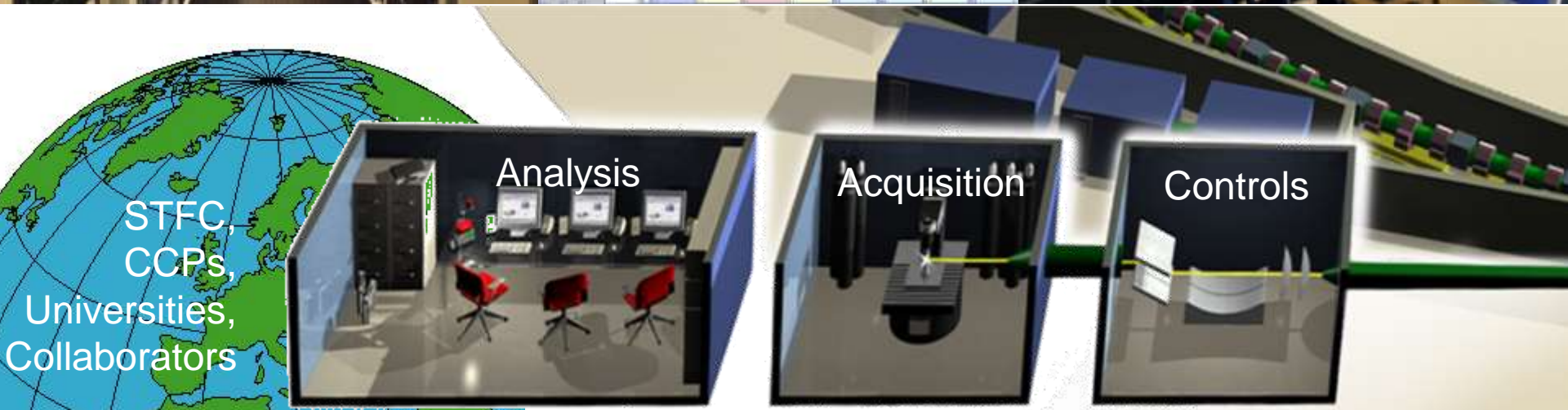
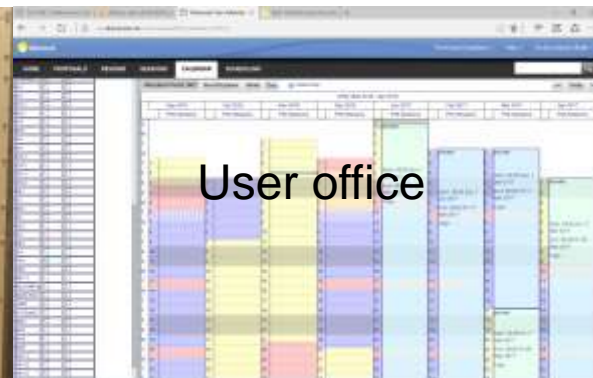


# Processing at eBIC/Diamond Light Source

Alun Ashton

# Computing/Software Support Groups



# Data Analysis – Experimental Challenges

<u><b>BEFORE</b></u> (FROM DLS/USER'S INSTITUTION)	<u><b>IMMEDIATE</b></u> (DURING EXPERIMENTS)	<u><b>SHORT TERM</b></u> (BEFORE THE USER GOES HOME)	<u><b>LONG TERM</b></u> (FROM DLS/USER'S INSTITUTION)
Simulations  Processing of older datasets	"Real time" data processing, analysis and visualisation – to make experimental decisions	Data reduction and processing – Users go home with clean data free of instrument artefacts.  Preliminary data analysis – helpful, but may require significant processing power and know-how	Detailed analysis – from data to information.  Incorporating results from other techniques.  Experiments: <ul style="list-style-type: none"><li>➤ Provide parameters for a model.</li><li>➤ Test/verify a model or theory.</li><li>➤ Show where a new theory or model is required.</li></ul>

Database ADM



Core Software



Core Analysis



MX,X



Scattering & Spectroscopy



Imaging



Software Engineers

Software Scientists

Other Funding

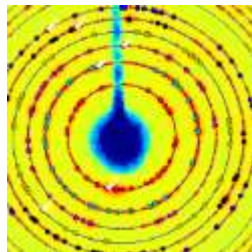
# Imaging Team

The logo for Savu, featuring a large yellow 'S' and the word 'avu' in blue.

Multi-dimensional and Multi-Modal  
tomographic reconstruction pipeline



Machine learning assisted annotation and  
segmentation for volumetric data.

The logo for ptypy, featuring the word 'ptypy' in blue with a stylized circular graphic behind it.

Extensions to Savu for  
Multimodal data collections

SCIPION RELION CCP-EM

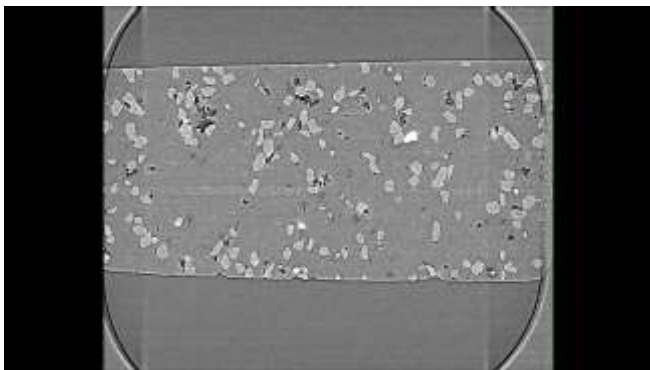


FEI  
part of Thermo Fisher Scientific

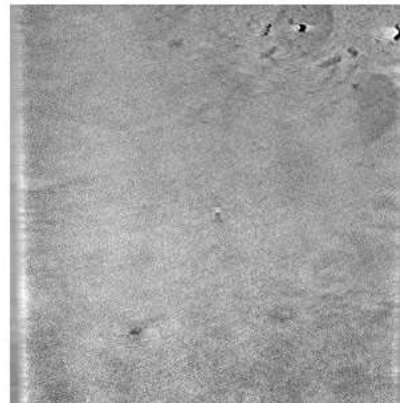
Automation of data processing for electron  
tomography and single particle analysis.



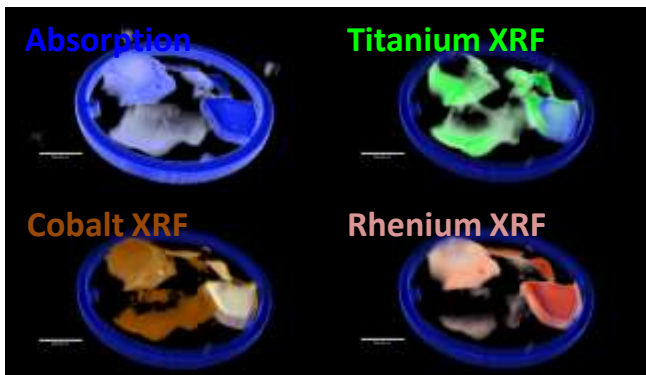
# Imaging Team



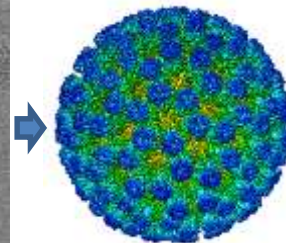
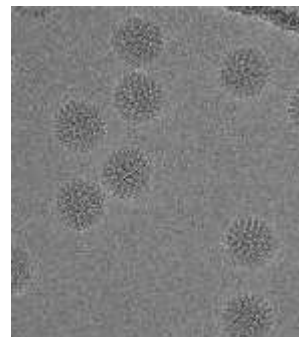
Thanks to Robert Atwood



