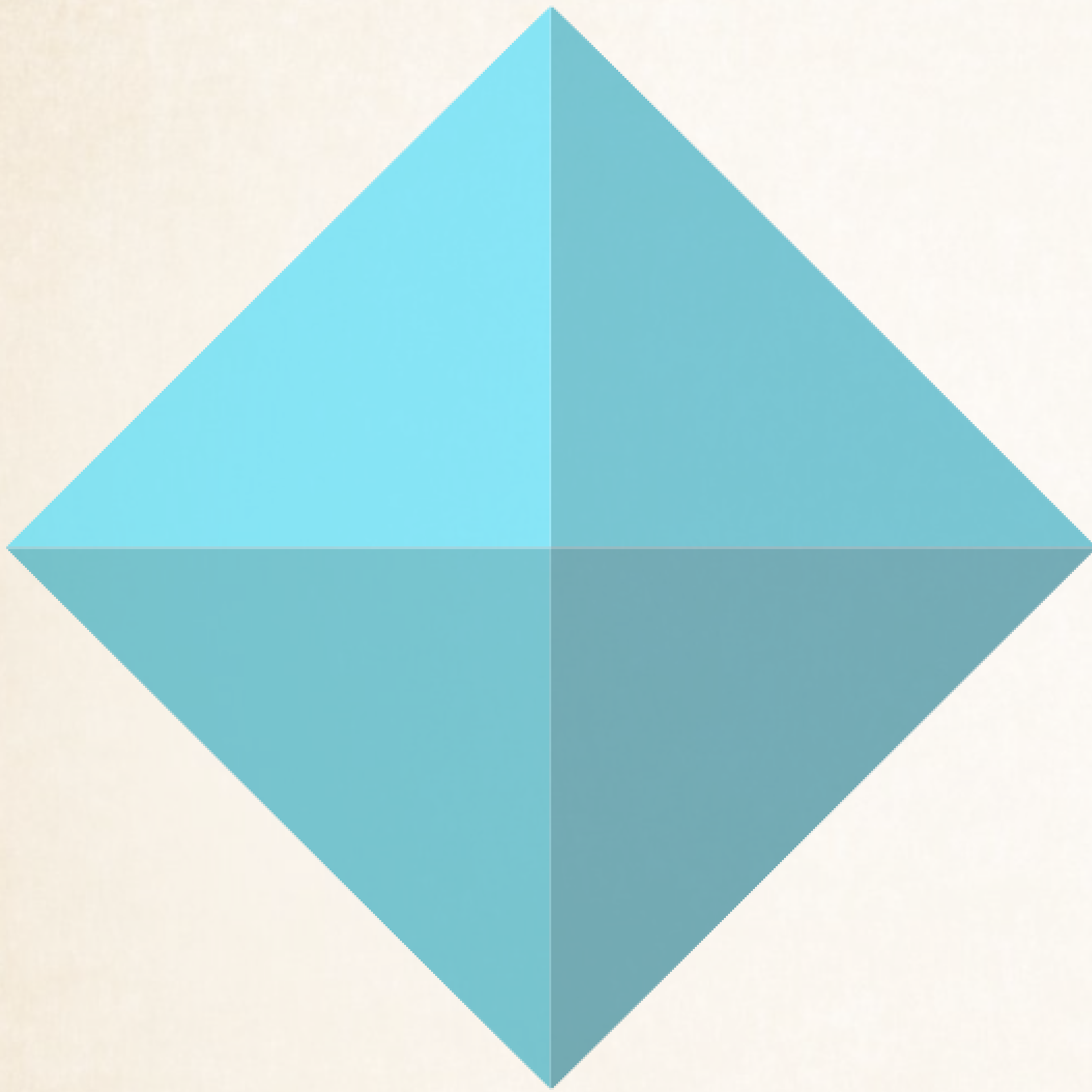# Raw Data Reuse: What Does this Mean for CCP4

