

Winter School on Structural Biology, CEITEC, Brno, 13 February 2015

# Applied common sense

The why, what and how of validation

Gerard J. Kleywegt  
Protein Data Bank in Europe (pdbe.org)  
EMBL-EBI, Cambridge, UK

## What *is* validation?

## Validation according to the dictionary

- Validation = establishing or checking the truth or accuracy of (something)
  - Theory
  - Hypothesis
  - Model
  - Assertion, claim, statement
- Integral part of scientific activity!

“Science is a way of trying not to fool yourself. The first principle is that you must not fool yourself, and you are the easiest person to fool.” (Richard Feynman)

## Critical thinking

- Essential “24/7” skill for every scientist
  - And, in fact, for every non-scientist too
- Important aspect of validation

## Critical thinking

### Gun deaths in Florida

Number of murders committed using firearms

Source: Florida Department of Law Enforcement  
© Chen 26/02/2014



Unskewed Graphs Presents...

### SCIENTISTS EVENLY SPLIT ON CLIMATE CHANGE!

@Christians for Michele Bachmann, LLC

### Critical thinking

- What is wrong here?
  - The tacR gene regulates the human nervous system
  - The tacQ gene is similar to tacR but is found in *E. coli*
  - ==> The tacQ gene regulates the nervous system in *E. coli*!






And here?  
 "The tetramer has a total surface area of 81,616Å²"  
 (Implies: +/- 0.5Å² ...)

**PDBe** EMBL-EBI


### Validation = critical assessment

- How good is my model, really?
- At the very least:
  - Does it explain all the data that I used?
  - Does it explain all the prior knowledge that I had?
- More importantly:
  - Does my model explain all the data that I didn't use?
  - Does my model explain all the prior knowledge that I didn't use?
  - Is my model the best possible, most parsimonious explanation for the data?
  - Are the testable predictions based on my model correct?
- If any of these questions is answered with "no", you have a problem!

**PDBe** Occam's razor Popper's falsifiability principle EMBL-EBI





### The why of validation



Essentially, all models are wrong, but some are useful.  
 (George E. P. Box)

**PDBe** EMBL-EBI

### Crystallography is great!!

- Crystallography can provide important biological insight and understanding!!

**PDBe** And NMR, 3DEM, SAS etc. too, of course! EMBL-EBI

### Crystallography is great!!





- Crystallography can result in an all-expenses-paid trip to Stockholm (albeit in December)!!

**PDBe** EMBL-EBI

### Nightmare before Christmas

... but sometimes we get it horribly wrong



#### Retraction

WE WISH TO RETRACT OUR RESEARCH ARTICLE "STRUCTURE OF MtbA FROM *E. coli*: A homolog of the multidrug resistance ATP binding cassette (ABC) transporters" and both of our Reports "Structure of the ABC transporter MtbA in complex with ADP+vanadate and lipopoly-saccharide" and "X-ray structure of the EmrE multidrug transporter in complex with a substrate" (1-3).

The recently reported structure of Sm 1866 (4) indicated that our MtbA structures (1, 2, 3) were incorrect in both the hand of the structure and the topology. Thus, our biological interpretations based on these inverted models for MtbA are invalid.

An in-house data reduction program introduced a change in sign for anomalous differences. This program, which was not part of a conventional data processing package, converted the anomalous pairs (I+ and I-) to (I- and I+) thereby introducing a sign change. As the diffraction data collected for each set of MtbA crystals and for the EmrE crystals were processed with the same program, the structures reported in (1-3, 5, 6) had the wrong hand.

The error in the topology of the original MtbA structure was a consequence of the low resolution of the data as well as breaks in the electron density for the connecting loop regions. Unfortunately, the use of the multikey refinement procedure still allowed us to obtain reasonable refinement values for the wrong structures.

