

# CCP4 Cloud: Potential for a European-Wide Resource in Computational Crystallography

**Eugene Krissinel**

Research Complex at Harwell, Rutherford-Appleton Laboratory,  
Science & Technology Facilities Council, Harwell, UK

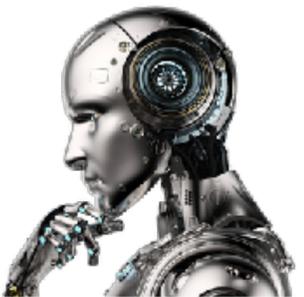
*eugene.krissinel@stfc.ac.uk*

# Science and Data Today

- ❖ In many fields, research becomes progressively data-driven
  - *taking on complex systems (e.g., climate, molecular biology, medicine)*
  - *puristic models not fitting real-world situations*
  - *imperfect metrics and rules for decision making*
  - *incomplete understanding and knowledge base*
  - *“I need to get things working first, and understand how they work second”*
- ❖ Combination of data, AI and computing power starts giving a break-through in many directions

*Recent example: predicting protein structures with AlphaFold 2 from DeepMind Technologies*

- *unprecedented accuracy in predicting protein structures in 3d*
- *expanding proteomics to genomics scales (from 180K known to 130M new structures)*
- *implications for structural biology research are breathtaking and still not fully understood*



