



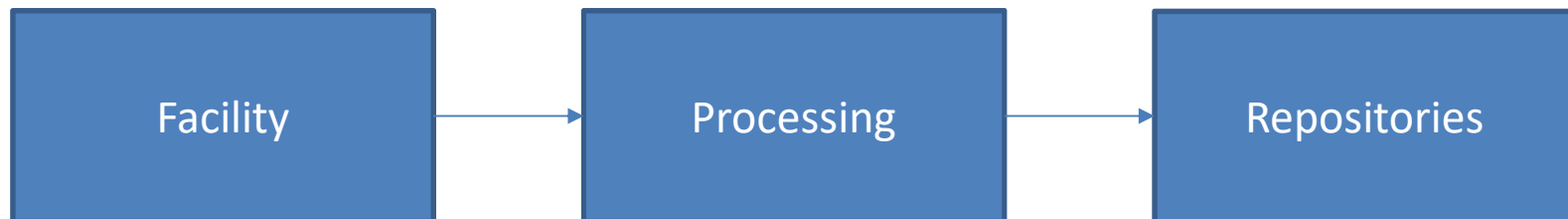
CryoEM DataModels for processing

Biocomputing Unit, CNB-CSIC
Instruct Image Processing Center

Data Models Use

- Public repositories:
 - EMDB: https://www.ebi.ac.uk/emdb/documentation#data_model
 - EMPIAR: <https://www.ebi.ac.uk/empiar/faq>
- Facilities and LIMS:
 - ISpyB: <https://github.com/ispyb/ispyb-database-modeling>
 - EMAdmin: <https://github.com/I2PC/EMadmin/tree/master/EMadmin>
- Data processing workflows:
 - Scipion: <https://github.com/scipion-em/scipion-em/blob/devel/pwem/objects/data.py>
 - (CryoSparc, Relion, ...)

Data models use

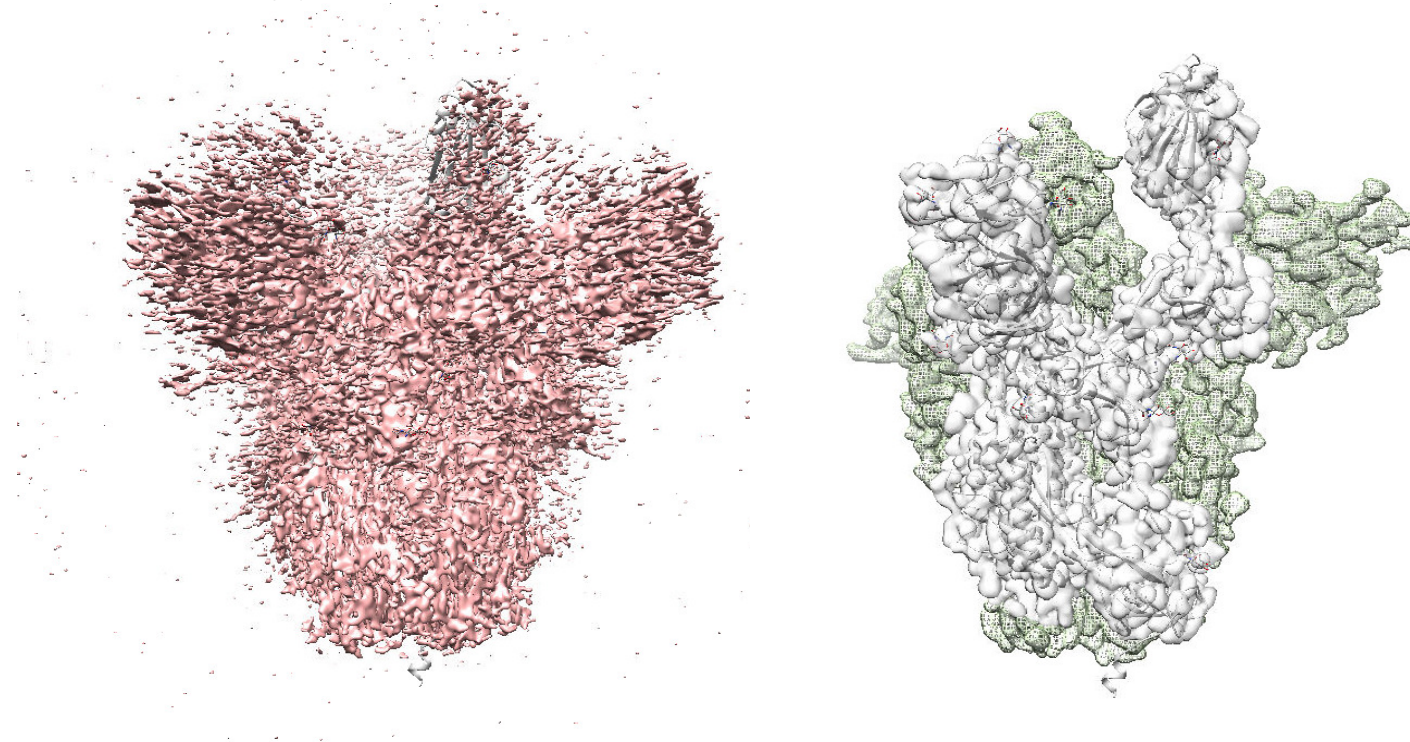


Data Models Workflows

We should distinguish the purpose:

- To run on **a new** dataset (WorkflowHub)
- To **understand** a given result (EMPIAR)

EMDB 22301



Workflows

View Protocols SPA

- Protocols SPA
 - Imports
 - import movies
 - import micrographs
 - import particles
 - import volumes
 - Movies
 - Micrographs
 - CTF estimation
 - Preprocess
 - Particles
 - Picking
 - Extract
 - Preprocess
 - Filter
 - Mask
 - 2D
 - Alien
 - Classifu
 - 3D
 - Initial volume
 - Preprocess
 - Classifu
 - Refine
 - Postprocess
 - Analysis
 - Reconstruct
 - Tools
 - Exports

Protocol Run: XmippProtMovieMaxShift

Protocol: xmipp3 - movie maxshift finished Cite Help

Run

Run name: xmipp3 - movie maxshift Comment: _____

Run mode: Continue Restart Host: localhost

Use queue? Yes No

Wait for: _____

Input

Input Movies: xmipp3 - FlexAlign.outputMovies

Rejection type: by frame or movie

Max. frame shift (A):