*Eugene Krissinel*

*CCP4, Research Complex at Harwell, RAL, UK*

*eugene.krissinel@stfc.ac.uk*

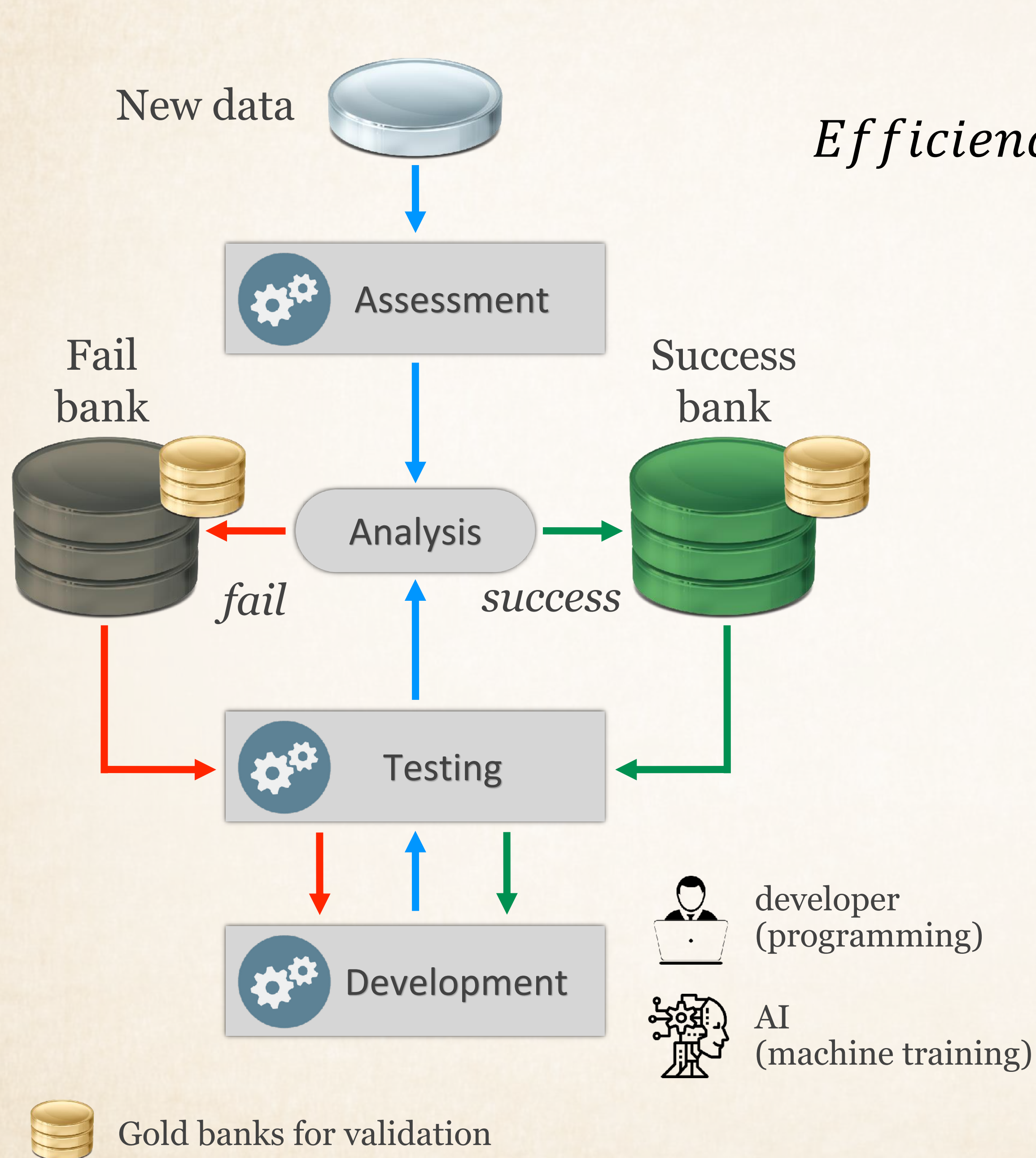*Data effect on Software **development***

- data reuse for software improvement
- data reuse for software testing
- data reuse for AI methods

*Data and Software **maintenance***

- data and software are in symbiosis
- data and software do age
- data and software have a cost to maintain

*Data and Software **legacy***

- data and software must be made available in publicly funded research
- data and software must be available for revisiting and revising results
- unlike software, data make scientific evidence

