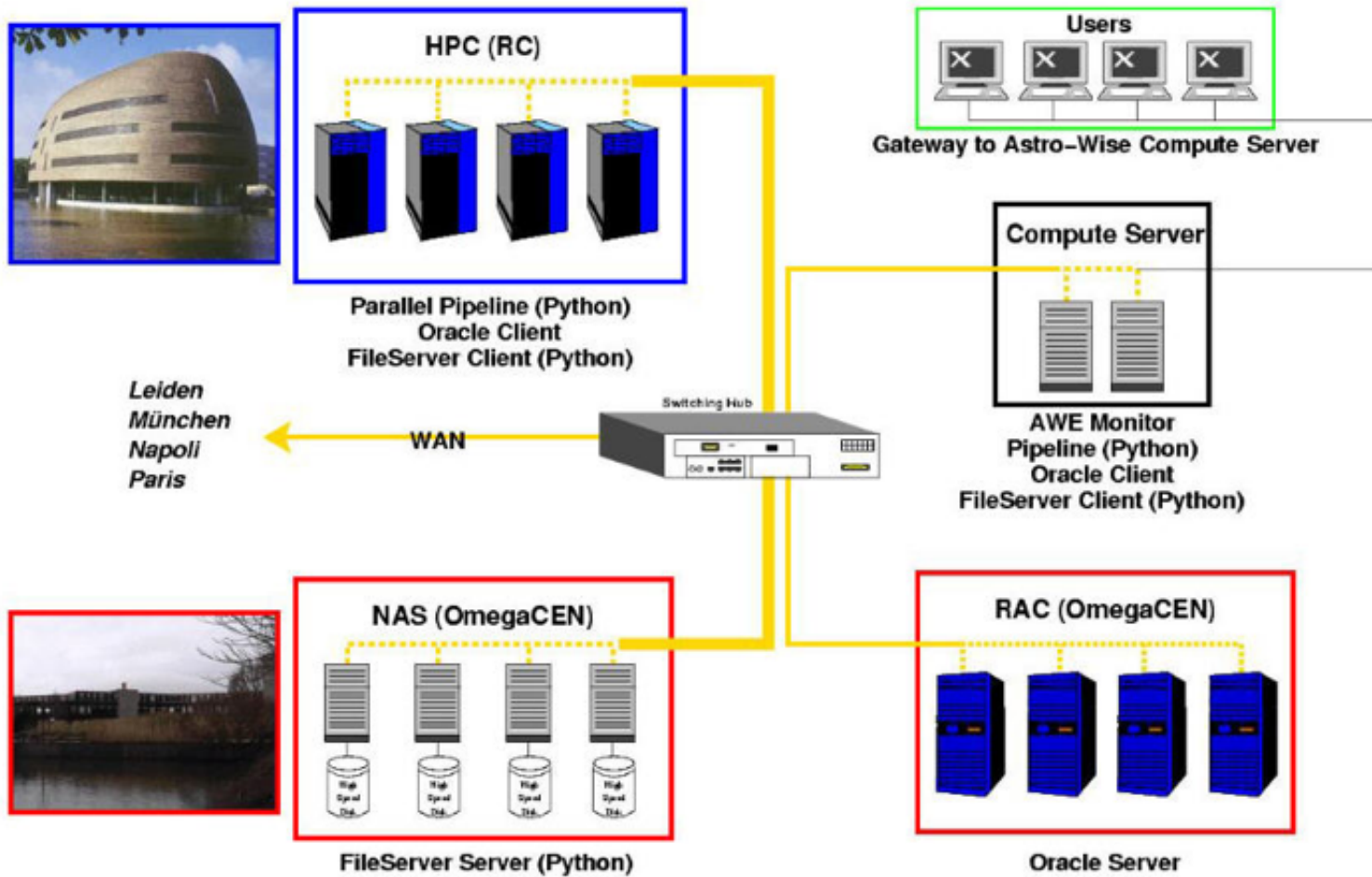
- Results validation
- Result debugging
- Reproducibility
- Repeatability
- Explanation (derivations, traces, proof trees)
- Runtime monitoring (Profiling, benchmarking)
- Performance Optimization (smart rerun)
- Fault tolerance, crash recovery
- System re-design

- Goal: Connect raw data to end user in VO environment
- an e-science infrastructure with fully distributed resources, which allows teams distributed over Europe
- handles optical, IR and radio.
  - Current: WFI@2.2m, INT, Subaru, ACS@ST
  - Near future: Panstarrs, VST, VISTA, Lofar
  - Far future: EUCLID
- Services
  - pipeline data reduction, calibration and re-calibration
  - Image comparisons and combinations
  - Working with source lists
  - Visualization of astronomical data
  - publishing results to the Virtual Observatory Archives