Max. movie shift (A):

Close Save Execute

View: Tree

```

PROJECT
├── pwem - import movies finished
├── xmipp3 - movie gain finished
├── xmipp3 - FlexAlign finished
└── xmipp3 - movie maxshift finished
    ├── pwem - ctffind4 finished
    └── gctf - ctf estimation finished
            
```

Summary Methods Output Log

Input

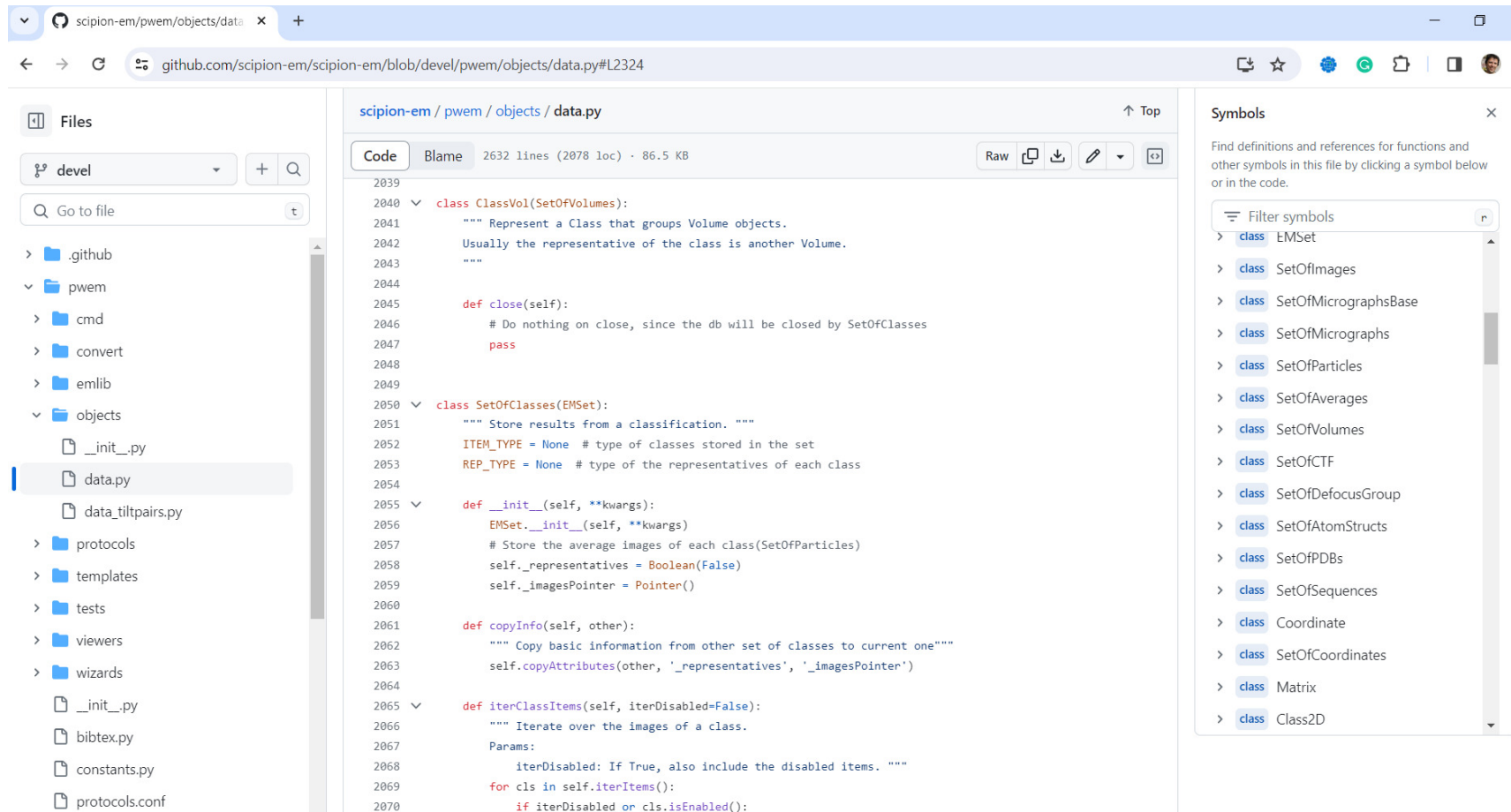
inputMovies (from xmipp3 - FlexAlign -> outputMovies [outputMovies]) SetOfMovies (30 items, 3710 x 3838 x 50 [1-50], 0)

Output

▶ xmipp3 - movie maxshift -> outputMovies SetOfMovies (29 items, 3710 x 3838 x 50 [1-50], 0)
▶ xmipp3 - movie maxshift -> outputMicrographs SetOfMicrographs (29 items, 3710 x 3838, 0.49)



Data Models



scipion-em / pwem / objects / data.py

Code Blame 2632 lines (2078 loc) · 86.5 KB

```
2039
2040 class ClassVol(SetOfVolumes):
2041     """ Represent a Class that groups Volume objects.
2042         Usually the representative of the class is another Volume.
2043     """
2044
2045     def close(self):
2046         # Do nothing on close, since the db will be closed by SetOfClasses
2047         pass
2048
2049
2050 class SetOfClasses(EMSet):
2051     """ Store results from a classification. """
2052     ITEM_TYPE = None # type of classes stored in the set
2053     REP_TYPE = None # type of the representatives of each class
2054
2055     def __init__(self, **kwargs):
2056         EMSet.__init__(self, **kwargs)
2057         # Store the average images of each class(SetOfParticles)
2058         self._representatives = Boolean(False)
2059         self._imagesPointer = Pointer()
2060
2061     def copyInfo(self, other):
2062         """ Copy basic information from other set of classes to current one"""
2063         self.copyAttributes(other, '_representatives', '_imagesPointer')
2064
2065     def iterClassItems(self, iterDisabled=False):
2066         """ Iterate over the images of a class.
2067         Params:
2068             iterDisabled: If True, also include the disabled items. """
2069         for cls in self.iterItems():
2070             if iterDisabled or cls.isEnabled():
```