Thanks to Michele Darrow



Thanks to Stephen Price





# LiMS:ISPyB/SynchWeb



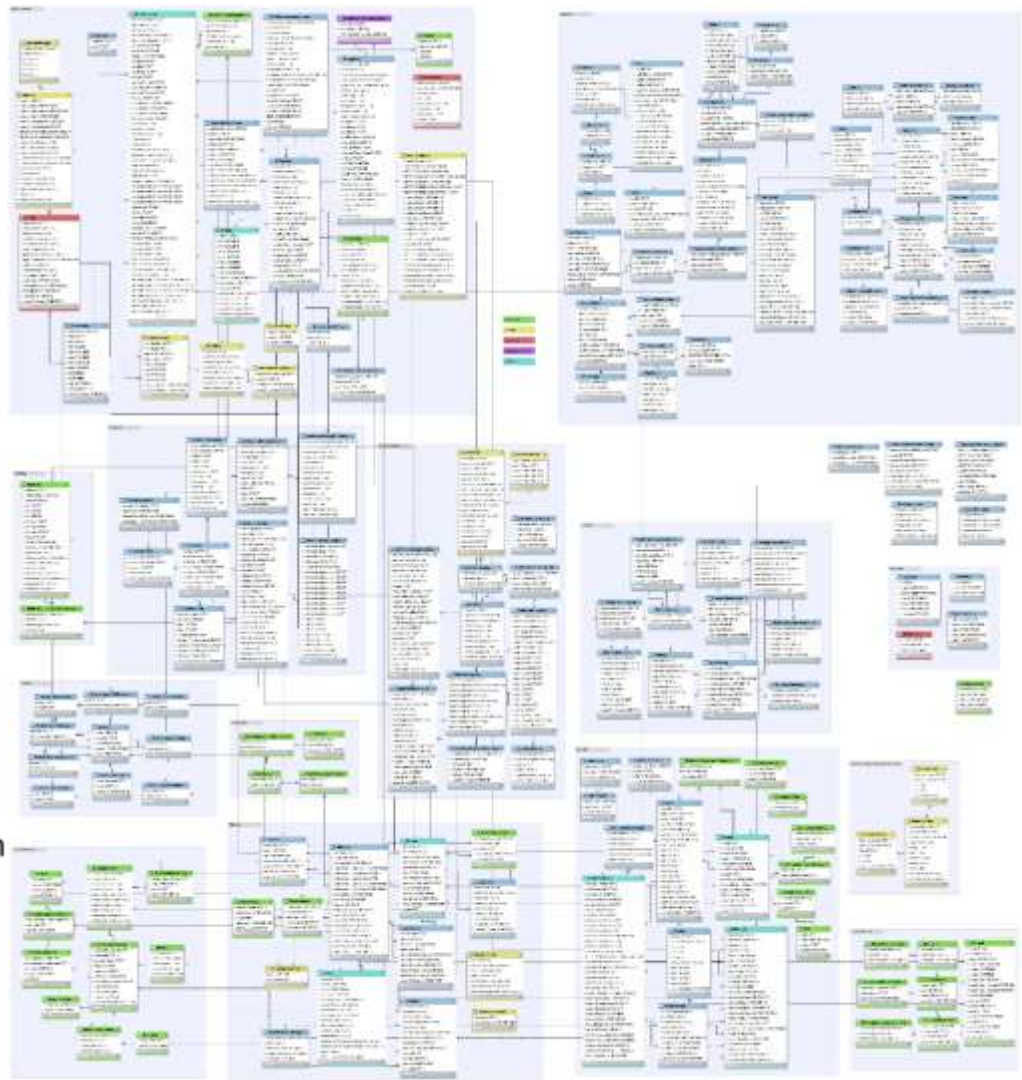
- MariaDB/MySQL-aware intelligent proxy
- Load balancing
- Monitoring mechanisms
- Logging, Filtering
- Multi-master – 3 nodes
- Each node has the full dataset
- Synchronous replication



# LiMS:ISPyB/SynchWeb

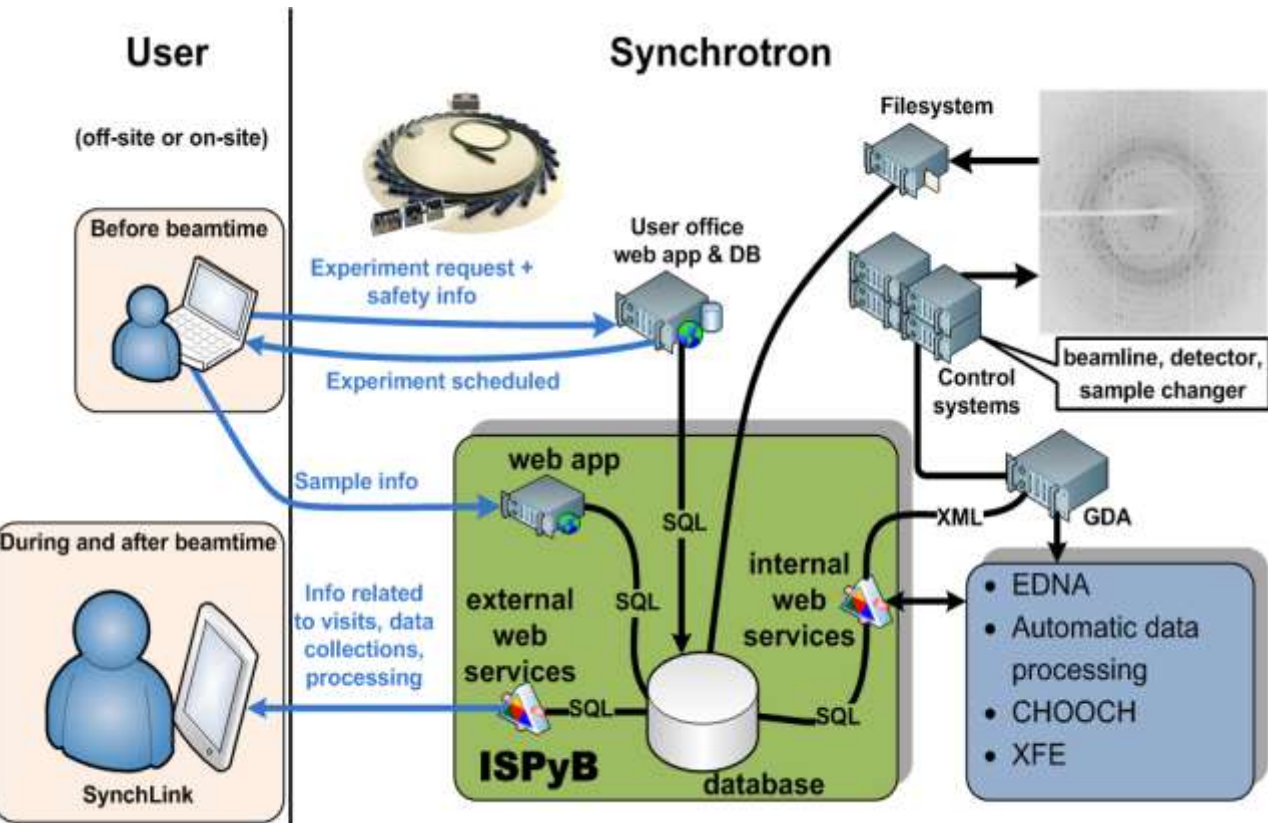


- MariaDB/MySQL-aware intelligent proxy
- Load balancing
- Monitoring mechanisms
- Logging, Filtering
- Multi-master – 3 nodes
- Each node has the full dataset
- Synchronous replication





# LiMS:ISPyB/SynchWeb



# LiMS:ISPyB/SynchWeb

User

(off-site or on-site)

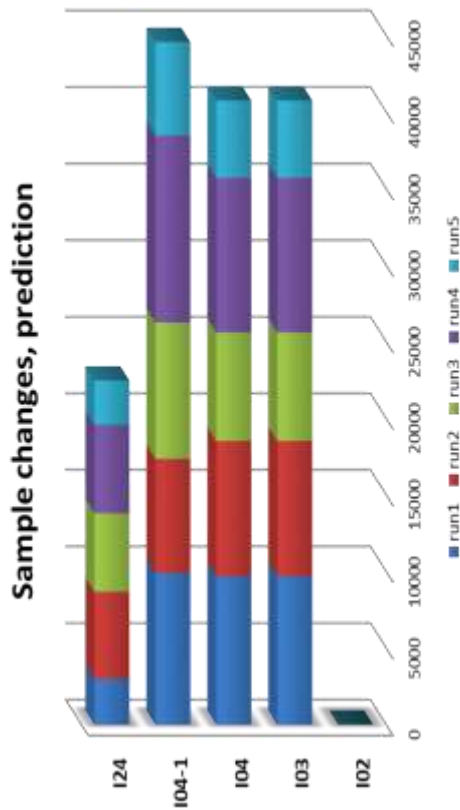
Before beamtime

Synchrotron

User of



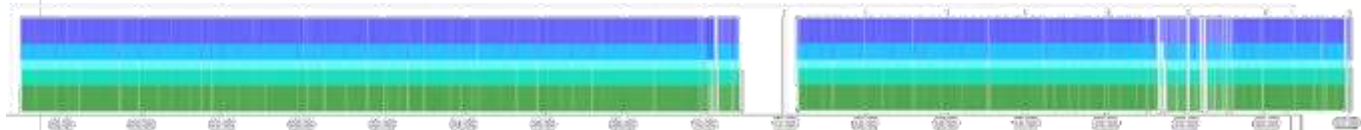
# Macromolecular Crystallography Automation



- Use case: fragment screening
  - Crystals mounted at synchrotron (100s/day)
  - Auto-collect (visual centring)
  - Critical: Capacity, speed, robustness
  - Not critical: small footprint...