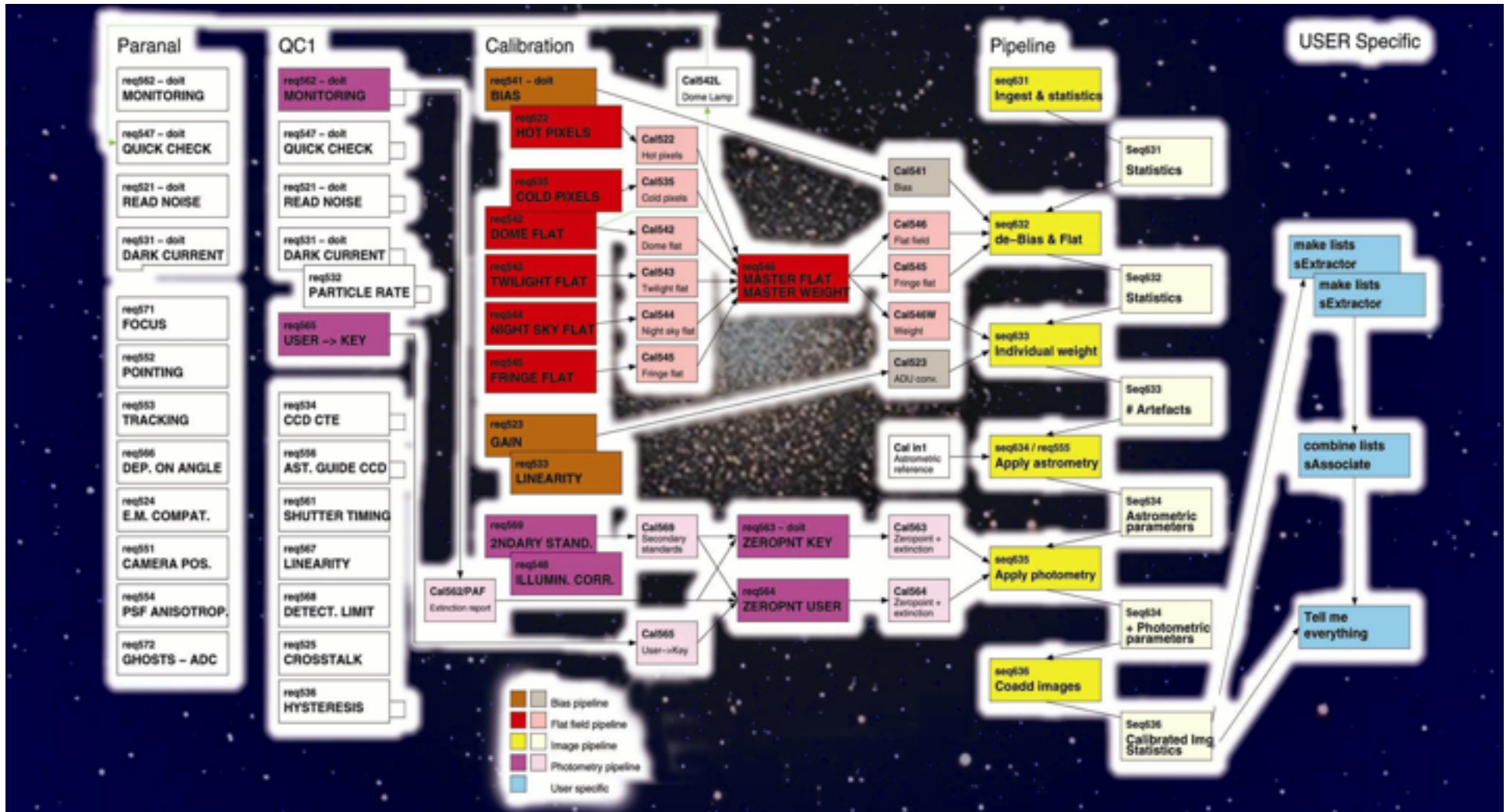
- Wide-field imaging = vast amounts of data
  - VST sees equivalent of Southern sky with 0.2 pixels in 3 years.
  - A project like KIDS (1000 Sq Deg) has  $\approx 10^6$  8Mpix raw data images
- Handling of the data is non-trivial
  - Pipeline data reduction
  - Calibration with very limited resources
  - Things change in time:
    - Physical changes (atmosphere, various gains)
    - Code (new methods, bugs)
    - Human insight in changes
  - Working with source lists