# Data and Software *Development*



bank mass

$$Efficiency(\text{⚙}) = \frac{M_{success}}{M_{fail} + M_{success}}$$

maximise by development

New data

Assessment

Fail bank

Success bank

Analysis

*fail*          *success*

Testing

Development

developer (programming)

AI (machine training)

Gold banks for validation
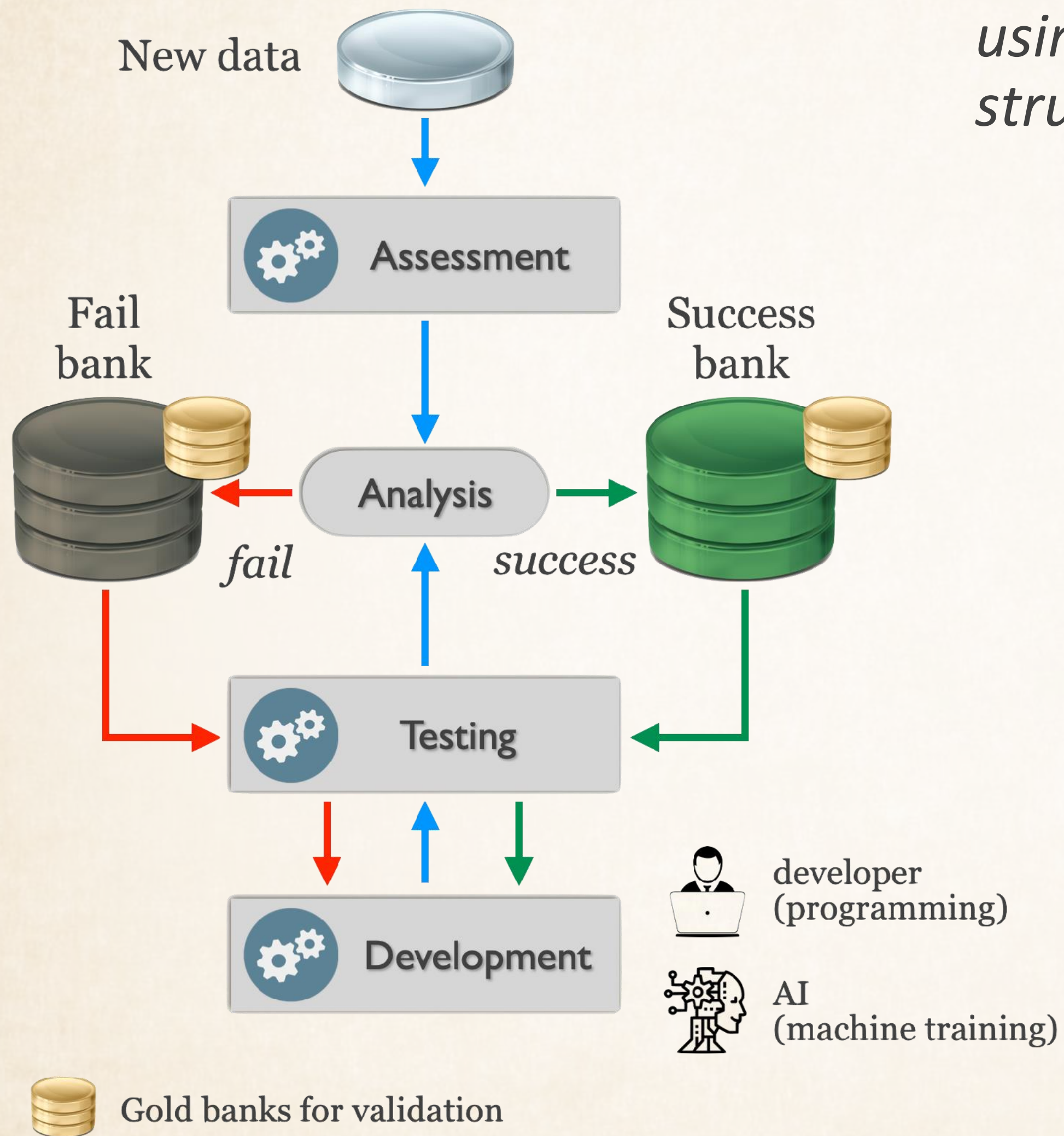
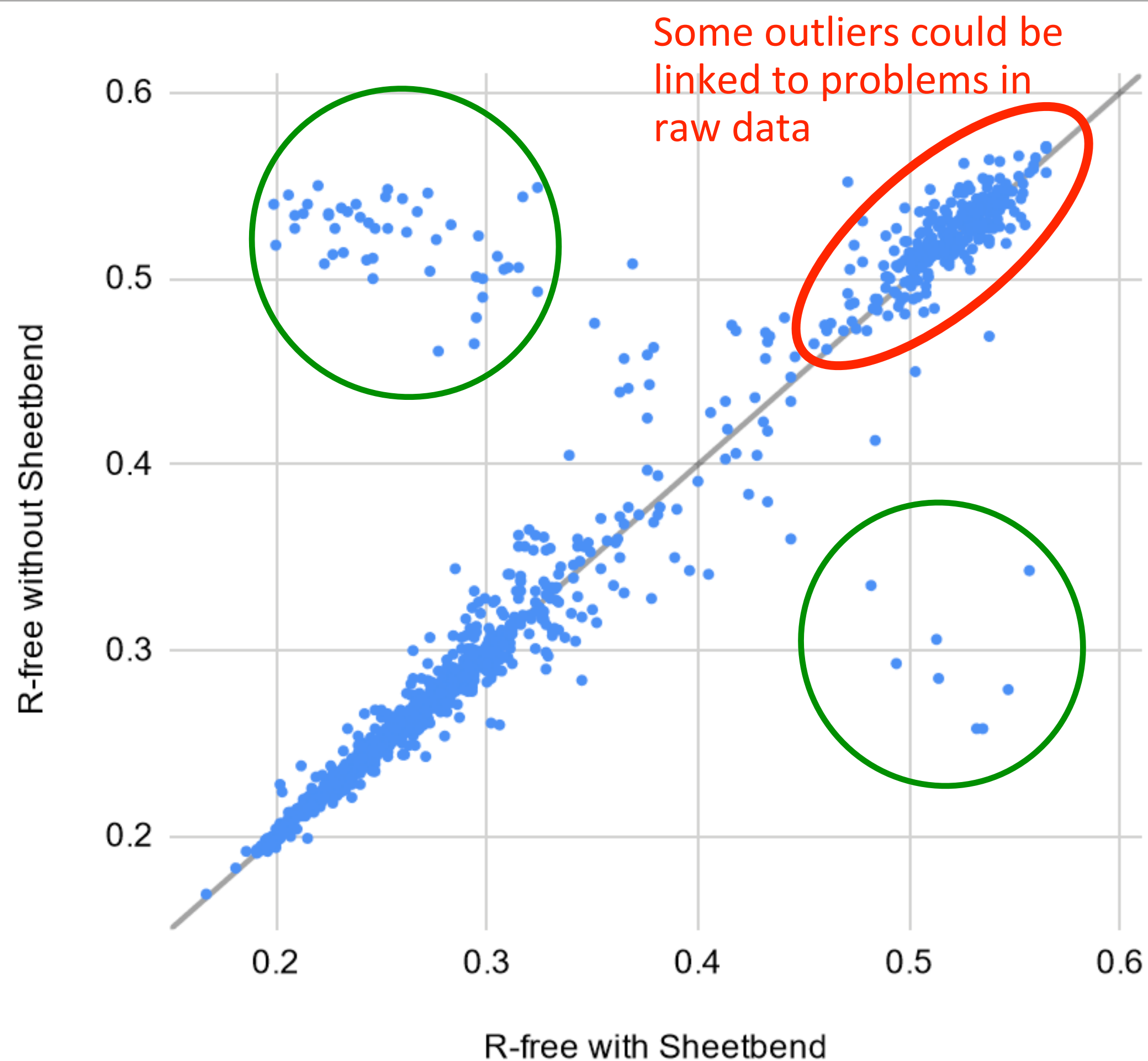*For efficient data (re)use*