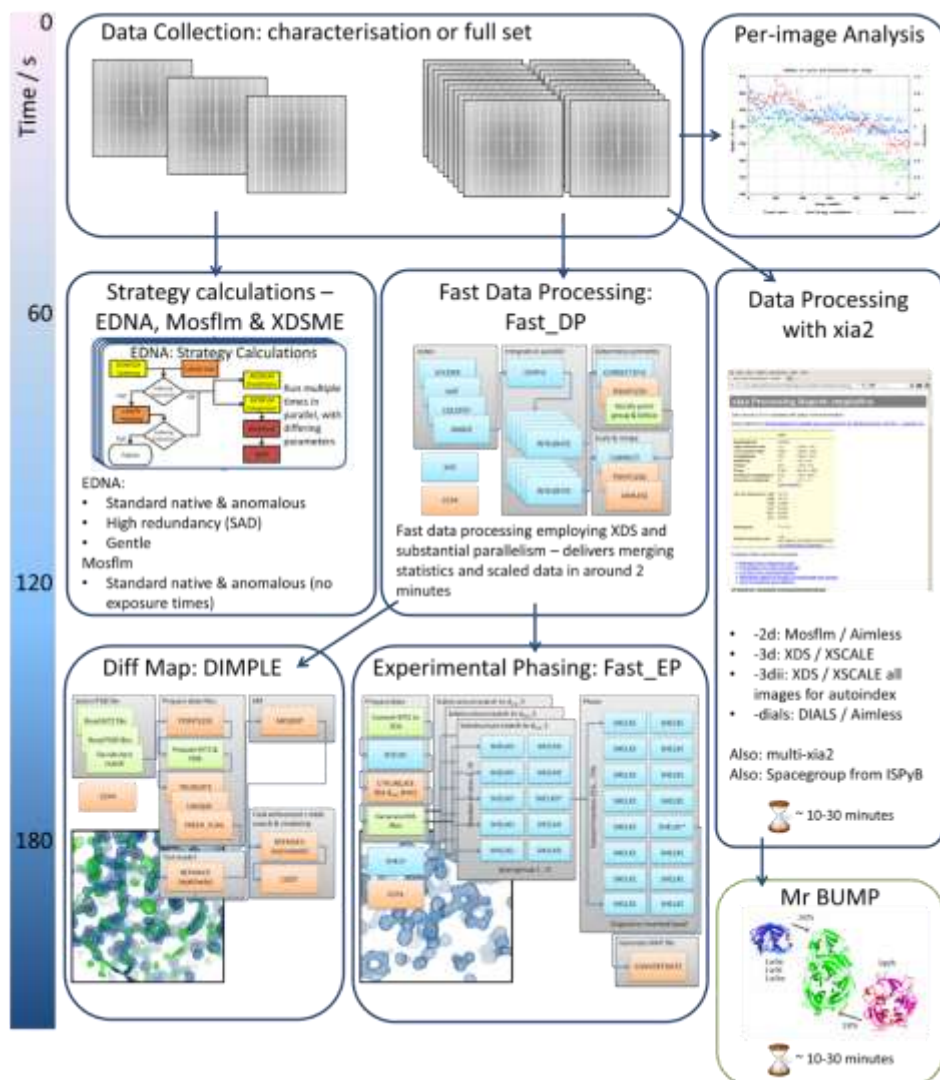
*“67 pucks, 978 data collections, zero problems - in about 40 hours”*



Before: ~20 datasets / hr

After: ~31/hr





## For eBIC

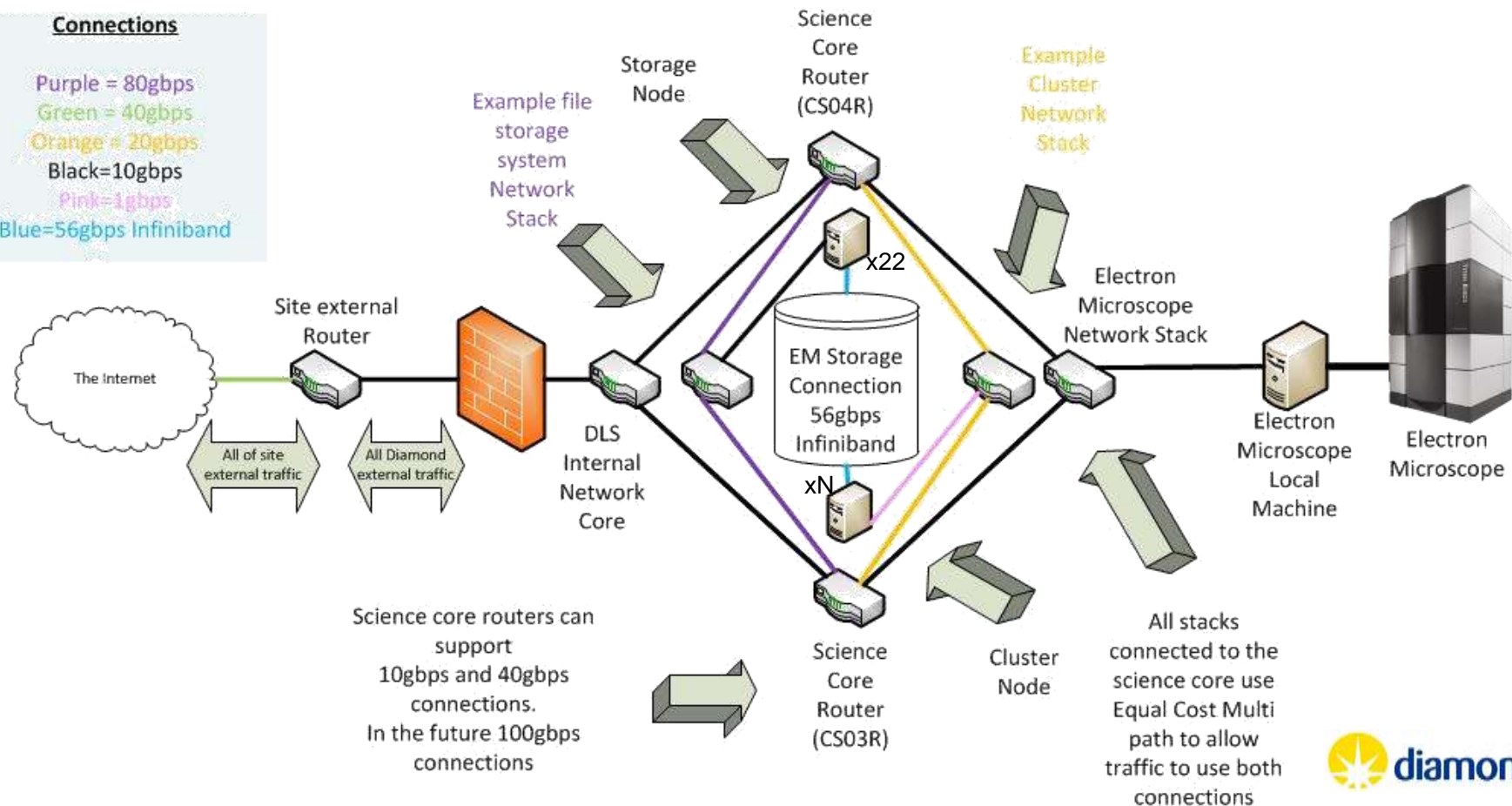
- Building on what we already have
- Produce automated pipelines for
  - Single particle analysis
  - Tomography
- Prioritise feedback useful at time of data collection
- Collect data into ISPyB to facilitate interoperability of data between measurement types
- View data online from SynchWeb