## OmegaCEN & HPC



© 2003 Astro-Wise



- The core of the system exploits three properties in a database environment.
  - Firstly, we apply the principle of inheritance using Object Oriented Programming (Python), where all Astro-WISE objects inherit key properties for database access, such as persistency of attributes.
  - Secondly, the linking (associations or references) between instances of objects in the database is completely maintained.
  - Thirdly, continuous growth of the database through the addition of new information or improvements made to existing information.
    - archiving or regenerate dynamically
    - exchanging methods, scripts and configuration

- Persistence -
  - persistence objects is the core of Astro-WISE pipeline processing and data lineage
  - Processing is done through invocation of methods on these objects.
  - All the I/O of the pipeline processes makes use of a database
  - Data can therefore only be manipulated through interaction with the database
- Extendable Schema
  - allows extending and modifying object attributes and method definitions through inheritance and polymorphism
- Version Control (Object Versioning) and Preservation Management
  - the connection between these classes and the created objects
  - addresses object inconsistencies (mismatches) created as a result of schema modification

- For each persistent class, a number of methods are defined which interact with the federated database. All class instantiations are automatically made persistent in the database forming an archive of all targets. To achieve this, the following major classes are implemented.
  - DBObject is the root class of the hierarchy of the persistent classes. This class defines the primary key `object_id` of all objects.
  - DBObjectMeta is the metaclass of DBObject. This is the class that is responsible for class creation and object instantiation. (class factory)
  - DBProperties is the module that defines all persistent attribute types (data dictionary) which are defined by persistent.
  - DBSelect implements a query language that is a natural extension to Python and that incorporates data lineage in the query syntax.
  - 
  - *sourcelist.coadd\_frame.regr\_frame.reduced.raw.filename*
  - DBProxy is an abstract interface to database vendor specific operations.

- lineage is captured as dependencies between persistent objects that refer to another persistent objects.
- The simplest persistent class has only one persistent attribute, like

```
class ClassA(DBObject):  
    attribute = persistent('Description',  
                           attribute_type, default_value)
```
- attribute will have the type attribute\_type.
- Atomic attribute types, such as integer, string or float, are translated into their database equivalents when an object is made persistent.
- If attribute\_type is a subclass of a DBObject, the attribute will be a link to an object of that subclass.
- If default\_value is an empty list, then the attribute will be an array of objects of attribute\_type.

- another Example

```
class ClassB(Object):  
    e = persistent('A link', ClassA, None)  
    f = persistent('List of links', ClassA, [])  
    g = persistent('Link to object of ClassB')
```