- maintain **both** Success and Fail banks

- for Analysis, **annotate and classify** data; at least:
  - quality
  - statistical significance
  - redundancy

- assume annotation **changes** with time

- **maintain** annotation software

*Reusable data bank*  ≫  *pile of files*

# Data Reuse in CCP4

*In CCP4, data-driven development is based mostly on using processed (merged and unmerged) data and structure models from the PDB*

- simplified development loop is often used

- active developments include (but not limited to)

    - model building
    - MR, auto-MR, auto-EP
    - validation

- most of raw data reuse for development occurs in data processing project (DIALS)

    - may be hindered by the absence of comprehensive annotated raw data repository

# Data Reuse in CCP4



Some outliers could be linked to problems in raw data

The effect of shift field refinement on the performance of Modelcraft model building software. Some failures in red area could be due to data processing effects, not verifiable without access to raw data.

Picture provided by Paul Bond and Kevin Cowtan, University of York.

*Raw data processing is first suspect for problems in structure solution, model building and refinement*

- cannot be verified without access to raw data - benchmarking/training can go wrong

- access to raw data is essential for diagnosing data problems and development of corresponding tools

  - e.g., crystal pathologies; note that it is the **Fail bank** that is needed here

- the access must be programmatic, via **API** rather manual file fetching

- access to raw data is important for crystallography **education**

# Data and Software *Maintenance*

*Suppose raw data was archived from when CCP4 was founded (1979). Could it be used today? Not unless special care was taken.*
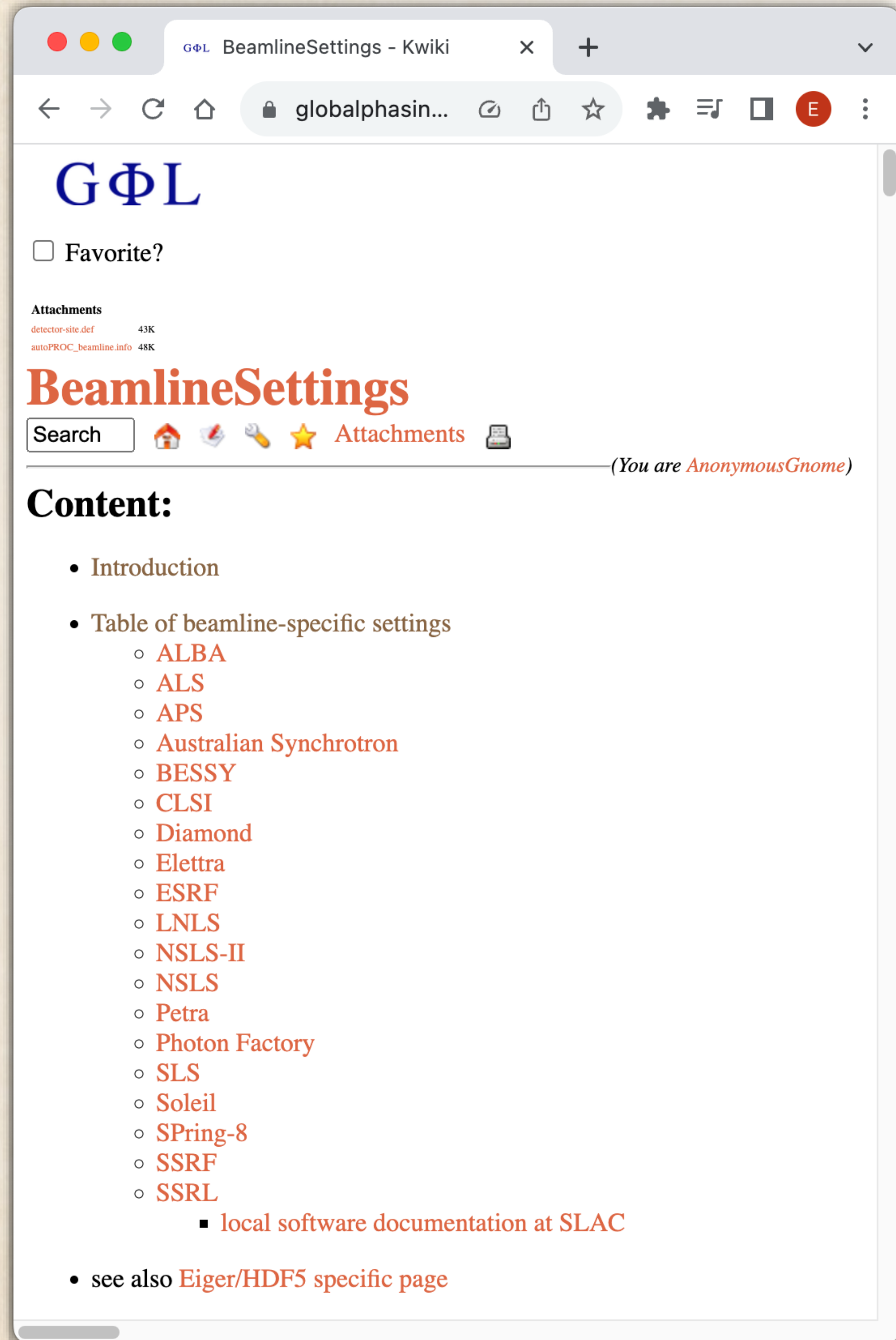
- evolution of detectors (effect on processing parameters)