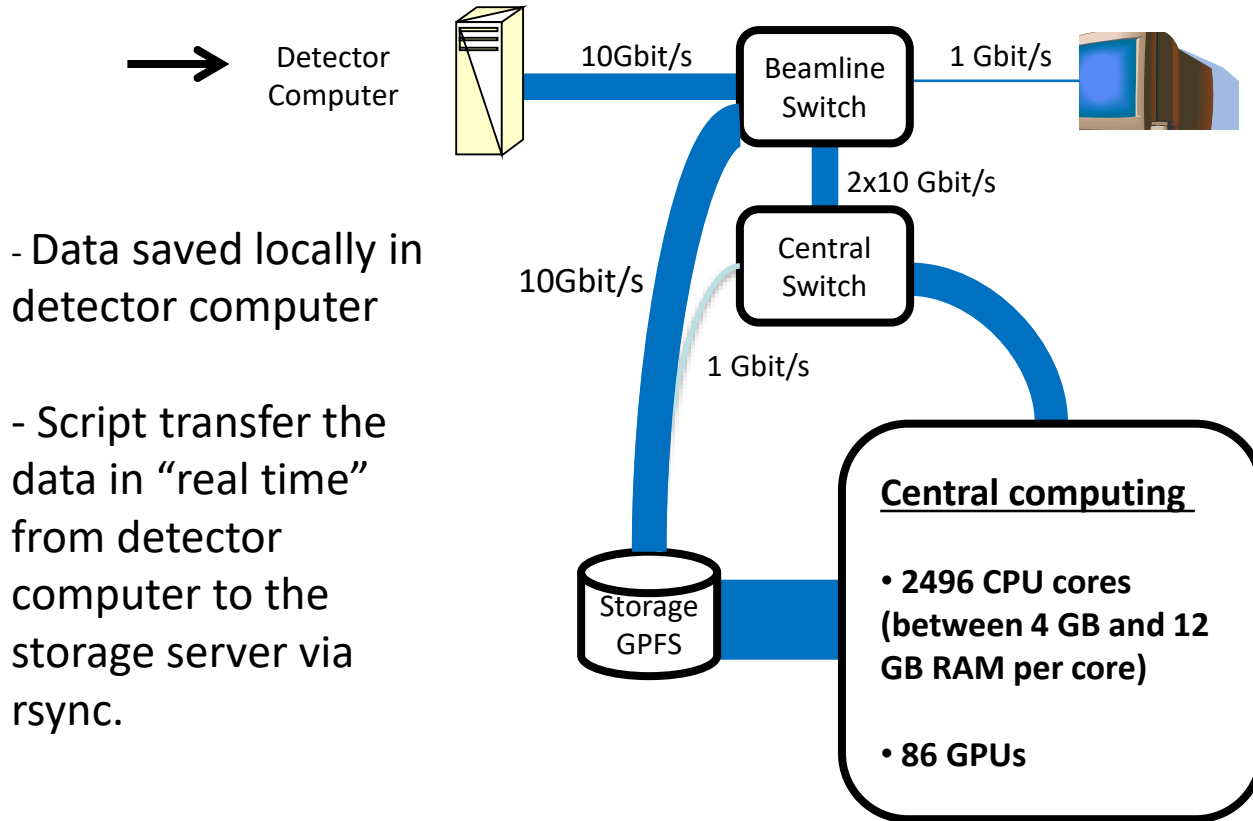
# Electron Microscope Connectivity

## Connections

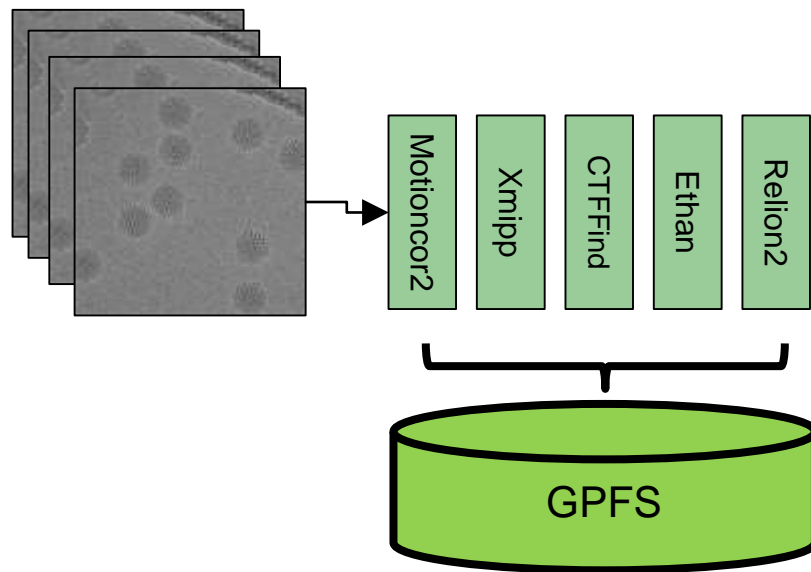
Purple = 80gbps  
Green = 40gbps  
Orange = 20gbps  
Black = 10gbps  
Pink = 1gbps  
Blue = 56gbps Infiniband



# Outline of current computing resources



# Single Particle Analysis



# Single Particle Analysis

Project TutorialIntro (nnp95151 on ws109.diamond.ac.uk)

Project Help

SCIPION devel (2016-08-02) cce98b9 Project TutorialIntro Protocols | Data

View: Protocols SPA

- Imports
  - import movies
  - import micrographs
  - import particles
  - import volumes
- more
- Micrographs
  - xmipp3 - optical alignment
  - grignonellab - unblur
  - grignonellab - summovie
  - xmipp3 - preprocess micrographs
- CTF estimation
- Particles
  - Picking
  - Extract
  - Preprocess
  - Filter
  - Mask
- 2D
  - Align
  - Classify
- 3D
  - Initial volume
  - Preprocess
  - Refine
  - Classify
  - Analysis
  - Reconstruct
- Tools

1. Import mics 2 finished

2. downsample x5 (copy) finished

xmipp3 - ctf estimation finished

xmipp3 - manual-picking (step 1) interactive

xmipp3 - auto-picking (step 2) finished

xmipp3 - extract particles finished

relicon - 2D classification finished

Analyse Results

Summary Methods Output Log

Input

InputMicrographs (from 2. downsample x5 (copy) -> out) SetOfMicrographs (3 items, 1843 x 1888, 5.0)

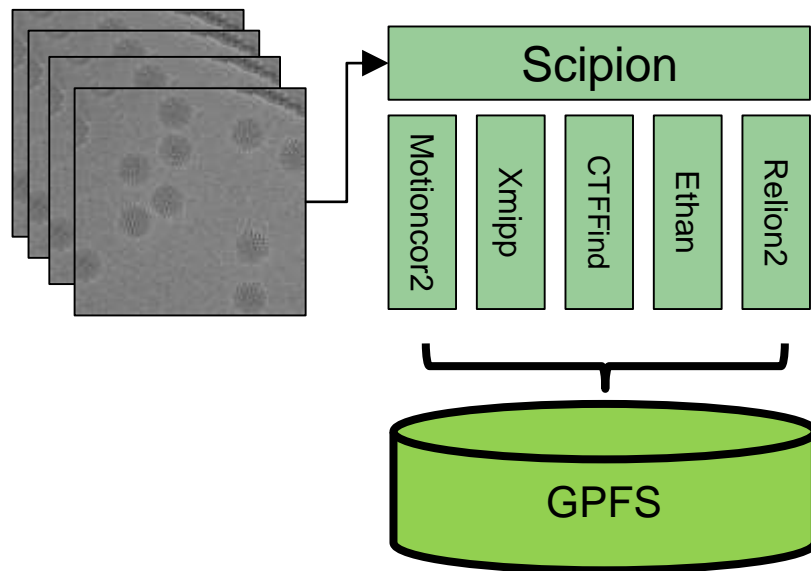
Output

xmipp3 - ctf estimation -> outputCTF SetOfCTF (3 items)

SUMMARY

CTF estimation of 3 micrographs.

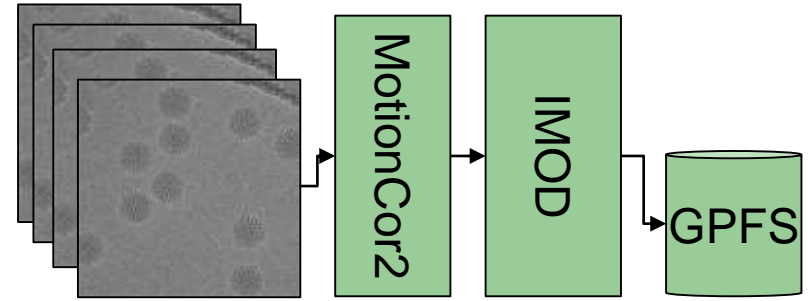
The range of micrograph's experimental defocus was 1.000 - 2.243 microns.