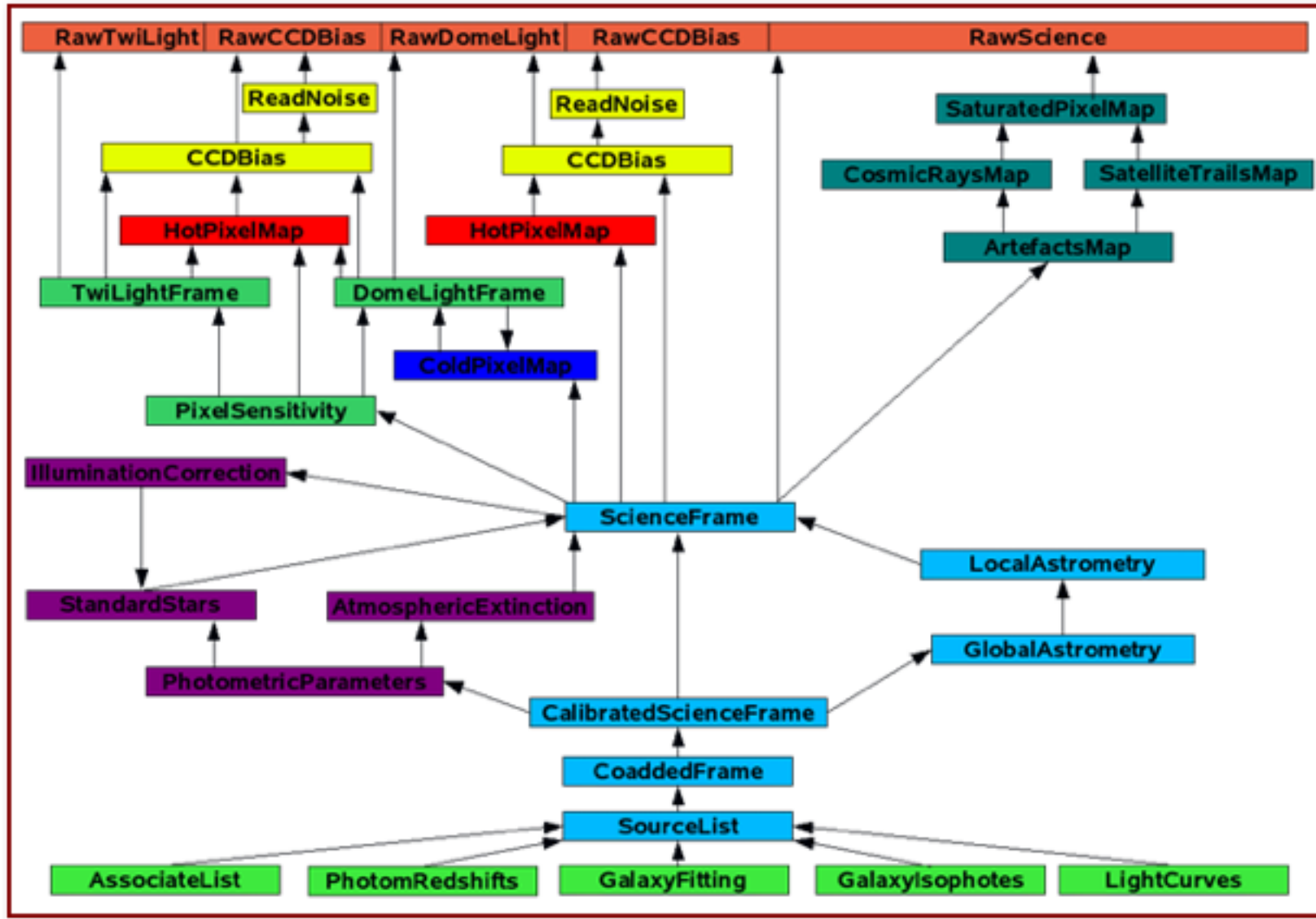
- ClassB defines three persistent attributes:
- e is a link to an instance of ClassA,
- f is a array of links to instances of ClassA (default empty), and
- g is a link to another instance of ClassB.

- To support the storage of files in a similar way, the DataObject class is used.

```
class DataObject(DBObject):  
    filename = persistent('File part', str, '')
```

- The filename attribute is used to implement store() and retrieve() methods to transfer files to and from dataservers.
- Because the filename is kept in the database the storage and retrieval of files is transparent to the application.

- backward-chaining approach while processing data
- The make metaphor -specify the result instead of the input
  - This allows the end user to trace the data product, following all its dependencies up to the raw observational data and, if necessary, to re-derive the result with better data, and/or improved methods
- Uses target Processing (TP)
  - The *Target* processor employs the dependency logic that is constructed using lineage data
  - Targets are rebuilt in a recursive cascade only if the dependencies have changed (i.e workflow reduction).