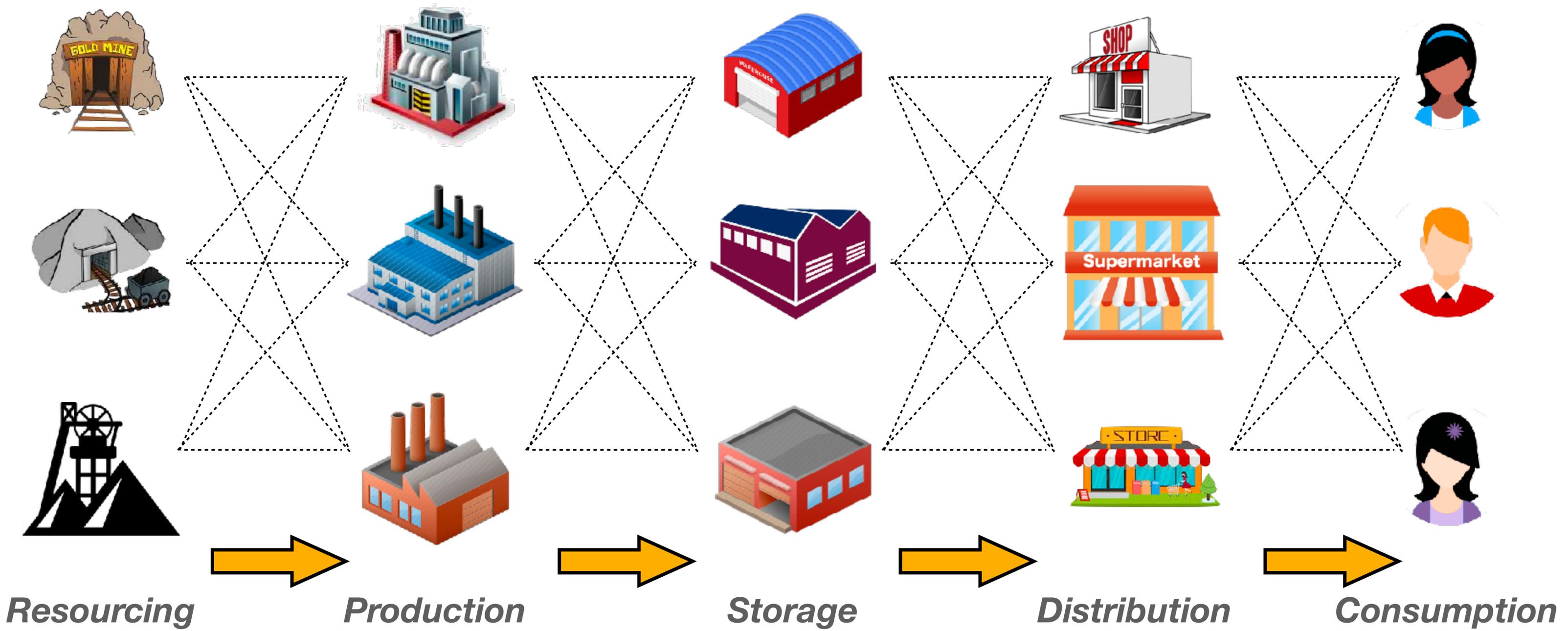
# Data and Computing Infrastructures



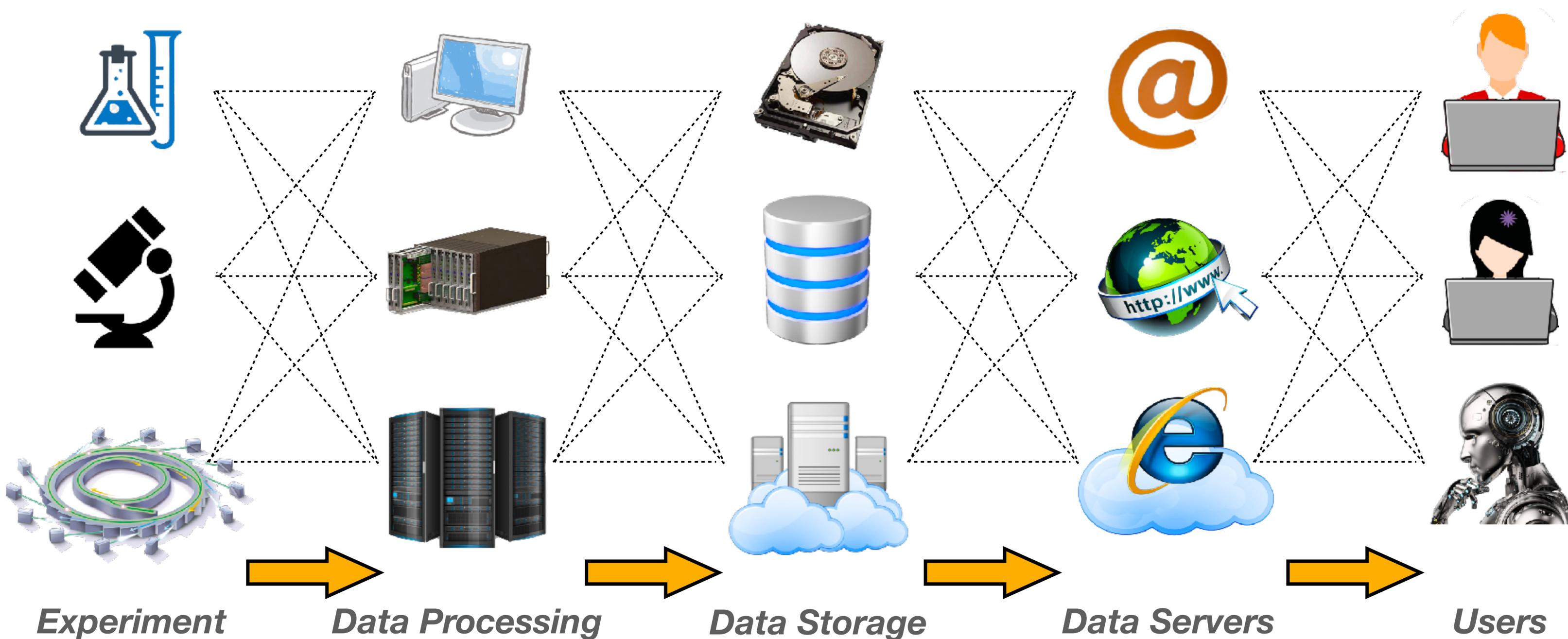
# EGI and Research Infrastructures



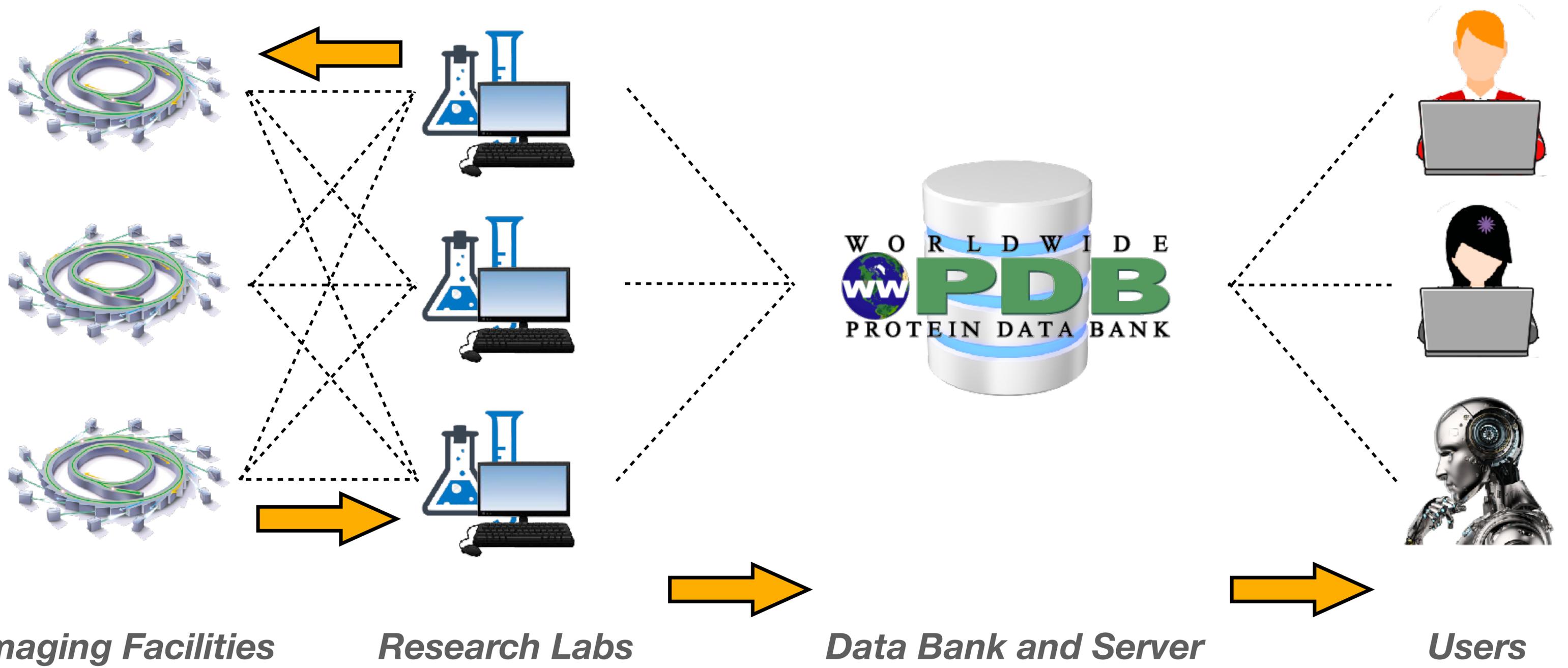
# Gross Idea: Distributed Production with Unified Access