iNEXT

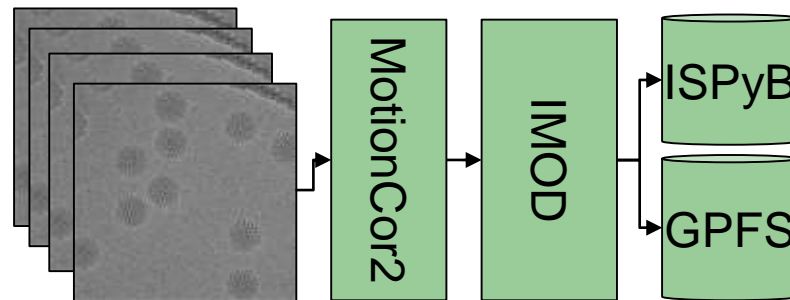


# Tomography

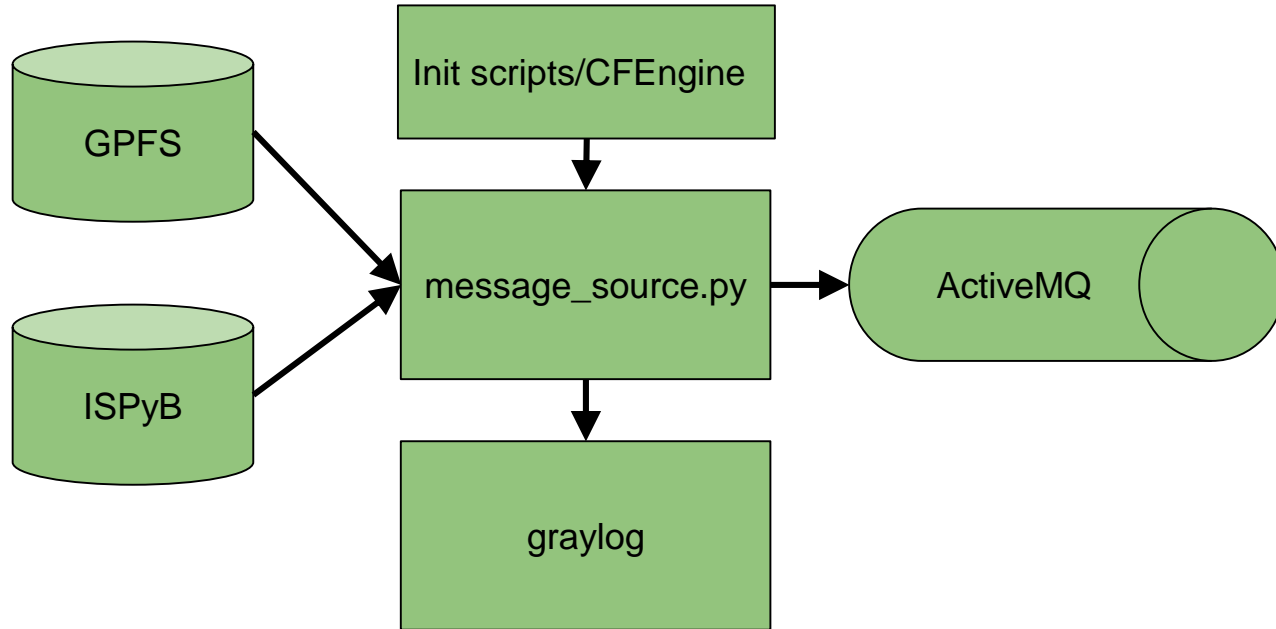




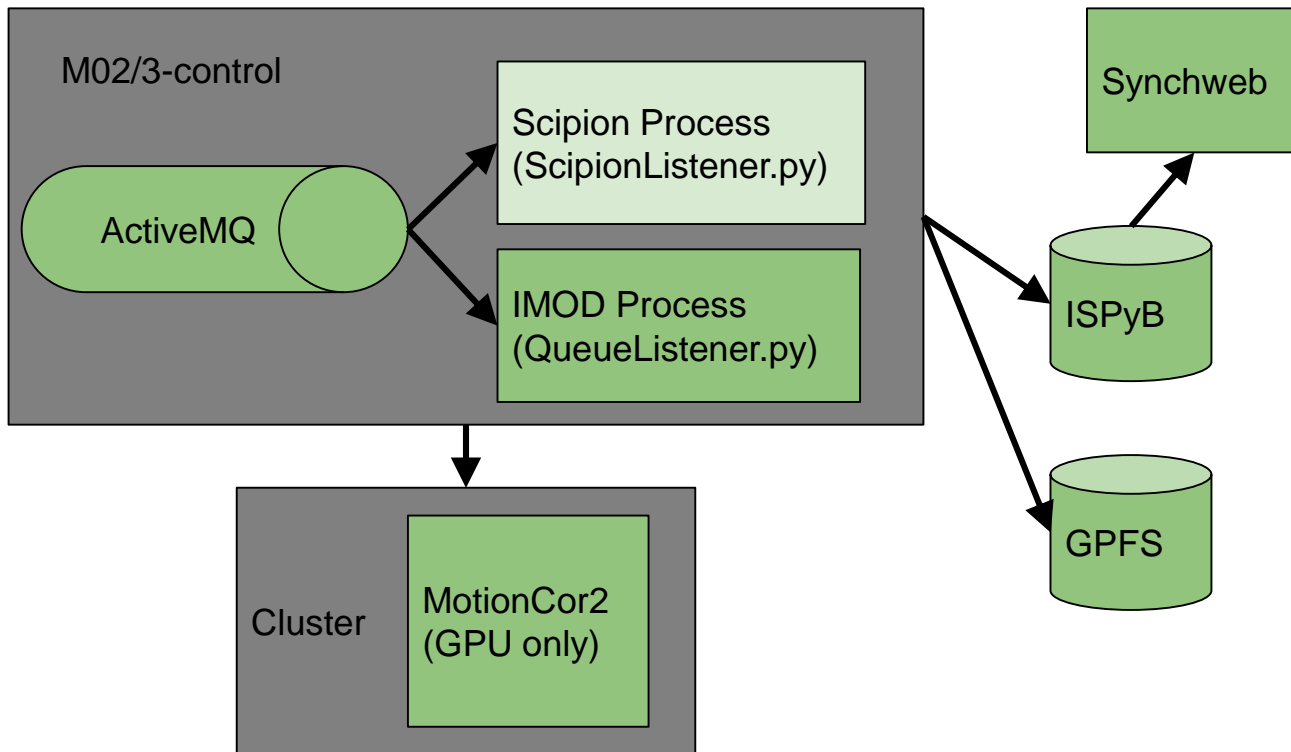
# Tomography: IMOD



# Automating : Notification



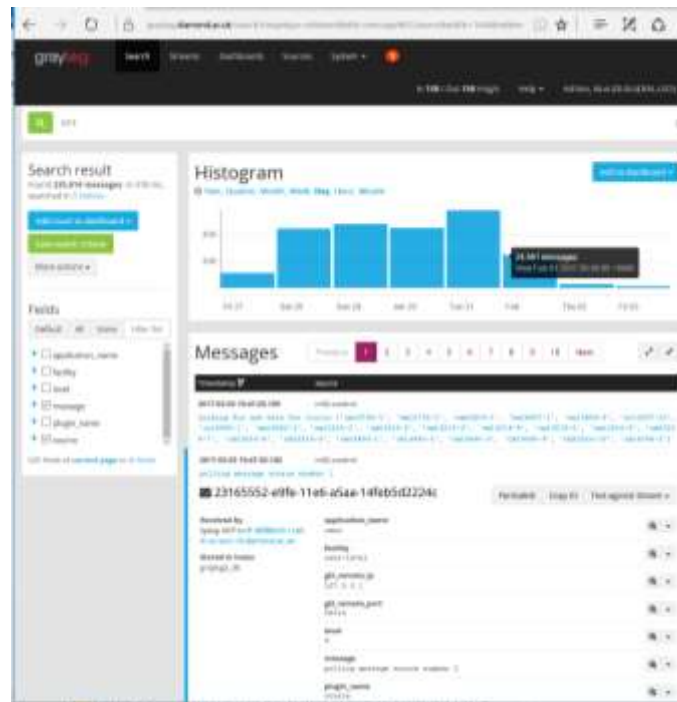
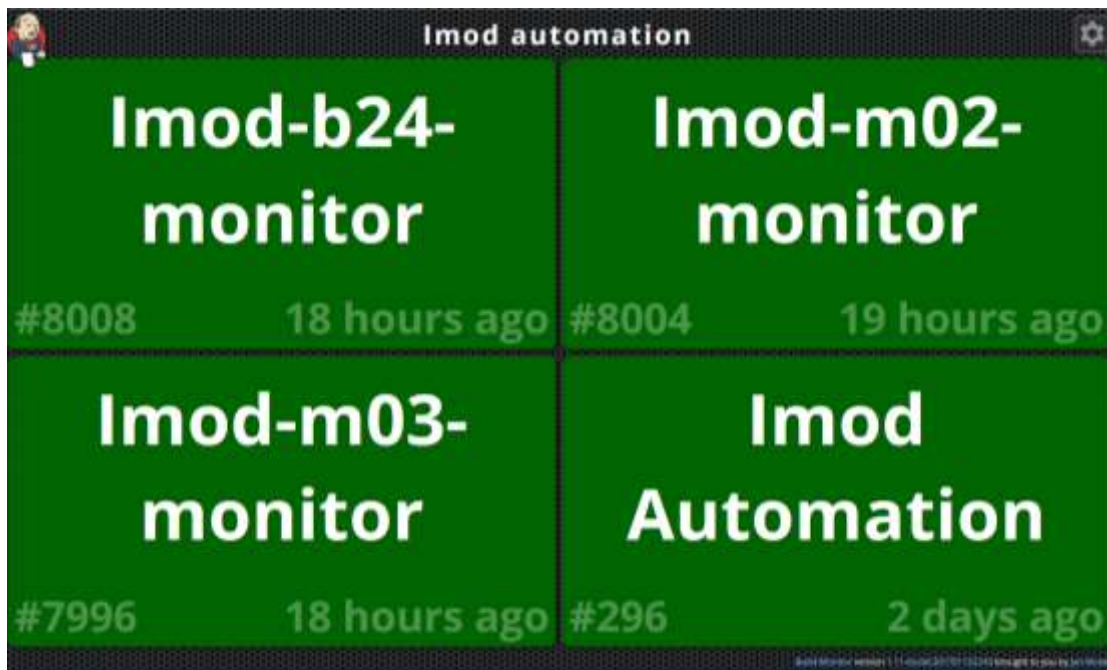
# Automating : Processing