**Astro-Wise Target Processing**

**Contact**  
 Willem-Jan Vriend

**DB User**  
 awevalentyn

**Help**  
 Getting Started

**Project**  
 WFI@2.2m

**Instrument**  
 WFI

**Single host (status)**  
 dpu.hpc.rug.astro...

**Parallel host (status)**  
 dpu.hpc.rug.astro...

**Processing**  
 Image pipeline  
 Depth 99  
 Full processing

**Options**  
[Job overview](#)  
 Upload Code  
 Popup Info  
 Object options  
[Renew cookie](#)

**Target**  
 MasterBias  
 MasterFlat  
**RegriddedFrame**  
 CoAddedFrame  
 SourceList

Advanced

**Querying**  
 Serial processing  
 Image pipeline  
 Depth 2  
 Full  
 Configure

**Filter**  
 #844 CousinsR  
 Date Obs: 2002-06-08  
 Time: 23:25:45

**Chip** >>  
 ccd50

**Parameters**  
 Configure

**Object** select    **RA**    **DEC**    +/- 0.5 select

**Search**

**Possible targets**

0	#844 CousinsR	08 Jun 2002 23:25:45	<a href="#">view all chips</a>	<a href="#">Process</a>
(+) (-)	0.0	RegriddedFrame (outdated)		
	0.1	AstrometricParameters (outdated)		
	0.1	ReducedScienceFrame (outdated)		
	1.1	BiasFrame		
	1.2	ColdPixelMap (new version available)		
	1.3	MasterFlatFrame (outdated)		
	1.4	FringeFrame (null)		
	1.5	HotPixelMap (outdated)		
	2.1	BiasFrame (new version available)		
	1.6	IlluminationCorrectionFrame		
	0.7	RawScienceFrame		
	1.2	GainLinearity		
	1.3	PhotometricParameters		
	0.4	ReducedScienceFrame (outdated)		
	1.1	BiasFrame		
	1.2	ColdPixelMap (new version available)		
	1.3	MasterFlatFrame (outdated)		
	1.4	FringeFrame (null)		
	1.5	HotPixelMap (outdated)		
	1.6	IlluminationCorrectionFrame		

```
+--Name of SourceList : SL-0000010011
+--COMBINE_METHOD: -1
+--OBJECT: field07B
+--SLID: 10011
+--associatelist: -1
+--astrom_params: None
+--chip: None
+--creation_date: 2006-03-20 13:57:25
+--detection_frame: None
+--filename: ldac.cat
+--filter: <astro.main.Filter.Filter object at 0x8df22d0>
+--filters:
+--frame: <astro.main.RegriddedFrame.CoaddedRegriddedFrame object at 0x8dd1b10>
+--globalname:
+--instrument: <astro.main.Instrument.Instrument object at 0x8df2050>
+--name: SL-0000010011
+--number_of_sources: 2236
+--object_id: '0F6E9292574FC6EEE0407D81C5061B26'
+--process_params: <astro.main.SourceList.SourceListParameters object at 0x8df20d0>
+--sexconf: <astro.main.Config.SExtractorConfig object at 0x8df2150>
+--sexparam: <class 'common.database.typed_list.typed_list'>(<type 'str'>, ['DELTA_J2000', 'ALPHA_J2000'])
+--sources: {}
+ - [some output omitted]
```

● Queries e.g `sourcelist.frame.info()`



```

Main-Object : SourceList      qualityview
|+ OBJECT                      | STD
|- name                        | GAS-Sci-WVRIEND-WFI-----#842-ccd50-Red---Sci-54964.6029650-8d91ef050a6e71840937ba4a40f538
|+ number_of_sources          | 36
|+ sexparam                    | (, ['FLUX_RADIUS', 'FWHM_IMAGE', 'MAG_ISO'])
| +astrom_params              | [level : 2] AstrometricParameters --> QC-WVRIEND-WFI-----#842-ccd50-AstromParams---Residu
| | +associateconf            | [level : 3] AssociateConfig
| | +astromconf                | [level : 3] AstromConfig
| | +process_params           | [level : 3] AstrometricParametersParameters
| | +reduced                   | [level : 3] ReducedScienceFrame      qualityview
| | |- filename                | (image)Sci-WVRIEND-WFI-----#842-ccd50-Red---Sci-54964.6029650.....fits
| | | +astrom                  | [level : 4] Astrom
| | | +bias                     | [level : 4] BiasFrame
| | | +chip                     | [level : 4] Chip
| | | +cold                     | [level : 4] ColdPixelMap
| | | +filter                   | [level : 4] Filter
| | | +flat                     | [level : 4] MasterFlatFrame
| | | +hot                      | [level : 4] HotPixelMap
| | | +illumination             | [level : 4] IlluminationCorrectionFrame
| | | +process_params           | [level : 4] ReducedScienceFrameParameters
| | | +raw                      | [level : 4] RawScienceFrame      qualityview
| | | +template                 | [level : 4] Template
| | | +weight                   | [level : 4] WeightFrame
| | +sexconf                   | [level : 3] SextractorConfig
| | |- WEIGHT_IMAGE            | (image) Cal-WVRIEND-WFI-----#842-ccd50-Red---Wei-54964.6031358-....fits
| |- filename                  | (image) Sci-WVRIEND-WFI-----#842-ccd50-Red---Sci-54964.6029650-....fits

```

- we can trace to a pixel level
- this can be used to find extremely rare astronomical events, such as a fast moving star, near-earth solar system objects, or variable objects like ultra-compact binary systems or distant supernovae

	RawFrame	ReducedFrame	RegriddedFrame	CoaddedRegriddedFrame	BiasFrame	ColdPixelMap	MasterFlatFrame	FringeFrame	HotPixelMap	IlluminationCorrect
SLID=4147 SID=0 RA=11.3289 DEC=-29.3084 X=1765 Y=84										
SLID=136151 SID=27 RA=9.5151 DEC=-28.9031 X=883 Y=45								None		
SLID=136151 SID=29 RA=9.6949 DEC=-28.9023 X=538 Y=126								None		
SLID=136151 SID=28 RA=9.1784 DEC=-28.9041 X=247 Y=96								None		
SLID=4147 SID=40 RA=11.4650 DEC=-29.3785 X=284 Y=187										