# Gross Idea: Distributed Production with Unified Access



# Data Production in Structural Biology



# Data Maintenance in Structural Biology

- ❖ Single repository
  - Protein Data Bank established in 1971
  - Underpins academic and industrial research on global scale
- ❖ Does not keep all data
  - originally kept only atomic coordinates of biological macromolecules
  - reduced experimental data accepted from 1999 and made mandatory in 2008
  - processed experimental data accepted from 2020 but not mandatory
  - raw experimental data are not deposited
- ❖ No framework for keeping structure solution projects
  - valuable information about how the structure was solved, is not retained
  - structure solution cost varies between \$100K and \$2M
  - the total cost of re-creating the PDB, if lost, is estimated at over \$4B
  - life-time expectancy of data kept in local labs may be as short as few years



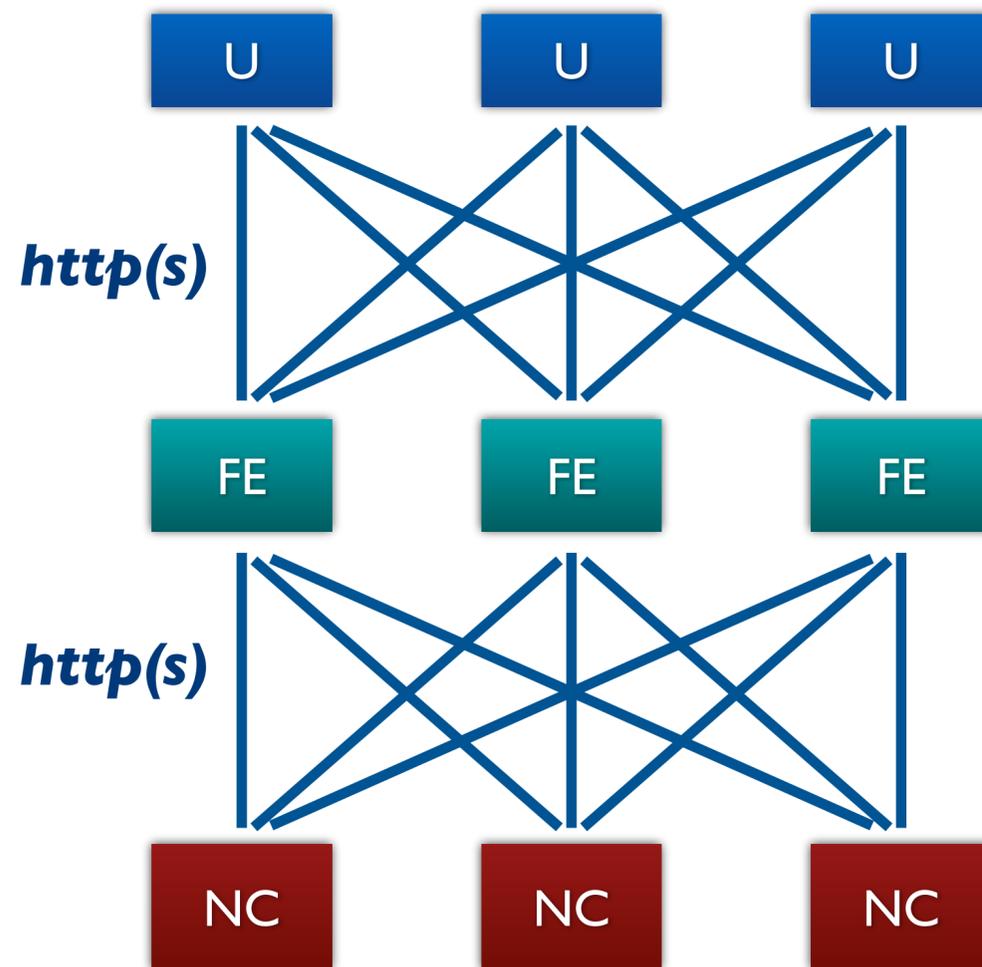
# CCP4 Cloud Initiative



- ❖ Conceived in 2016
  - *Funded by BBSRC UK and CCP4*
- ❖ Response to demands and trends rapidly emerging in the field
  - *CPU power (due to increased automation)*
  - *Centralised database support (due to expansion of methods based on data templates)*
  - *Software as a service (due to increased size and complexity of software setups)*
  - *Supporting distributed projects for team work*
  - *Cloud model for geographically-agnostic access and project data safety*
  - *Supporting personal mobile platforms (tablets and smartphones)*
  - *Communication with data facilities (synchrotrons, PDB, AFDB, etc)*
- ❖ Released in 2018
  - *Add-on in CCP4 7.0 (September 2018), fully integrated in CCP4 7.1 (April 2020)*
  - *Timely in pandemic situation*
  - *2,000 registered users*

# CCP4 Cloud Architecture

User Nodes  
(Integration)



Front End Nodes  
(Data and communication)



Number Cruncher Nodes  
(Computations)