# Current Status

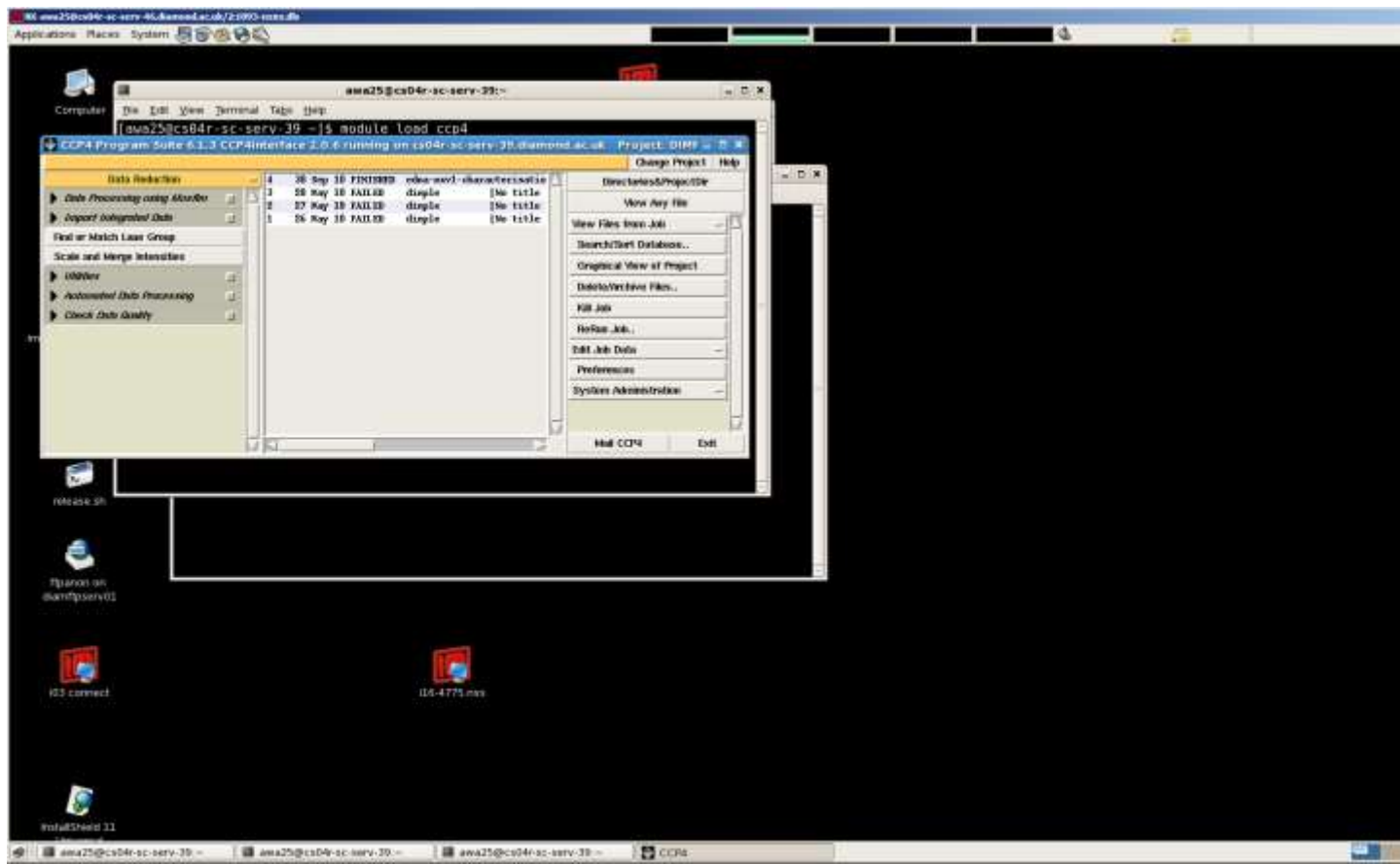
- Tomography running for some time
  - also running on full field cryo-Xray microscopy beamline
  - Scipion setup, modified and well used
- Automatically triggered Scipion went live 30<sup>th</sup> January 2017
- Some data in ISPyB

# Builds and Services



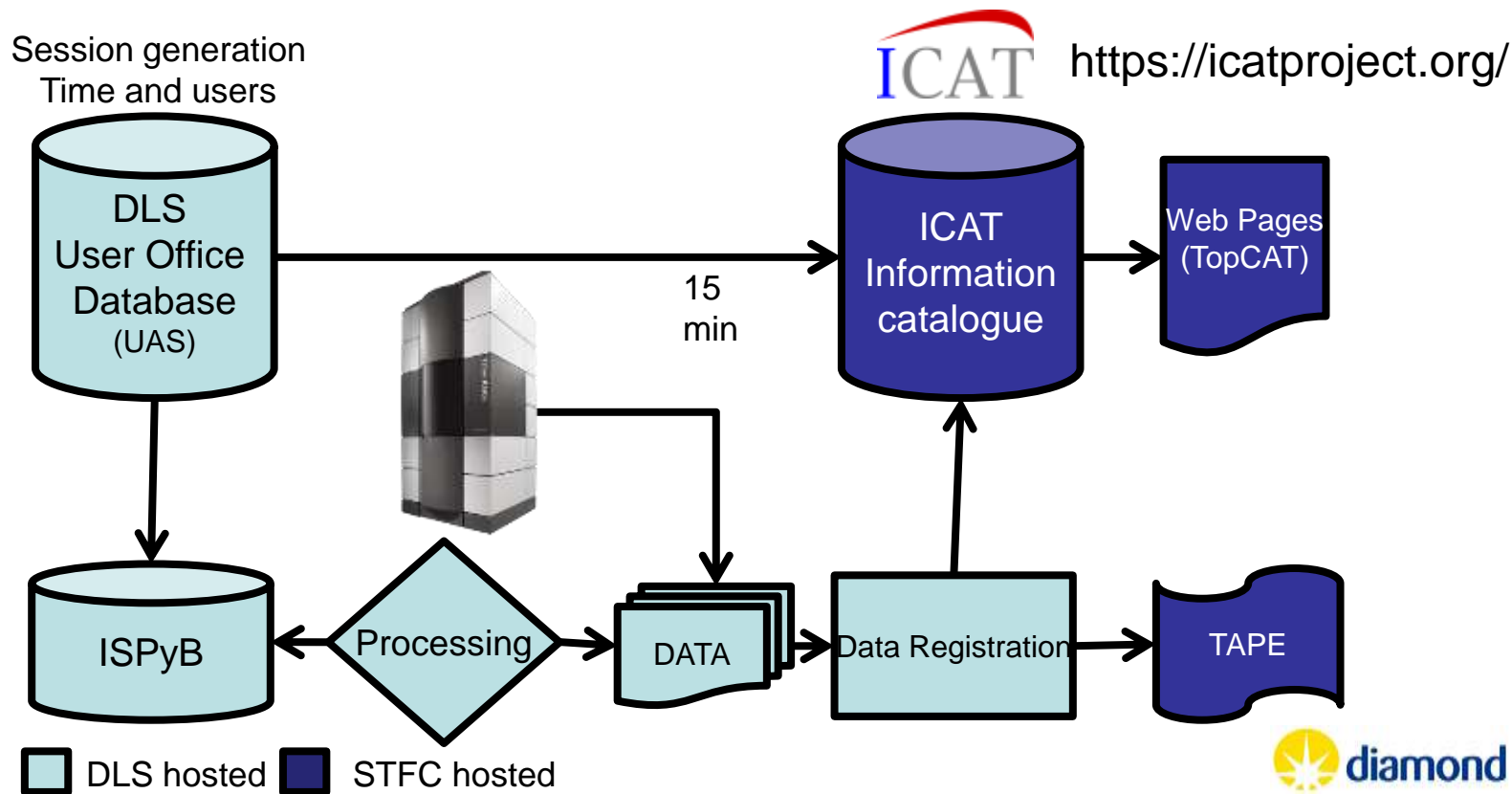


# Remote Reprocessing



Fair play policy.