Files

devel

Go to file

- .github
- pwem
 - cmd
 - convert
 - emlib
 - objects
 - __init__.py
 - data.py
 - data_tiltpairs.py
 - protocols
 - templates
 - tests
 - viewers
 - wizards
 - __init__.py
 - bibtex.py
 - constants.py
 - protocols.conf

Symbols

Filter symbols

- class EMSet
- class SetOfImages
- class SetOfMicrographsBase
- class SetOfMicrographs
- class SetOfParticles
- class SetOfAverages
- class SetOfVolumes
- class SetOfCTF
- class SetOfDefocusGroup
- class SetOfAtomStructs
- class SetOfPDBs
- class SetOfSequences
- class Coordinate
- class SetOfCoordinates
- class Matrix
- class Class2D

Data Models

scipion-em-xmipp / xmipp3 / protocols / protocol_movie_max_shift.py

↑ Top

Code

Blame

391 lines (337 loc) · 17.3 KB

Raw



```
42     class XmippProtMovieMaxShift(ProtProcessMovies):
66         def __init__(self, **args):
73
74         # ----- DEFINE param functions -----
75     def _defineParams(self, form):
76         form.addSection(label=Message.LABEL_INPUT)
77         form.addParam('inputMovies', PointerParam, important=True,
78                       label=Message.LABEL_INPUT_MOVS,
79                       pointerClass='SetOfMovies',
80                       help='Select a set of previously aligned Movies.')
81
82         form.addParam('rejType', params.EnumParam, choices=self.REJ_TYPES,
83                       label='Rejection type', default=self.REJ_OR,
84                       help='Rejection criteria:\n'
85                            ' - *by frame*: Rejects movies with drifts between '\n'
86                            'frames bigger than a certain maximum.\n'
87                            ' - *by whole movie*: Rejects movies with a total '\n'
88                            'travel bigger than a certain maximum.\n'
89                            ' - *by frame and movie*: Rejects movies if both '\n'
90                            'conditions above are met.\n'
91                            ' - *by frame or movie*: Rejects movies if one of '\n'
92                            'the conditions above are met.')
93
94         form.addParam('maxFrameShift', params.FloatParam, default=5,
95                       label='Max. frame shift (A)',
96                       condition='rejType==%s or rejType==%s or rejType==%s'
97                               % (self.REJ_FRAME, self.REJ_AND, self.REJ_OR),
98                       help='Maximum drift between consecutive frames '\n'
99                            'to evaluate the frame condition.')
100        form.addParam('maxMovieShift', params.FloatParam, default=15,
101                      label='Max. movie shift (A)',
102                      condition='rejType==%s or rejType==%s or rejType==%s'
```


Data Models

scipion-em-xmipp / xmipp3 / protocols / protocol_movie_max_shift.py

↑ Top

Code

Blame

391 lines (337 loc) · 17.3 KB

Raw



```
42     class XmippProtMovieMaxShift(ProtProcessMovies):
205         def fillOutput(newDoneList, firstTime, AccOrDisc='Accepted'):
247             # Update output movies
248
249             if movieSet.getSize() > 0:
250                 self._updateOutputSet('outputMovies%s' % suffix, movieSet,
251                                     streamMode)
252
253             if self.inputMics is not None and micsSet.getSize() > 0:
254                 self._updateOutputSet('outputMicrographs%s' % suffix, micsSet,
255                                     streamMode)
256
257             if self.inputDwMics is not None and micsDwSet.getSize() > 0:
258                 self._updateOutputSet('outputMicrographsDoseWeighted%s' % suffix,
259                                     micsDwSet, streamMode)
259             if firstTime: # define relation just the first time
260                 if movieSet.getSize() > 0:
261                     self._defineTransformRelation(self.inputMovies.get(), movieSet)
262                 if self.inputMics is not None and micsSet.getSize() > 0:
263                     self._defineTransformRelation(self.inputMics, micsSet)
264                 if self.inputDwMics is not None and micsDwSet.getSize() > 0:
265                     self._defineTransformRelation(self.inputDwMics, micsDwSet)
266
267             movieSet.close()
268             if self.inputMics is not None:
269                 micsSet.close()
270             if self.inputDwMics is not None:
271                 micsDwSet.close()
272
```


Scipion data model

DB Browser for SQLite - /home/coss/data/Dropbox/H/ScipionUserData/projects/CoursePasteur/Runs/007895_ProtMotionCorr/microgr...

File Edit View Tools Help

New Database Open Database Write Changes Revert Changes Open Project Save Project Attach Database

Database Structure Browse Data Edit Pragmas Execute SQL

Table: Classes Filter...

id	label_property	column_name	class_name
1	self	c00	Micrograph
2	_index	c01	Integer
3	_filename	c02	String
4	_samplingRate	c03	Float
5	_acquisition	c04	Acquisition
6	_acquisition._magnification	c05	Float
7	_acquisition._voltage	c06	Float
8	_acquisition._sphericalAberration	c07	Float
9	_acquisition._amplitudeContrast	c08	Float
10	_acquisition._doseInital	c09	Float
11	_acquisition._dosePerFrame	c10	Float
12	_acquisition.opticsGroupInfo	c11	String
13	_micName	c12	String
14	plotGlobal	c13	Image
15	plotGlobal._index	c14	Integer
16	plotGlobal._filename	c15	String
17	plotGlobal._samplingRate	c16	Float
18	_rInAccumMotionTotal	c17	Float
19	_rInAccumMotionEarly	c18	Float
20	_rInAccumMotionLate	c19	Float

Mode: Text

1

Type of data currently in cell: Text / Numeric
1 character(s)

Apply

Remote

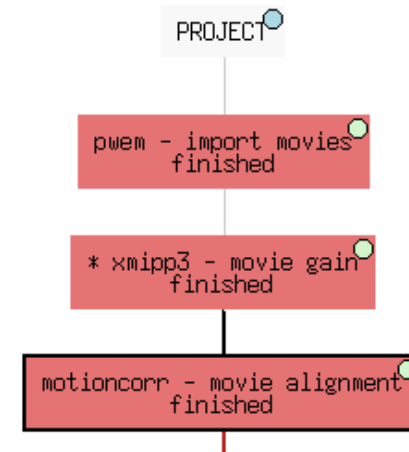
Identity Select an identity to connect

DBHub.io Local Current Database

Name	Last modified	Size	Co
------	---------------	------	----

SQL Log Plot DB Schema Remote

UTF-8



Summary Methods Output Log

Input

inputMovies (from xmipp3 - movie gain -> outputMovies [outputMovies])

Output

motioncorr - movie alignment -> outputMovies

motioncorr - movie alignment -> outputMicrographsDoseWeighted

Movies (30 items, 3710 x 3838 x 50 [1-50], 0.49 Å/px)

Movies (30 items, 3710 x 3838 x 50 [1-50], 0.49 Å/px)

Micrographs (30 items, 3710 x 3838, 0.49 Å/px)

Extensible data model

- SetOfMicrographs:
 - Compulsory fields
 - Extended fields

CryoEM ontology model

Search CRYOEM...