  - photographic film (known from 1885)
  - charge-coupled devices (CCDs, known from 1970s)
  - pixel area detectors (known from 1980s, XRD use from 1990s)

- evolution of image formats (effect on software)

  CBF, MAR, TIFF, MAR345, CBF, SMV, RAXIS, CCP4, STOE, HARVARD, BRUKER, EDF ...

*Software should be kept backward-compatible; usually ages much faster than data*

- topic is at heart of CCP4 project; software maintenance proves to be a challenge

- programs do get retired for various reasons

| Package | Year | Status | Vendor/Distributor |
|---------|------|--------|--------------------|
| XDS | **1988** | active | Max-Plank |
| HKL | 1997 | active | HKL Research Inc. |
| DIALS | 2018 | active | Diamond/CCP4 |
| Mosflm | 1999 | sunset | CCP4 |
| d*Trek | 1999 | sunset | Rigaku |

# Data and Software *Maintenance*



*Having just image files is better than nothing but not enough!*

- instrument effects

  - damaged pixels and beamline settings vary from place to place, and may change in time.

  - essential for data re-use

  - pixel maps are included in some formats but not in all

  - Global Phasing Ltd. makes excellent job maintaining this information

*For efficient, long-term data re-use, consider:*

- deposition (proxy) data format
  - must be extendable; mind the future!

  - annotation framework
    - extendable for future
    - as automatic as possible (curators may be needed)

    - dedicated software support
      - data processing software
      - format adapters
      - annotation software
      - tests and integrity checks

      - central gateway
        - repositories may (and probably should) be distributed
        - access API