# Data is archived and metadata captured from every stage of the process

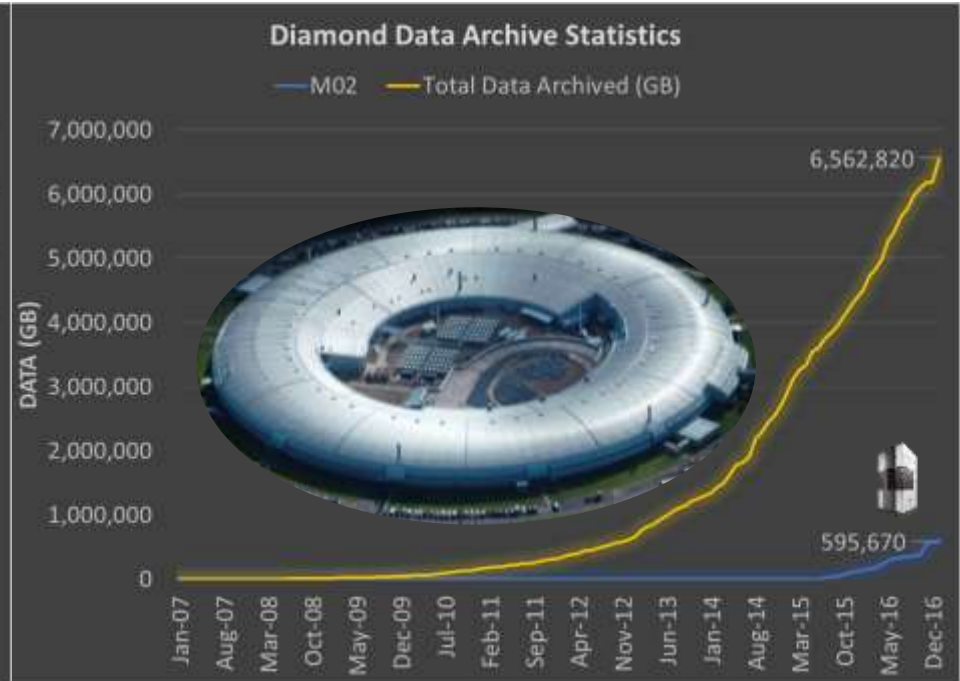


# EM data flow



- Data capture at full rate on both direct electron detectors (e.g. 17 fps for Falcon II or 40 fps for the K2).
- Falcon II data rate ~50-100 movies/hr
- K2 data rate ~25-75 movies/hr
- All data are directly written to high speed central computing/storage facility.
- All data are archived to tape and stored for the lifetime of the media.
- Diamond cluster available to external users particularly during beamline shut downs....

# EM data flow



# Lessons learnt

- Persistent services : CFEngine
  - can make sure that a process is always running
  - Scientific computing helped setting it up
- Parallel file system : GPFS
  - Really fast at loading data
  - Does not like polling (we are going to change that bit)
- Workflow : Scipion
  - Designed as a GUI based tool
  - We were able to make modifications to Scipion to run “headless”
  - It’s nice on the inside!



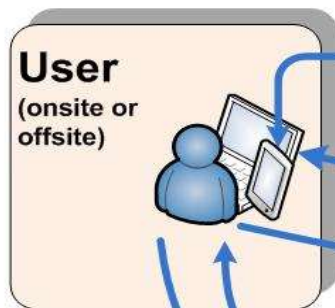
# Current challenges

- Scipion headless is very 'fresh'
  - Working on the actual workflow steps
- We would like ISPyB data to be richer
- Processing goes wrong:
  - Fiducial Alignment
  - Particle Picking
- More microscopes
  - Data volumes

## Home Laboratory



Samples



<sup>1</sup>Pre Experiment

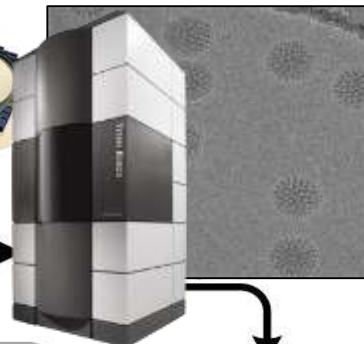
<sup>2</sup>During/After Experiment

## Synchrotron/EM Facility



Samples

Laboratory sample treatments

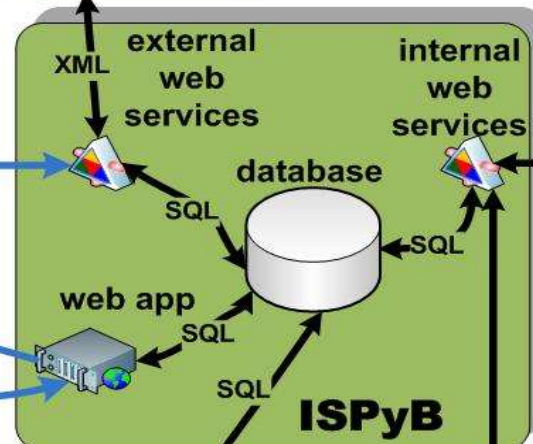


Filesystem



Event handler/  
experiment recording

Automatic/on  
the fly data  
processing



Info related  
to visits, data  
collections,  
processing<sup>2</sup>

Experiment files<sup>2</sup>

Sample info<sup>1</sup>

Experiment scheduled<sup>1</sup>

Experiment request<sup>1</sup>

User Admin  
System

# ISPyB Data Structure

Tomography is currently populating:

BLSAMPLEID, SESSIONID

startTime, endTime

ImageDirectory, fileTemplate

numberOfImages

XtalSnapshotFullPath[1-4] – converted to jpeg to display in SynchWeb

axisStart, axisEnd – angles

# ISPyB Data Structure

Single particle will populate:

**Sample** : name, code, comments (i.e. no additional columns)

**DC group** : sample ID, session ID, experiment type, start & end time (i.e. no additional columns)

## **DC - MX/code columns :**

dcg ID

session ID

sample ID

detector ID

imgdir

file\_template

xtal\_snapshot1

xtal\_snapshot2

xtal\_snapshot4

starttime

endtime

n\_images

exp\_time

run\_status

resolution

comments

## **DC - special EM columns :**

dat\_file

magnification (unit: X)

total\_absorbed\_dose (Unit: e-/A<sup>2</sup> for EM)

binning (1 or 2. Number of pixels to process as 1. (Use mean value.)

particle\_diameter (Unit: nm)

box\_size\_ctf (Unit: pixels)

min\_resolution (Unit: Å)

min\_defocus (Unit: Å)

max\_defocus (Unit: Å)

defocus\_step\_size (Unit: Å)

amount\_astigmatism (Unit: Å)

extract\_size (Unit: pixels)

bg\_radius (Unit: nm)

voltage (Unit: kV)

obj\_aperture (Unit: um)

c1aperture (Unit: um)

c2aperture (Unit: um)

c3aperture (Unit: um)

c1lens (Unit: %)

c2lens (Unit: %)

c3lens (Unit: %)

# Potential areas for standardisation

- Metadata structures and API
- Sample and container IDs

# Single Particle Analysis

Project TutorialIntro (nnp95151 on ws109.diamond.ac.uk)

Project Help

SCIPION devel (2016-08-02) cce98b9 Project TutorialIntro Protocols Data

View: Protocols SPA

- Imports
  - import movies
  - import micrographs
  - import particles
  - import volumes
- more
- Micrographs
  - xmipp3 - optical alignment
  - grignonellab - unblur
  - grignonellab - summovie
  - xmipp3 - preprocess micrographs
- CTF estimation
- Particles
  - Picking
  - Extract
  - Preprocess
  - Filter
  - Mask
- 2D
  - Align
  - Classify
- 3D
  - Initial volume
  - Preprocess
  - Refine
  - Classify
  - Analysis
  - Reconstruct
- Tools

1. Import mics 2 finished

2. downsample x5 (copy) finished

xmipp3 - ctf estimation finished

xmipp3 - manual-picking (step 1) interactive

xmipp3 - auto-picking (step 2) finished

xmipp3 - extract particles finished

relicon - 2D classification finished

Analyse Results

Summary Methods Output Log

Input

InputMicrographs (from 2. downsample x5 (copy) -> out) SetOfMicrographs (3 items, 1843 x 1888, 5.0)

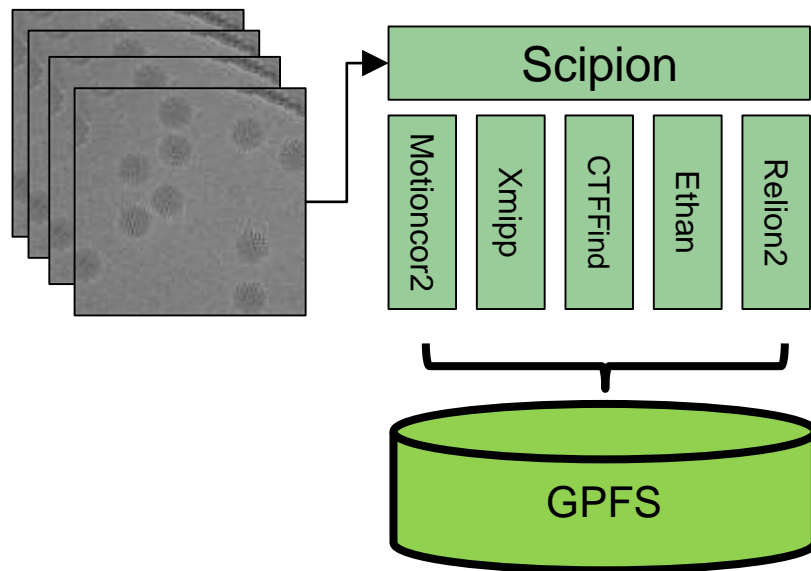
Output

xmipp3 - ctf estimation -> outputCTF SetOfCTF (3 items)

SUMMARY

CTF estimation of 3 micrographs.