Exact match Include obsolete terms Include imported terms

Tree

Graph

- EMThing (110)
 - EMObject (44)
 - Acquisition
 - AtomStruct
 - Coordinate
 - CTFModel
 - DefocusGroup
 - EMSet (23)
 - SetOfAtomStructs
 - SetOfClasses (3)
 - SetOfCoordinates
 - SetOfCTF
 - SetOfDefocusGroup
 - SetOfFSCs
 - SetOfImages (11)
 - SetOfImages2D (7)
 - SetOfAverages
 - SetOfMicrographs
 - SetOfMovies
 - SetOfParticles (3)
 - SetOfImages3D (2)
 - SetOfNormalModes
 - SetOfSequences
 - FSC
 - Image (10)
 - NormalMode
 - Sequence
 - Transform
 - EMProtocol (64)

Show counts
 Show obsolete terms
 Show all siblings

▼ Class Information

IsASetOf
Micrograph

▼ Class Relations

Subclass of
SetOfImages2D

<https://www.ebi.ac.uk/ols4/ontologies/cryoem>

Ontology

Actions achieved within the project: Workflow FAIRness

- CryoEM ontology
- Ontology Lookup Service: <https://www.ebi.ac.uk/ols/ontologies/cryoem>
- BioPortal: <https://bioportal.bioontology.org/ontologies/CRYOEM>
- FAIRsharing: <https://fairsharing.org/bsg-s001477/>
- RO-Crate describing the image processing process
- JSON and CWL workflow + diagram + metadata

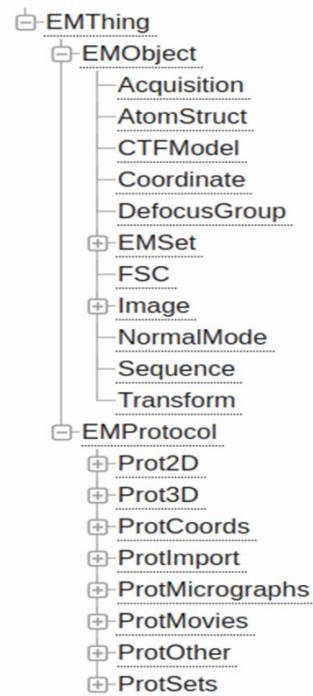


Figure 1. CryoEM ontology view from OLS

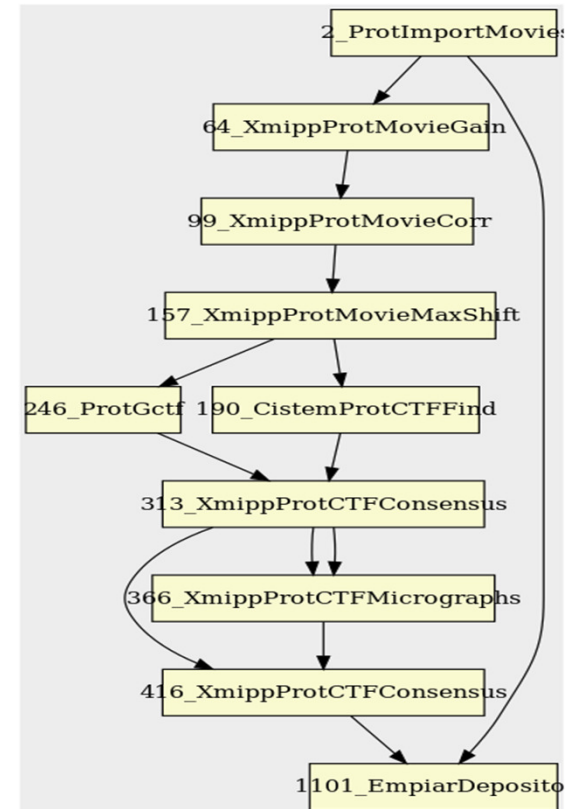


Figure 2. RO-Crate diagram

Mathematical meaning

Geometrical transformations can be represented by matrix operations between homogeneous coordinates:

$$\tilde{\mathbf{r}}_{\tilde{A}} = \tilde{A} \tilde{\mathbf{r}}, \quad (1.2)$$

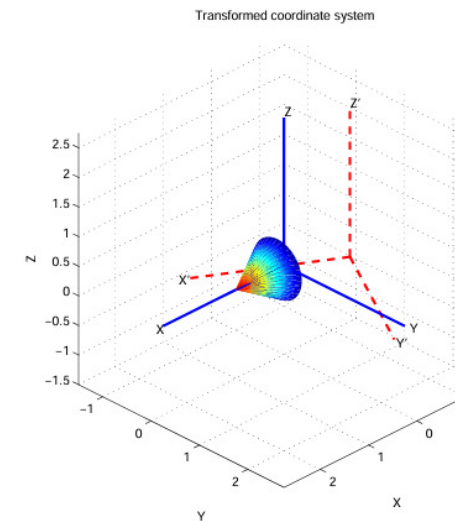
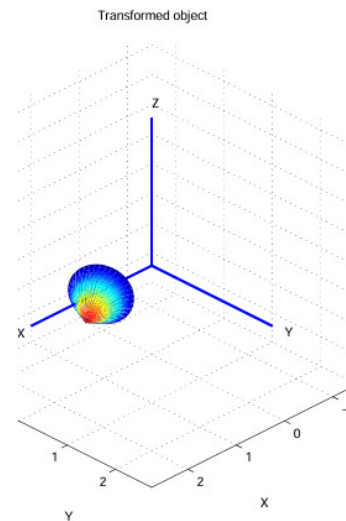
where $\tilde{\mathbf{r}} \in \mathbb{R}^3 \times \{1\}$ is the homogeneous coordinate of the point to transform, $\tilde{\mathbf{r}}_{\tilde{A}} \in \mathbb{R}^3 \times \{1\}$ is its transformed point in homogeneous coordinates, and \tilde{A} is a 4×4 invertible, real matrix of the form