*It is a Big and Costly Project*
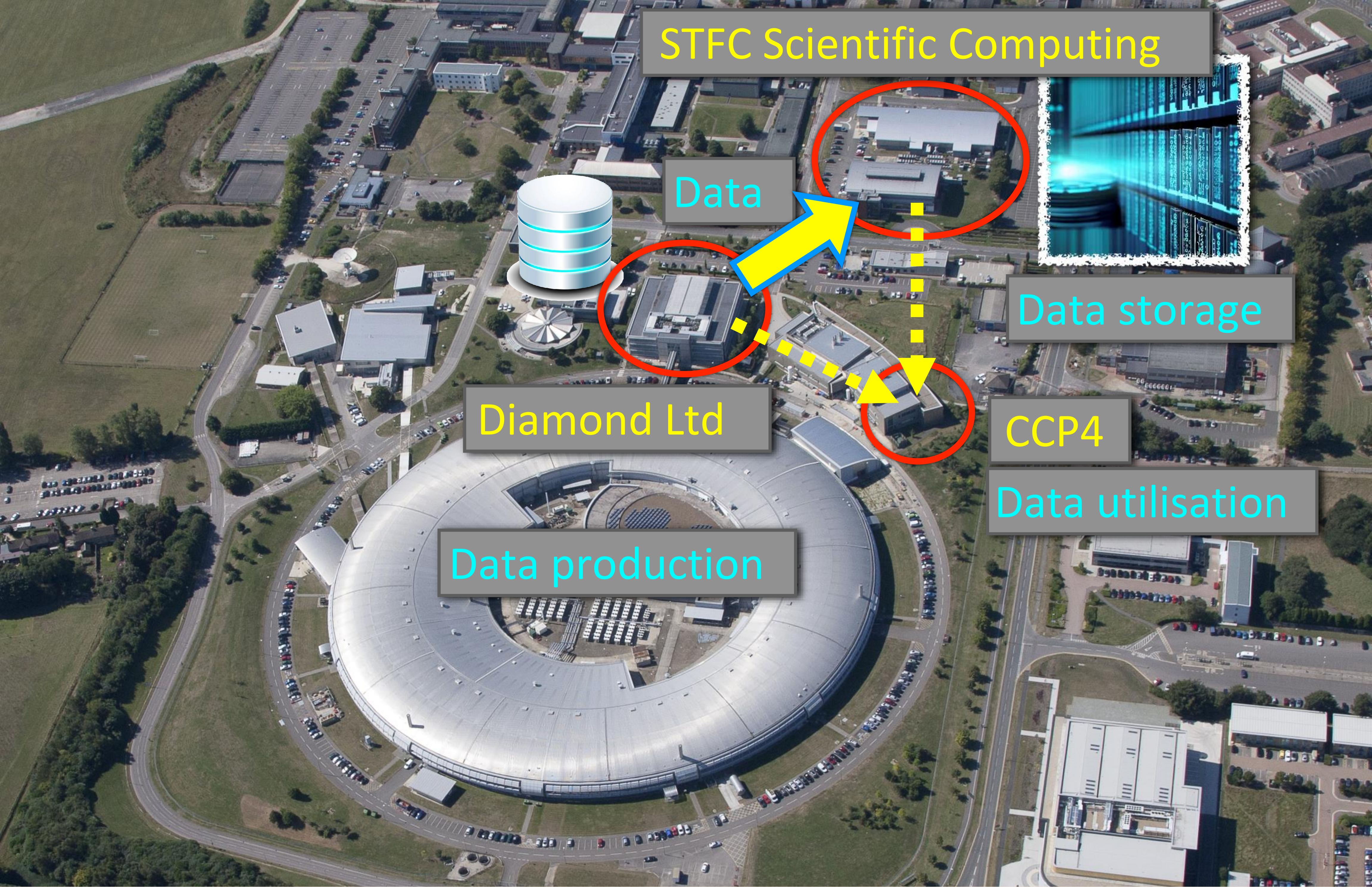
**buy cheap, buy twice**

STFC Scientific Computing

Data

Data storage

Diamond Ltd

CCP4

Data utilisation

Data production

# CCP4 and *Utilisation of Data*

Research Complex at Harwell

UKRI Science and Technology Facilities Council — Scientific Computing

*CCP4 Cloud for solving structures online*

## http://cloud.ccp4.ac.uk

- New way to solve **structures**
- New way to maintain **data**
- New way to maintain **projects**
- New way to use CCP4 **Software**
- Geographically agnostic

How To Overcome Fear of Change

CHANGE

**Structural Biology is a data-driven discipline: no data – no research**

**Public repository**

PDB Code: 1XYZ

Research

*structures*

~~*evidence*~~

# Current state of affairs: technique with limited reproducibility

**Experiment**

**Structure Solution**

**Public repository**

CCP4, Phenix, Global Phasing ...

*evidence*

*interpretation*

- *methods*
- *decisions*
- *assumptions*
- *doubts resolution*
- *validation*
- *alternatives*

*structures*

**PDB Code:** 1XYZ

As a rule, no or limited public access to raw data

**Reproducibility is limited**

As a rule, no public access to all details (publications and local computers)

**Reproducibility is limited**

**Research**

# *A better picture: raw data repository cross-linked with the PDB*

**Public**
**Experiment**   **repository**   **Structure Solution**

**Public**
**repository**

CCP4, Phenix,
Global Phasing ...