❖ Not an application, but rather a configurable and expandable framework

- All nodes are functionally independent entities

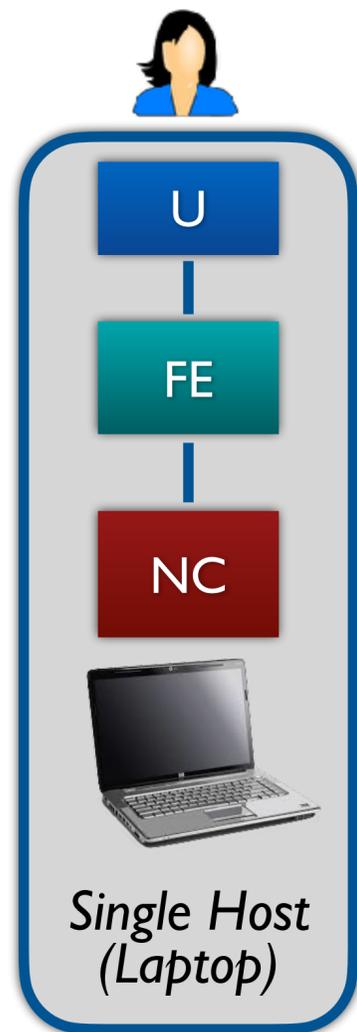
- Configuration done with light-weight JSON files

- Based on Node JS

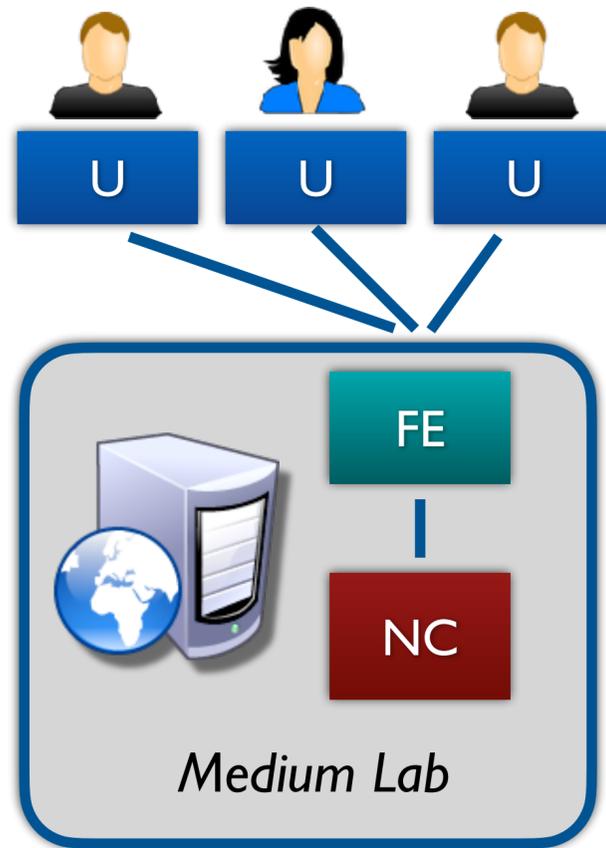
- Python adapters for executables

- No principal restrictions on the location and number of components

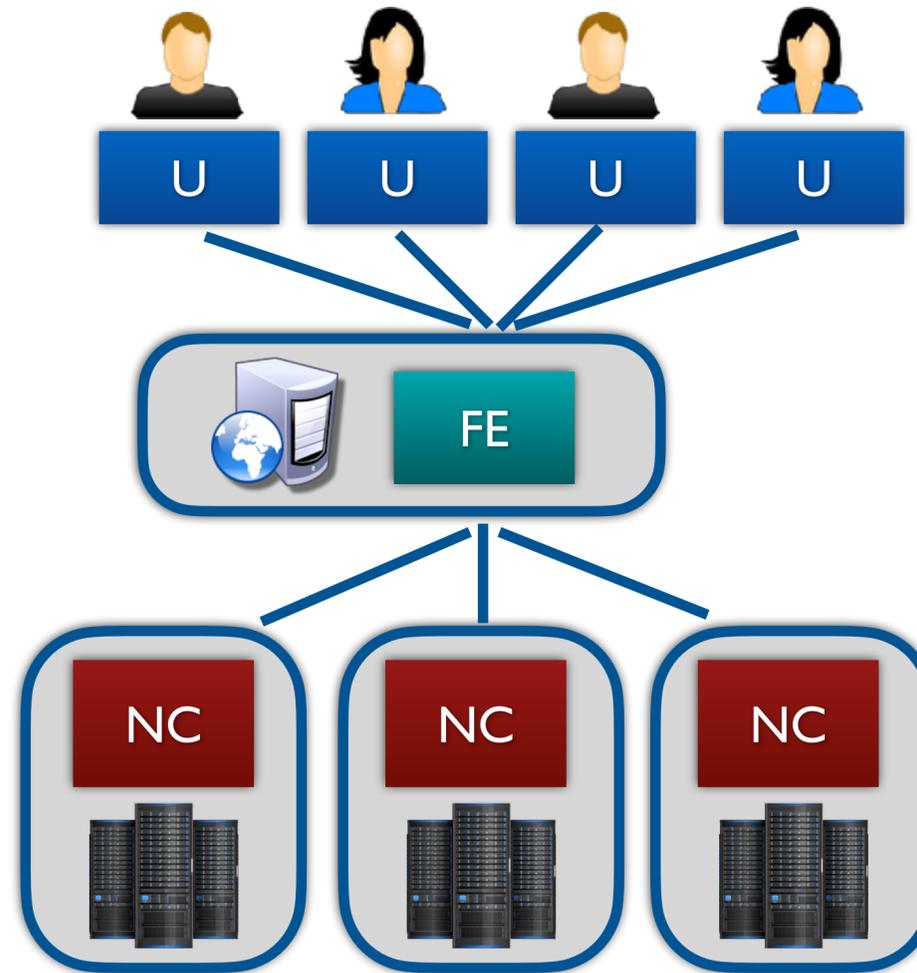
# Typical CCP4 Cloud Configurations



out-of-the-box



almost out-of-the-box



Major Hub (CCP4-Harwell)