$$\tilde{A} = \begin{pmatrix} r_{11} & r_{12} & r_{13} & t_x \\ r_{21} & r_{22} & r_{23} & t_y \\ r_{31} & r_{32} & r_{33} & t_z \\ 0 & 0 & 0 & 1 \end{pmatrix} = \begin{pmatrix} R & \mathbf{t} \\ \mathbf{0}^T & 1 \end{pmatrix}. \quad (1.3)$$

$$V_{\tilde{A}}(\tilde{\mathbf{r}}) = V(\tilde{A}^{-1} \tilde{\mathbf{r}})$$

$$I_{\tilde{A}}(\tilde{\mathbf{s}}) = \int_{-\infty}^{\infty} V_{\tilde{A}}(\tilde{H}^T \tilde{\mathbf{s}}) dt$$

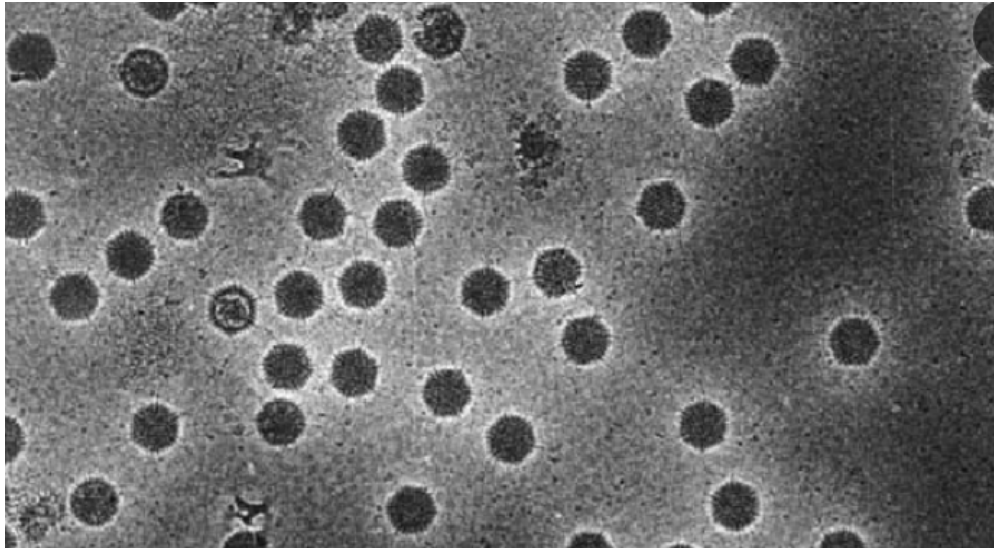
$$\tilde{H}^T = \begin{pmatrix} 1 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & t \\ 0 & 0 & 1 \end{pmatrix}.$$



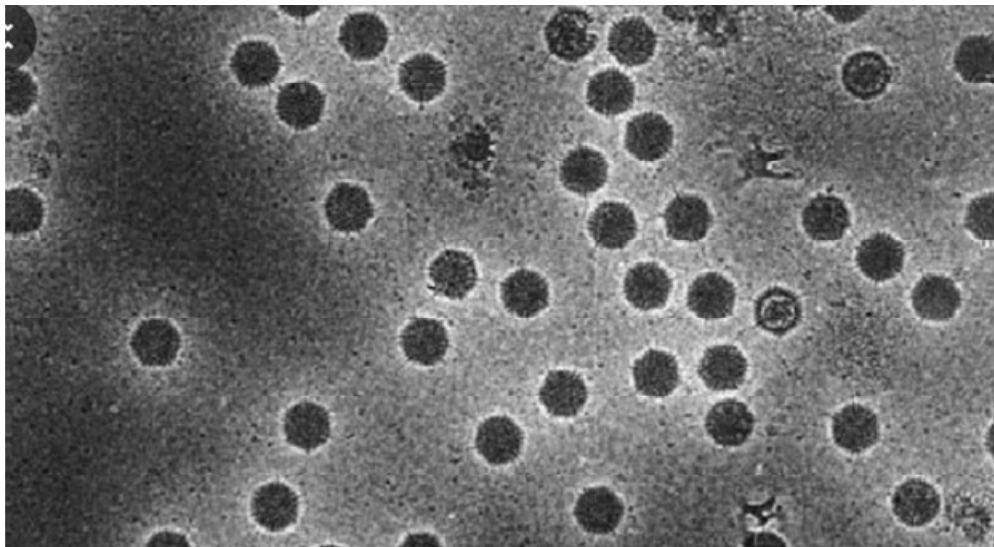
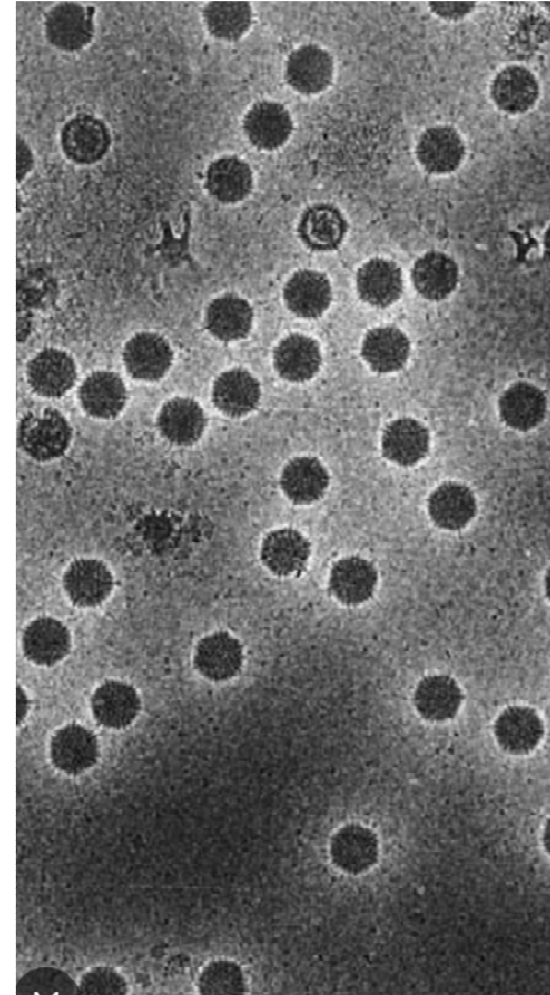
C.O.S. Sorzano, R. Marabini, J. Vargas, J. Oton, J. Cuenca-Alba, A. Quintana, J.M. de la Rosa-Trevin, J.M. Carazo. *Interchanging geometry information in electron microscopy single particle analysis: mathematical context for the development of a standard*. Computational Methods for Three-Dimensional Microscopy Reconstruction: 7-42 (2014) ([preprint](#))