WORLDWIDE
wwPDB
PROTEIN DATA BANK

*evidence*   *raw data*   *interpretation*   *structures*

**Image DOI**

- *methods*
- *decisions*
- *assumptions*
- *doubts resolution*
- *validation*
- *alternatives*

**PDB Code:** 1XYZ

As a rule, no
public access to
all details

**Research**

# Data Trail in Structural Biology

Research Complex at Harwell

UKRI Science and Technology Facilities Council
Scientific Computing

*Yet better picture*

**Experiment** → **Public repository** → **Structure Solution** → **Public repository**

*evidence* → *raw data* → **CCP4, Phenix, Global Phasing ...** → *structures*

**Image DOI**

*interpretation*

**Public repository**

*project archive*

**Archive ID:** Project-1XYZ

**PDB Code:** 1XYZ

**Research**

# Data Trail in Structural Biology

Research Complex at Harwell

UKRI Science and Technology Facilities Council — Scientific Computing

*Possible picture in future*

**Experiment** — *evidence*

**Public repository** — *raw data* — **Image DOI**

**Structure Solution in CCP4 Cloud** — *interpretation*

*Solving Insulin structure with Molecular Replacement*

- [insulin-mr] Solving Insulin structure with Molecular Replacement
  - xia2 [0001] process diffraction images with Xia-2 / created datasets: **Unmerged (2**
  - [0002] import from cloud storage -- *imported: Sequence (2)*
  - *[0003] prepare MR search models for sequences A and B*
  - [0004] find MR models for sequence A with mrparse -- *24 MR model(*
  - [0005] find MR models for sequence B with mrparse -- *32 MR mo*
  - *[0006] define ASU content*
  - [0007] asymmetric unit contents -- *2 molecules in ASU, Solv*
  - *[0008] perform Molecular Replacement*
  - [0009] molrep -- R=0.5515 R$_{free}$=0.5552
  - [0010] refmac5 -- R=0.5244 R$_{free}$=0.5415

**Public repository** — *structures* — **WORLDWIDE PDB PROTEIN DATA BANK**

**PDB Code:** 1XYZ

*emerging, partly available*

**Public repository** (from 01/2023)

*CCP4 Cloud archive*

**Archive ID:** CCP4-1XYZ

**Research**

*If Raw Data were available as routinely as PDB models, CCP4 could:*

- do better testing of some components

- further enforce data-driven development practices

- deliver software of higher quality

- help improving data provenance in structural biology

*For Raw Data to be available as routinely as PDB models, CCP4 would need to:*

- maintain backward compatibility of data processing software

- maintain raw data format converters

- keep in sync with relevant metadata frameworks

- maintain data links with data producing facilities and data repositories

- maintain raw data access facilities in CCP4 front-ends for users

# Acknowledgements

## CCP4 Core, STFC, Harwell, UK:

Charles Ballard
Ronan Keegan
David Waterman
Andrey Lebedev
Ville Uski
Maria Fando
Tarik Drevon
David McDonagh
Jools Wills
Daniel Garza

## MRC/LMB, Cambridge, UK:

Garib Murshudov
Paul Emsley
Robert Nicholls

## University of Cambridge, UK:

Randy Read
Airlie McCoy
Robert Oeffner
Tristan Croll

## University of York, UK:

Keith Wilson
Kevin Cowtan
Jon Agirre
Paul Bond
Stuart McNicholas
Filomeno Sanchez

## EMBL-EBI, Hinxton, UK:

Sameer Velankar
John Berrisford
Deborah Harrus

## Leiden University, The Netherlands:

Navraj Pannu
Pavol Skubak

## Global Phasing Ltd, Cambridge, UK:

Gerard Bricogne
Clemens Vonrhein
Marcin Wojdyr

## University of Liverpool, UK:

Dan Rigden
Adam Simpkin
Jens Thomas

## Newcastle University, UK:

Martin Noble
Arnaud Basle

## University of Southampton, UK:

Ivo Tews

## University of Exeter, UK:

Michael Isupov

## EMBL-Hamburg, Germany

Victor Lamzin
Grzegorz Chojnowski

## Francis Crick Institute, UK:

Andrew Purkiss