❖ From CCP4 7.1, any CCP4 setup may be used as one or few CCP4 Cloud node(s)

- All nodes have identical codebase

- Single Host comes as a standard (effectively a desktop GUI)

- User Node comes pre-configured for working with CCP4-Harwell instance

# On User Side

The screenshot shows a web browser window with the address bar displaying 'localhost:51253'. The main content area features a large light blue arrow pointing right. In the center of the arrow is the 'CCP4 Cloud Login' page. The page title is 'CCP4 Cloud Login' in a large, bold, italicized font. Below the title is the 'CCP4-Harwell' logo and the URL 'https://cloud.ccp4.ac.uk/'. The login form consists of two input fields: 'Login name:' with the value 'eugene' and 'Password:' with masked characters '.....'. Below the form are three buttons: 'Login' (with a key icon), 'Forgotten password' (with a padlock icon), and 'Registration' (with a plus icon). Below the buttons, there is a note: 'Check CCP4 Cloud [roadmap](#) for new users' and a link to the '[Privacy Statement](#)'. At the bottom of the page, there is a lightbulb icon and the text: 'Are you using the Task Dock yet? Read [here](#) about it.' The footer contains the text 'Powered by CCP4 v.7.1.016' on the left, the 'CCP4 on-line' logo in the center, and the 'iris' logo on the right. The version and date 'CCP4 Cloud v.1.6.023 [17.08.2021]' are displayed in the bottom right corner.

*A Web-Application with rich graphical front-end*

# On User Side

Projects can be exported for exchange or archival

Work arranged in Projects

Projects can be shared in real time

ID	Name	R <sub>free</sub>	Disk (MBytes)	CPU (hours)	Date Created	Last Opened
05.staged-ep	Staged solution of Insulin structure with EP on sulphur atoms	0.1886	104.8	0.7749	2019-12-10	2021-10-13
[Afshan]:MBRWT	2.37_AUTOSELECT	0.1994	8956.6	14.093	2021-08-10	2021-10-13
[schenkan]:beta-lac	workshop	0.2471	869.6	1.9878	2021-08-09	2021-10-11
rbase	RNase	0.2139	103.8	177.7341	2021-05-13	2021-10-09
dna	DNA	0.2258	75.1	1.3758	2019-12-03	2021-10-08
[fabi.sm]:hGST	hGST	0.2863	4039.1	33.7089	2021-09-28	2021-10-05
MR-prac	Molecular Replacement Practical		12	0.0043	2021-08-20	2021-10-03
[Karthi]:006	ROQ_112	0.4583	529.1	6.3008	2021-09-27	2021-10-01
CatV-PCI3-1	CatV-PCI3 Rescued	0.3057	71.1	0.0381	2021-09-23	2021-09-24
hop-on	Hop-on example	0.2369	1480	0.7371	2021-02-20	2021-09-23
gdemo	Gamma demo to check parallel workflowing	0.2452	710.9	6.2086	2021-07-13	2021-09-22
[andrey_aps]:mdm2-2	mdm2	0.2915	296	1.4473	2021-07-02	2021-09-22
mdm2	MDM2	0.2883	368.5	63.3673	2021-01-13	2021-09-20
gere	GERE	0.3012	2224.7	1.7358	2020-01-14	2021-09-16
[nicholls]: HemeTest	HemeTest	0.3404	79.4	0.0761	2021-08-10	2021-09-13
bl			127.8	0.295	2021-07-28	2021-08-24
test	Installation test		711.6	5.5284	2021-08-16	2021-08-19
insulin	Insulin	0.2128	868.1	3.6383	2020-01-08	2021-08-16
ep-1	Experimental phasing	0.3089	107.4	0.7522	2020-02-24	2021-08-11

# On User Side

localhost:51253

Eugene Krissinel

### Staged solution of Insulin structure with EP on sulphur atoms

- [05.staged-ep] Staged solution of Insulin structure with EP on sulphur atoms
  - [0001] Based on CCP4 Insulin example. Open this remark for input data location
    - xia2 [0003] created datasets: **Unmerged (2) HKL (1) -- completed.** ← **Data import**
    - [0004] imported: **Sequence (2) -- completed.**
    - [0005] asymmetric unit contents -- *Solv=64.4%*
      - [0006] shelx substructure search (SAD) -- *R=0.5788 R<sub>free</sub>=0.5731*
        - [0008] low value of CC suggests that solution is unlikely; try another space group
        - [0007] change space group (ASU) -- *SpG=l 21 3*
        - [0009] shelx substructure search (SAD) -- *R=0.5643 R<sub>free</sub>=0.5922*
        - [0010] phaser EP (SAD) -- *completed.*
          - [0011] original hand branch
            - [0021] parrot DM -- *completed.*
            - [0027] buccaneer -- *Compl=100.0% R=0.2447 R<sub>free</sub>=0.2611*
            - [0037] structure has been built - correct hand
              - [0030] fit waters -- *N<sub>waters</sub>=66*
              - [0034] refmac5 -- *R=0.1527 R<sub>free</sub>=0.1886*
              - [0036] deposition -- *package prepared, pdb report obtained* ← **Sending results to the PDB**
              - [0038] needs more refinement
          - [0012] inverted hand branch
            - [0024] parrot DM -- *completed.*
            - [0025] buccaneer -- *Compl=0.0%*
            - [0026] structure cannot be built - wrong hand ← **Remarks for project annotation**