Mathematical meaning

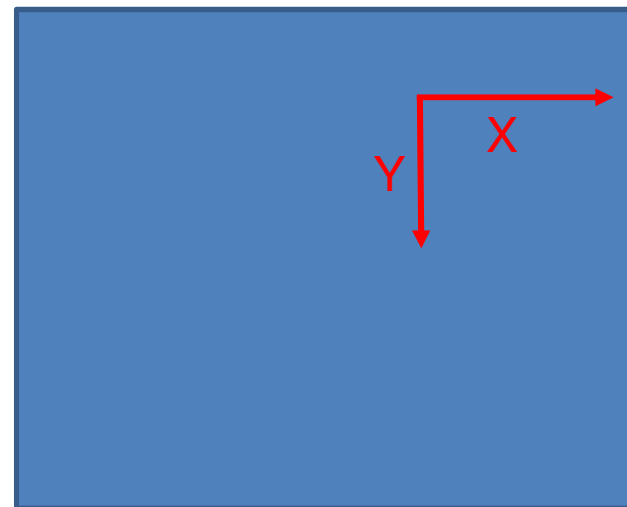
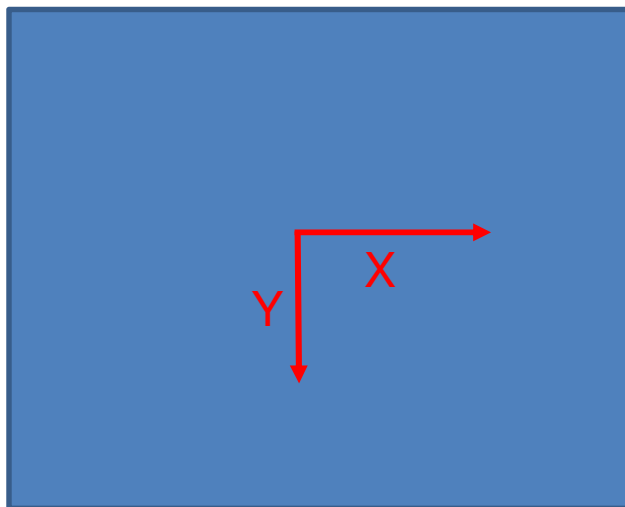
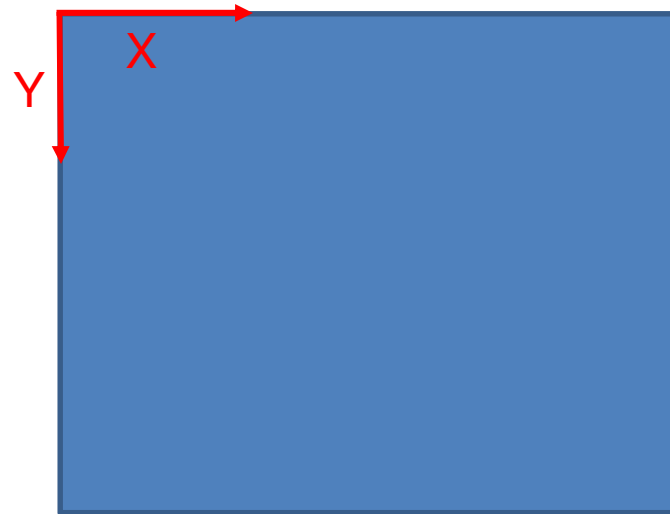
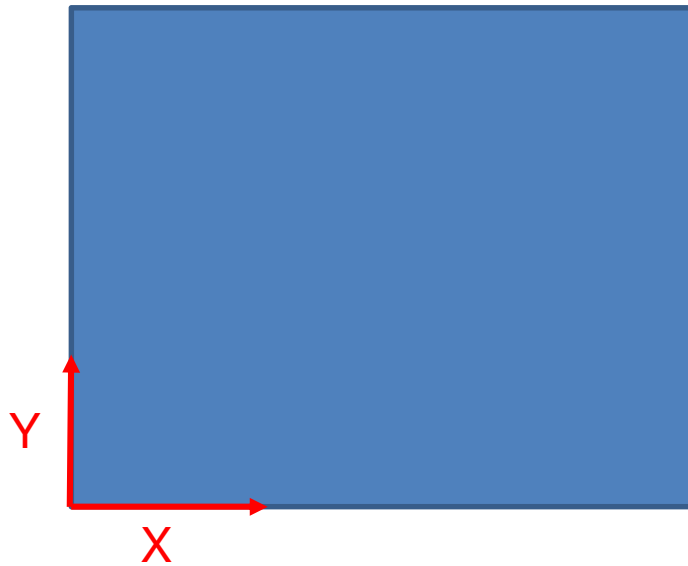
myMicrograph.mrc