The range of micrograph's experimental defocus was 1.000 - 2.243 microns.



iNEXT





# Relion2

- 12 nodes (20 CPU, 256GB RAM, 2xK80 NVIDIA cards, infiniband connections)

- Standard 3D classification tests:

- 120 CPUs (6 x com12 nodes): 20:50 hours
    - 20 CPUs, 4 GPUs (1 x com10 node): 4:49 hours



- Currently planning future hardware resources

# Cluster Scheduling Software : Job burst out

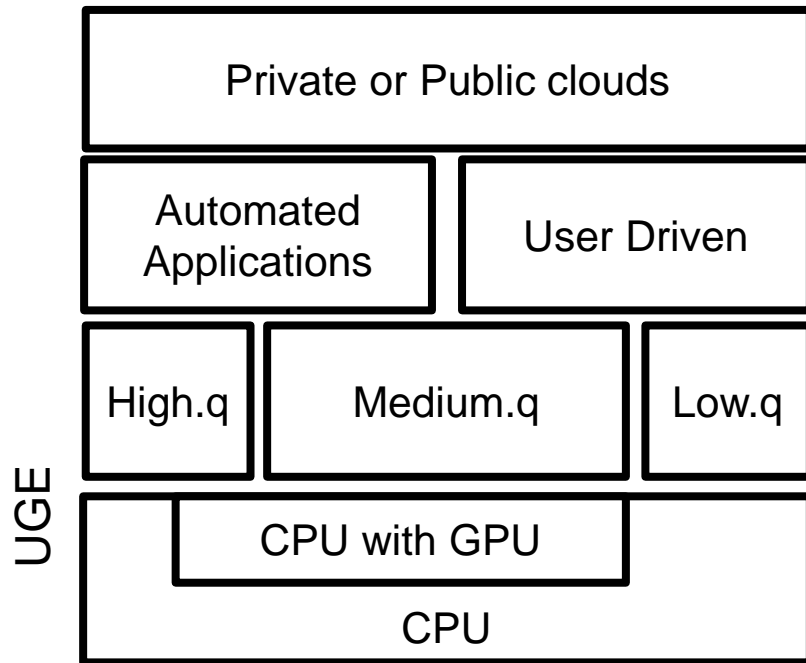
UGE scalability

- Fair play policy per project
- Quota resources

Burst out tests with Azure

- Manual triggering IMOD reconstruction (incl data transfer) gives comparable results to local cloud

Resource	Data transfer	Processing
Beamline Server		7m19.9s
Azure	1m24s (1.7GB)	4m56.9s
Cluster node		5m49.5s



# Cluster Scheduling Software : Job burst out

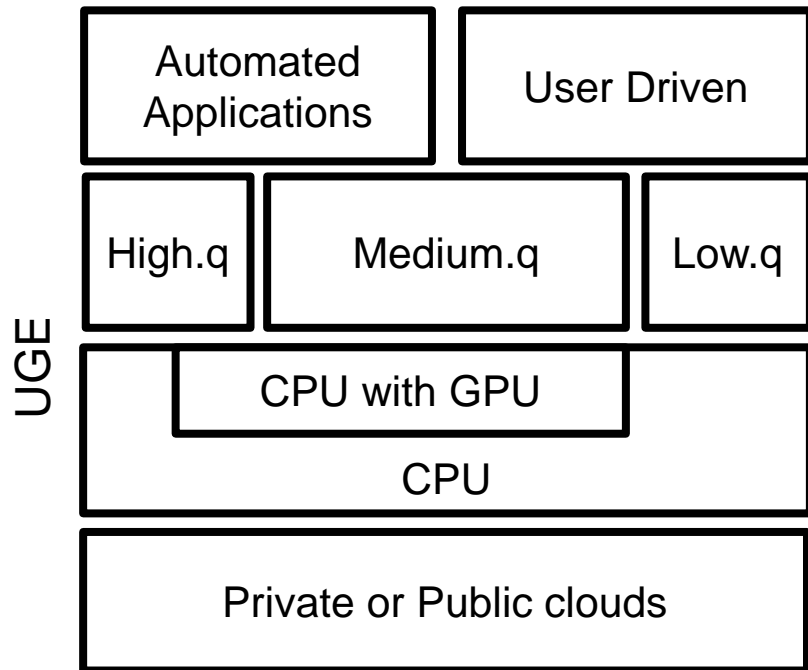
UGE scalability

- Fair play policy per project
- Quota resources

Burst out tests with Azure

- Manual triggering IMOD reconstruction (incl data transfer) gives comparable results to local cloud

Resource	Data transfer	Processing
Beamline Server		7m19.9s
Azure	1m24s (1.7GB)	4m56.9s
Cluster node		5m49.5s





1. Log onto SCD cloud using fedid & password

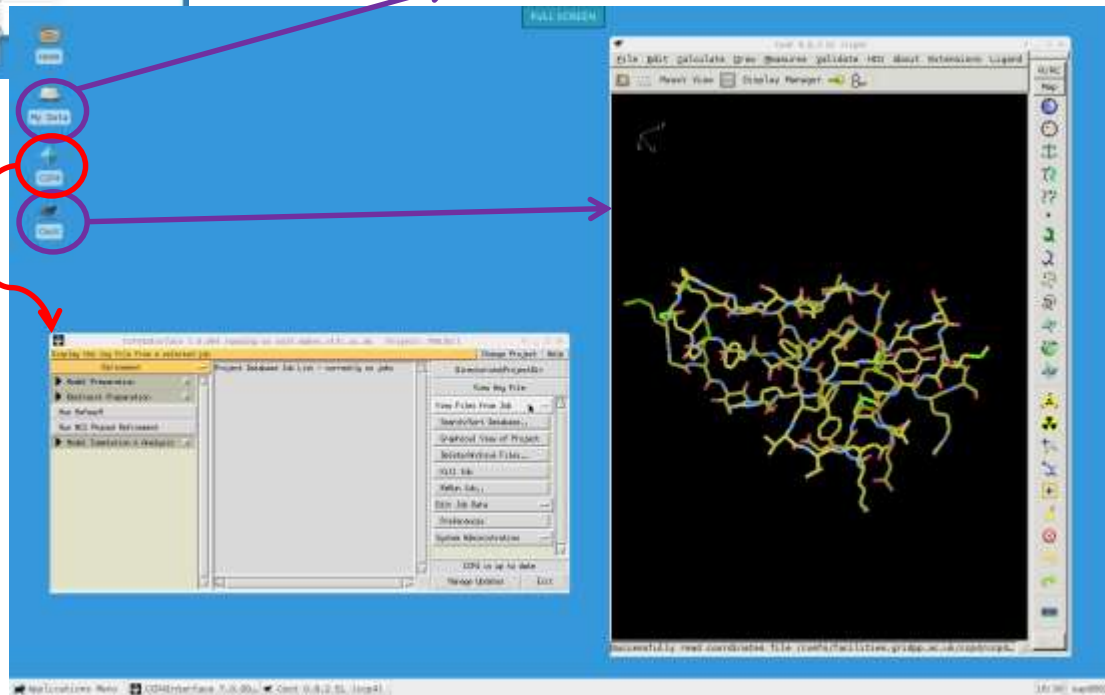


2. Select work environment



4. Browser view of the user's persistent storage (mock-up)

3. VM has latest CCP4 & useful software pre-installed, pulled from CVMFS repo



# Acknowledgements

- Karl Levik, Tina Fredrich, Ala Al-Afeef, Alun Ashton, Dave Stuart & @ Diamond
- Dan Clare, Alistair Siebert, Corey Hecksel, Peijun Zhang at eBIC
- SCIPION team: Jose Miguel de la Rosa Trevin, Jose Maria Carazo from IIPC, Madrid
- Juha Huiskonen from Strubi, Oxford
- CCPEM management and developers: Colin Palmer, Tom Burnley and Martyn Winn