Projects develop as branching trees

Data flow automated  
Data types enforced

Branching

Sending results to the PDB

Remarks for project annotation

# On User Side

localhost:51253

[0034] refmac5 -- completed

Input Output

Report Main Log Service Log Errors

	Initial	Final
R factor	0.2278	0.1527
R free	0.2661	0.1886
Rms BondLength	0.0174	0.0196
Rms BondAngle	2.6902	2.3573
Rms ChirVolume	0.2161	0.2021

**Graph Data**

- ▼ Cycle 1. Rfactor analysis, F distribution v resln
  - ▶ Cycle 1. <Rfactor> v. resln
  - ▶ Cycle 1. <Fobs> and <Fc> v. resln
  - ▶ Cycle 1. % observed v. resln
- ▶ Cycle 1. FSC and Fom(<cos(DelPhi)>-acentric, centric, overall v resln
- ▶ Cycle 20. Rfactor analysis, F distribution v resln
- ▶ Cycle 20. FSC and Fom(<cos(DelPhi)>-acentric, centric, overall v resln
- ▶ Cycle 21. Rfactor analysis, F distribution v resln

Print

Resolution (Å)

Legend: Rf\_used, WR\_used, Rf\_free, WR\_free

Powered by CCP4 v.7.1.016 .2021]

*Rich graphical reports  
in-browser*

# On User Side

0034-01\_refmac Structure and electron density

Nothing here. Press H for help.

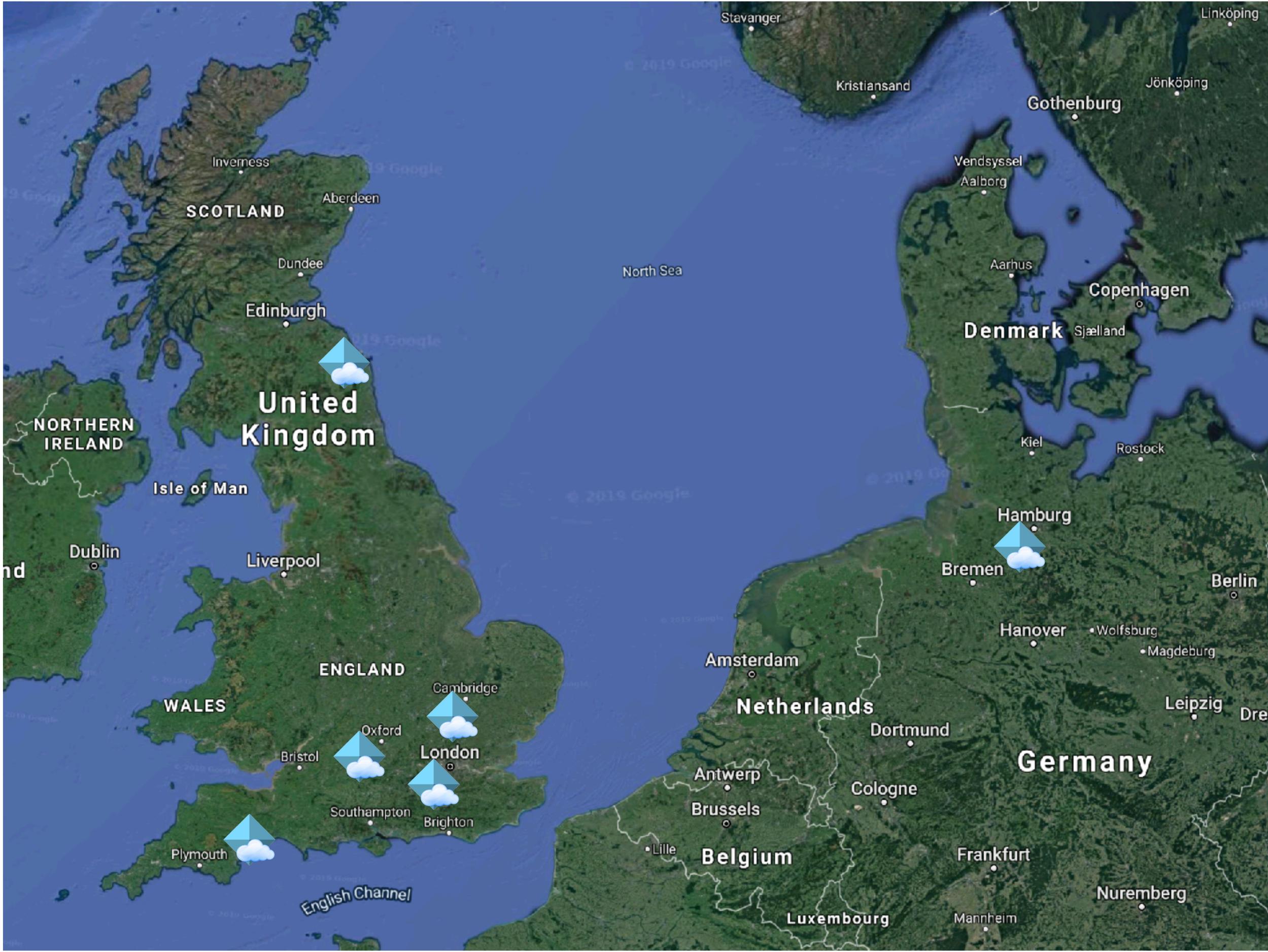
Embedded Molecular Graphics

Rf\_used  
WR\_used  
Rf\_free  
WR\_free

1.83 1.63

Powered by CCF

2021]