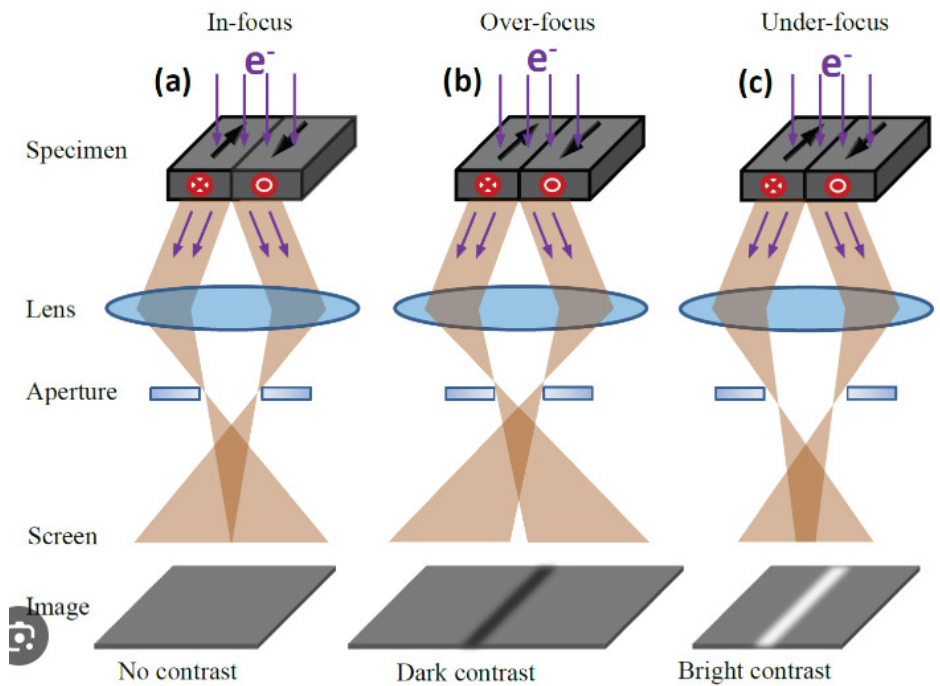
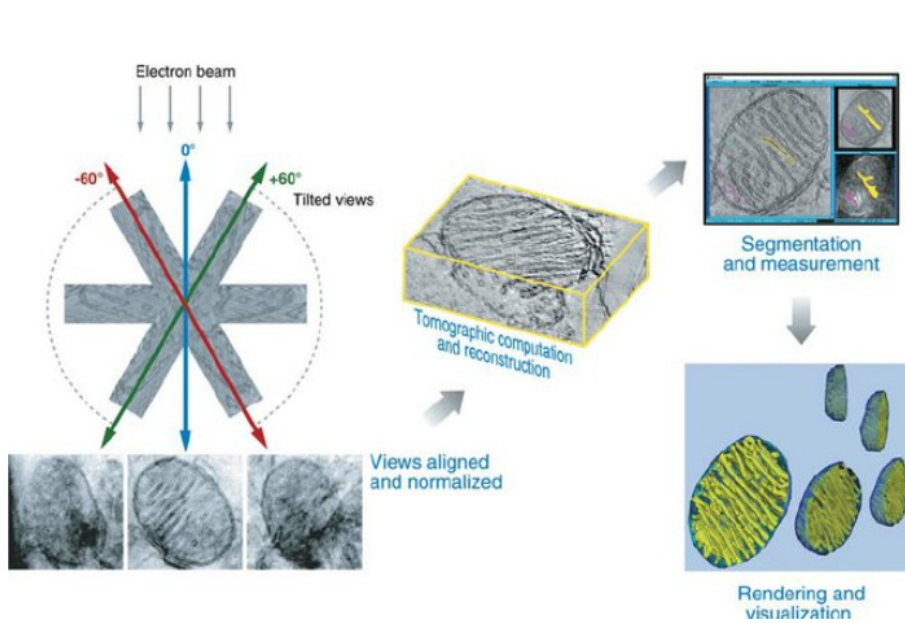
myMicrograph.tif