# Active CCP4 Cloud Instances

*Public (global)*

- ❖ CCP4-Harwell

*Institutional (local)*

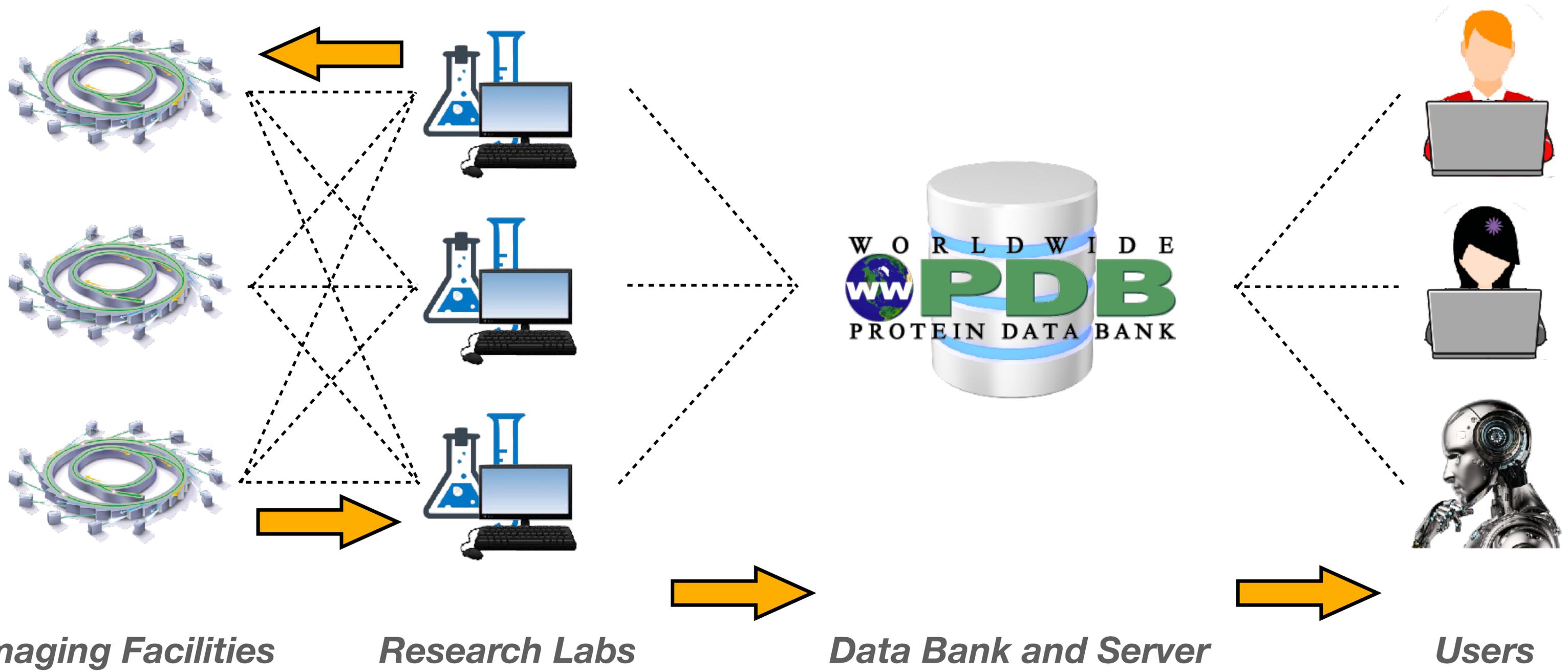
- ❖ Uni Exeter
- ❖ Francis Crick
- ❖ Uni Newcastle
- ❖ LMB Cambridge
- ❖ EMBL-Hamburg
- ❖ Incyte Inc. (USA)

# Main Features

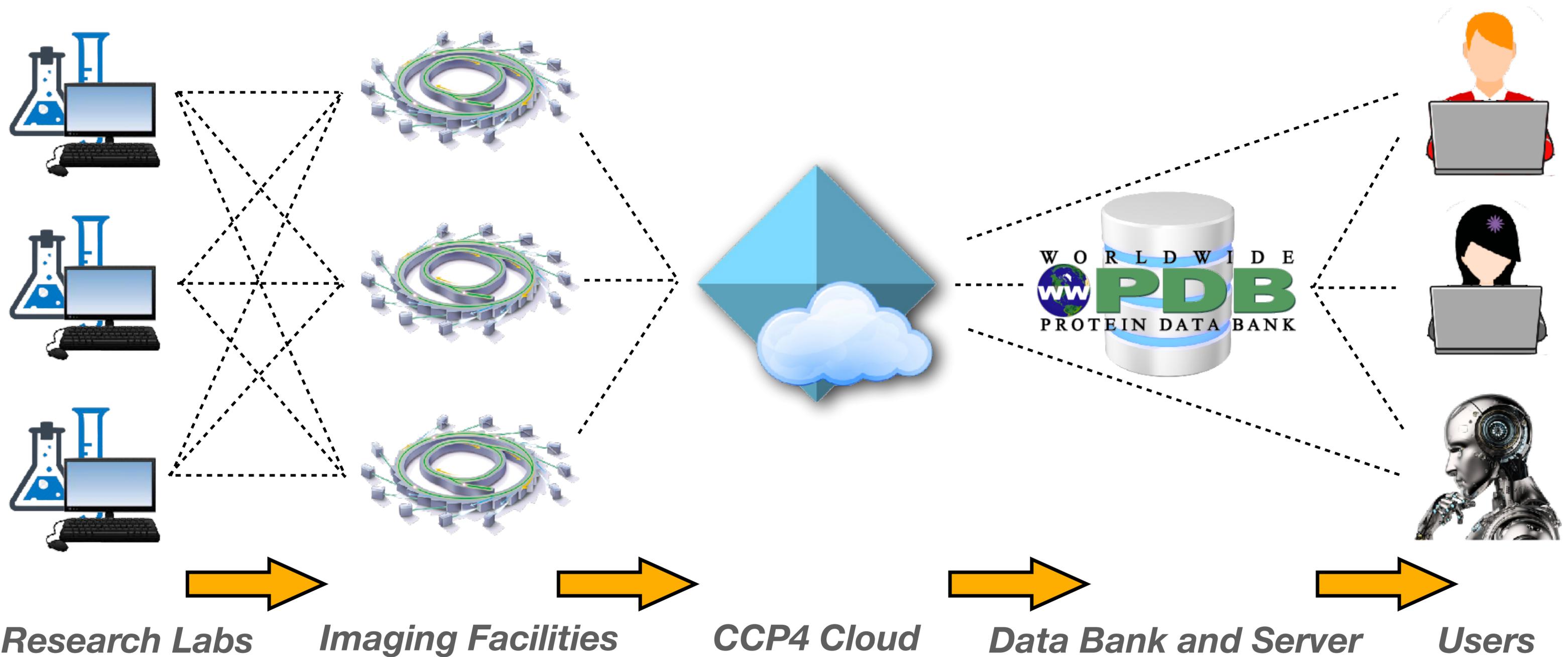
- ❖ Functionally complete
  - *includes all stages of structure solution*
- ❖ Designed to serve a number of users
  - *2,000 registered users*
  - *70,000 jobs/year at CCP4-Harwell*
- ❖ Designed to work with external data sources and services
  - *import from data facilities (currently iCAT, other in progress)*
  - *import from PDB and AFDB*
  - *integrated PDB deposition line*
- ❖ Designed to maintain structure solution projects
  - *backward compatible, versioned, metadating to ensure projects' integrity in future*
  - *all data is kept with projects*
  - *sharing projects in real time*



# Data Production in Structural Biology



# Data Production in Structural Biology



# Summary

- ❖ Data production in Science, in general, and Structure Biology, in particular, is taking industrial scales and approaches
  - *Data revolution, rise of AI*
  - *Automation of research and engineering*
  - *Approach to complex and fuzzy problems*
- ❖ Developing research infrastructures in form of centrally maintained data, software and computing services, linked to data producing and research centres, is a desirable way forward
  - *Efficiency through concentration and standardisation*
  - *Long-term data and software maintenance, knowledge keeping*
- ❖ Structural Biology infrastructure can be complemented with CCP4 Cloud
  - *Uniform approach to data management and structure solution across sites*
  - *Efficient data logistics*
  - *Retaining all data and metadata related to costly experiments*
- ❖ EGI is invited

# Acknowledgements

**CCP4, STFC & RCaH**

*Fantastic work environment, support and dissemination*



**CCP4 Collaboration,  
CCP4 School hosts and  
CCP4 developers**

*Contribution of task reports,  
general support and  
valuable feedback*

Research Complex  
at Harwell



**~300 test users**  
*Worldwide*

*Trial use and feedback on  
development versions of CCP4 Cloud*



**Andy Purkiss**

*Francis Crick Institute, London*

**Grzegorz Chojnowski**

*EMBL-Hamburg*

**Arnaud Basle**

*Newcastle University*

**Michael Isupov**

*University of Exeter*

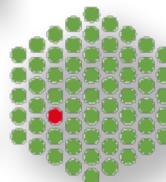
*Setup and  
maintenance  
of CCP4  
Cloud  
instances in  
their home  
labs*



Newcastle  
University

UNIVERSITY OF  
EXETER

EMBL



European Molecular  
Biology Laboratory

**Biotechnology and Biological  
Sciences Research Council  
(BBSRC) UK**

*Research grant  
BB/L0070317/1  
(2014-2019)*



## Developers and contributors

STFC, CCP4, Harwell, UK:

**Andrey Lebedev, Oleg Kovalevskiy,  
Charles Ballard, Ville Uski, Ronan Keegan**

RAS, Puschino, Russia:

**Maria Fando**

MRC/LMB, Cambridge, UK:

**Robert Nicholls**

EMBL-EBI, Hinxton, UK:

**John Berrisford**

Uni Leiden, The Netherlands:

**Navraj Pannu, Pavol Skubak**

Global Phasing Ltd, Cambridge, UK:

**Marcin Wojdyr, Clemens Vonrhein**

Uni York, UK:

**Stuart McNicholas**

Uni Liverpool, UK:

**Adam Simpkin, Jens Thomas**

Uni Birmingham, UK:

**Christopher Oliver**