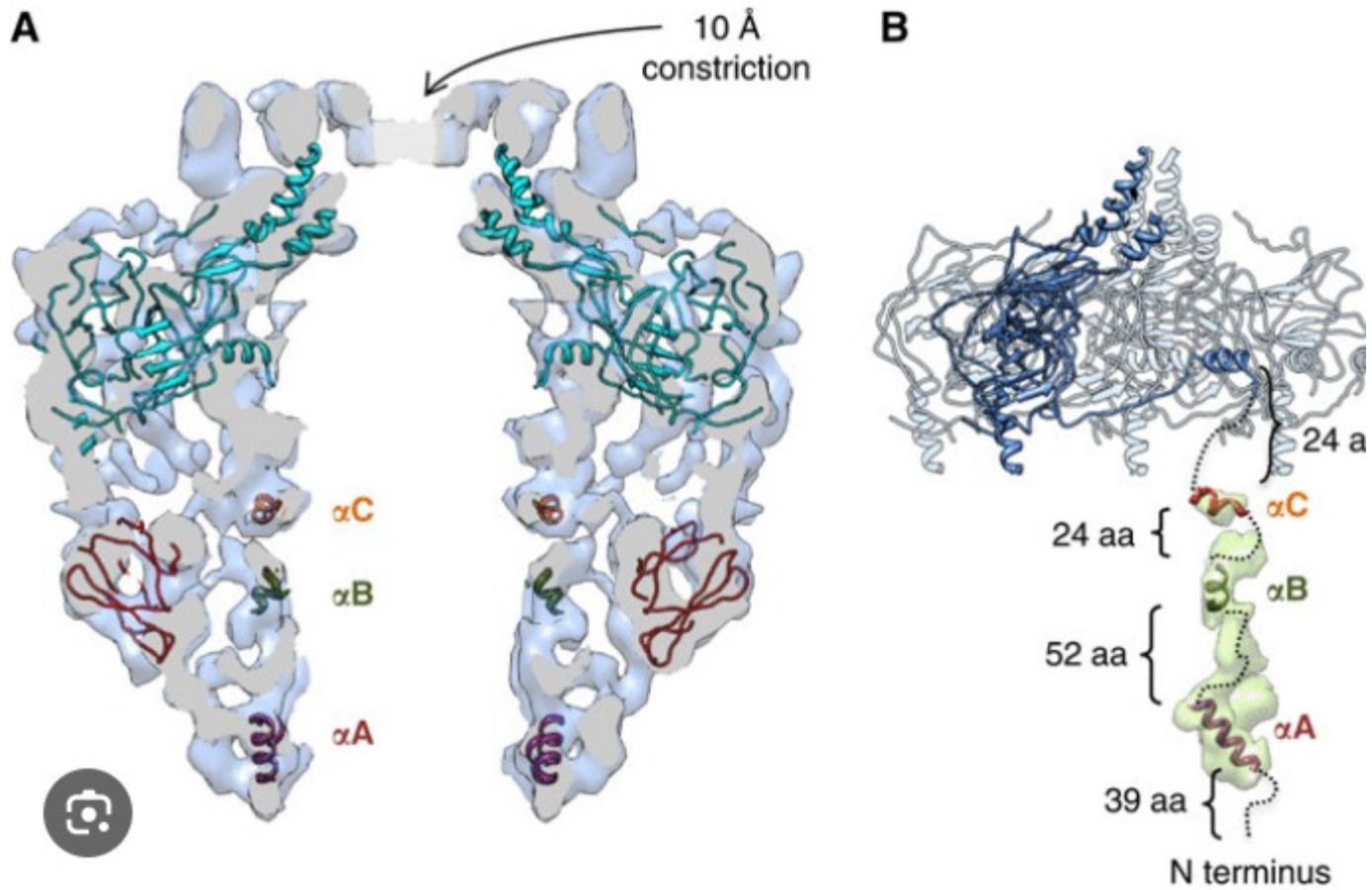
Mathematical meaning



Mathematical meaning



Mathematical modelling



EMPIAR Submission

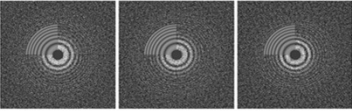
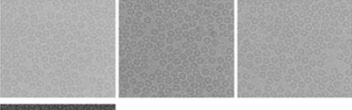
ebi.ac.uk/pdbe/emdb/empiar/entry/10516/scipion_workflow/data/SARS-CoV-2-spike/workflow.json

Scipion workflow viewer - EMPIAR-10516

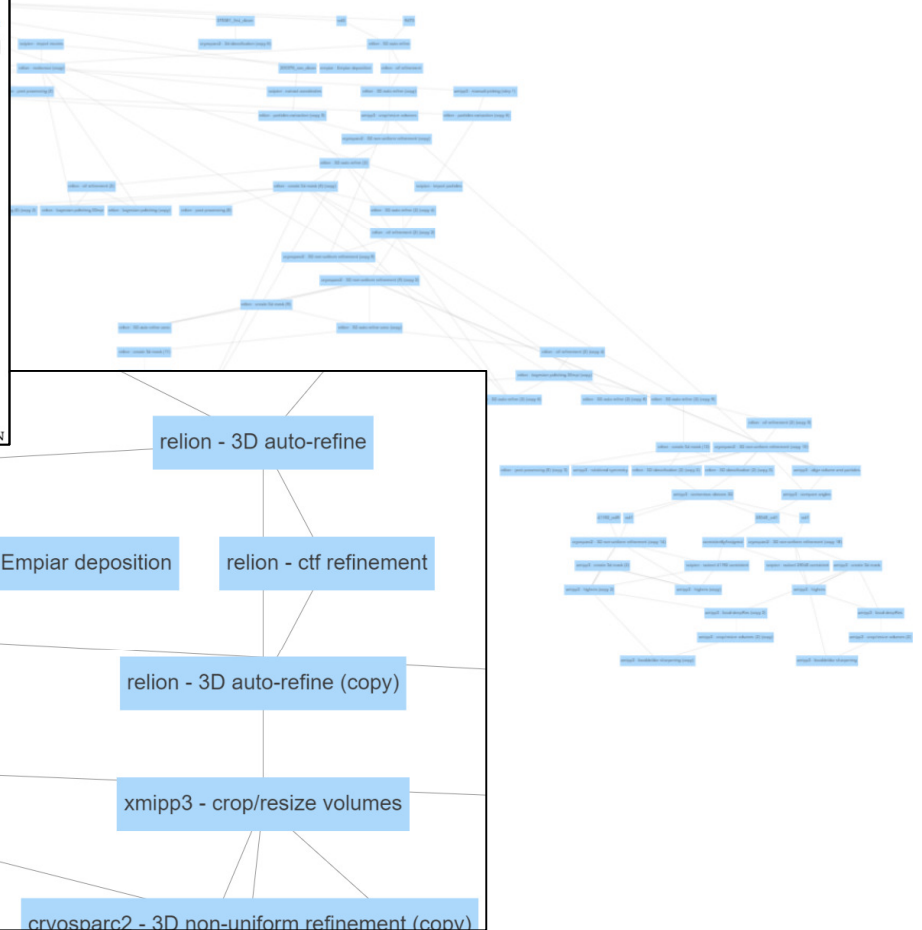
CryoEM workflow viewer

Select view mode: Summarized

xmipp3 - ctf consensus GCTF+CTFFind4

Param name	Param value
object.className:	XmippProcCTFConsensus
output:	
-6819.outputCTF:	
-6819.outputMicrographs:	

Powered by SCIPION



EMPIAR Submission

```
{
  "object.className": "ProtImportMovies",
  "object.id": "7262",
  "object.label": "pwem - import movies",
  "object.comment": "",
  "_useQueue": false,
  "_prerequisites": "",
  "_queueParams": null,
  "runName": null,
  "runMode": 0,
  "importFrom": 0,
  "filePath": "/home/coss/ScipionUserData/projects/Example_10248_Scipion3/EMPIAR/",
  "filesPattern": "*.tiff",
  "copyFiles": false,
  "haveDataBeenPhaseFlipped": false,
  "acquisitionWizard": null,
  "voltage": 300.0,
  "sphericalAberration": 2.7,
  "amplitudeContrast": 0.1,
  "magnification": 50000,
  "samplingRateMode": 0,
  "samplingRate": 0.495,
  "scannedPixelSize": 7.0,
  "doseInitial": 0.0,
  "dosePerFrame": 1.0,
  "gainFile": "/home/coss/ScipionUserData/projects/Example_10248_Scipion3/EMPIAR//gain.mrc",
  "darkFile": null,
  "dataStreaming": false,
  "timeout": 43200,
  "fileTimeout": 30,
  "blacklistDateFrom": null,
  "blacklistDateTo": null,
  "useRegexps": true,
  "blacklistFile": null,
  "inputIndividualFrames": false,
  "numberOfIndividualFrames": null,
  "stackFrames": false,
  "writeMoviesInProject": false,
  "movieSuffix": "_frames.mrcs",
  "deleteFrames": false
},
{
  "object.className": "XmippProtMovieGain",
  "object.id": "7393",
  "object.label": "xmipp3 - movie gain",
```

Conclusions

- The internal data model of each processing package is very detailed.
- They are incompatible among programs.
- Purpose is important: future or past
- There are multiple levels of detail
- The standard should be compatible with both (shared minimum and individual extensions)

Thanks



I²PC



Laura del Caño



Carolina Simón



Irene Sánchez