The Protein Data Bank (PDB) files 1J5Q, 1PF4, and 1ZZR for MtbA and 1S7B and 1F2M for EmrE have been moved to the archive of obsolete PDB entries. The MtbA and EmrE structures will be recalculated from the original data using the proper sign for the anomalous differences, and the new C<sub>α</sub> coordinates and structure factors will be deposited.

We very sincerely regret the confusion that these papers have caused and, in particular, subsequent research efforts that were unproductive as a result of our original findings.

GEORFFREY CHANG, CHRISTOPHER B. ROTH, CHRISTOPHER L. REYES, OWEN PORRILLIS, YEN-JU CHEN, ANDY P. CHEN

Department of Molecular Biology, The Scripps Research Institute, La Jolla, CA 92037, USA.

References


1. G. Chang, C. B. Roth, *Science* 298, 1763 (2002).
2. C. L. Reyes, G. Chang, *Science* 308, 3218 (2005).
3. G. Shenkin, Y. J. Chen, A. P. Chen, G. Chang, *Science* 318, 1910 (2005).
4. M. J. Dawson, K. J. Lesch, *Nature* 443, 180 (2006).
5. G. Chang, *J. Mol. Biol.* 306, 419 (2001).
6. C. Liu, G. Chang, *Proc. Natl. Acad. Sci. USA*, 103, 2852 (2006).

SCIENCE VOL 314 22 DECEMBER 2006 1875

**PDBe** EMBL-EBI

### Why do errors survive?

- "Why do errors make it into the literature and the PDB?"
- Suggestions from students
  - Cold Spring Harbor course, 2005
  - Copenhagen University course, 2006



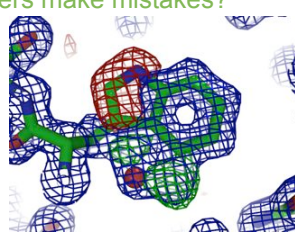
Who/what do YOU think is to blame?

### Playing the Blame Game ...

- Crystallographer
  - ignorance, lack of experience, incompetence, incorrect preconceptions/bias, cheating, laziness, "science by mouse-click", stress, can't be bothered to fix minor problems, no validation
- PI
  - pressure to publish/graduate fast, career interest, competition/scoops, grant writing, insufficient supervision
- Referees/editors
  - lazy, inadequate reviewing routines, no access to raw data, "validation by senior author name", lack of experience
- Software
  - misses or causes errors
- PDB
  - doesn't check
- "Nature"
  - limitations of the technique/resolution, errors hard to detect, poor data


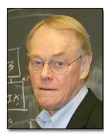
### Why do crystallographers make mistakes?

- Limitations to the data
  - Incomplete
  - Weak
  - Limited resolution
  - Space and time averaged
  - Phase errors
- The human factor
  - Subjectivity and bias involved in map interpretation and refinement (even at atomic resolution!)
  - Inexperienced people do the work, use of black boxes, ...
  - Not everybody is a good chemist
  - Even experienced people make mistakes



Kleywegt, *Acta Cryst. D65*, 134 (2009)

### Crystallographer = Super(wo)man?


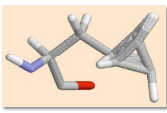




- The crystallographer ideally has
  - Knowledge of the history of the sample
  - Knowledge of the biology of the system
  - Knowledge of chemistry
  - Knowledge of physics
  - Understanding of data collection and processing
  - Understanding of the refinement process and software
  - Experience in map interpretation (preferably with a range of resolutions, space groups, etc.)
  - Read and remembered all the relevant literature
  - ...

(Wayne Hendrickson)


### The odds are stacked against us

- Crystallographers produce models of structures that will contain errors
  - High resolution AND skilled crystallographer → probably nothing major
  - High resolution XOR skilled crystallographer → possibly nothing major
  - NOT (High resolution OR skilled crystallographer) → pray for nothing major

"I know the human being and fish can coexist peacefully"

### A little experiment



- Hypothesis: "If a card has a vowel on one side, then it has an even number on the other side"
- Validate this hypothesis by turning as few cards as possible
- How many, and which, cards must you turn?

Wason selection task

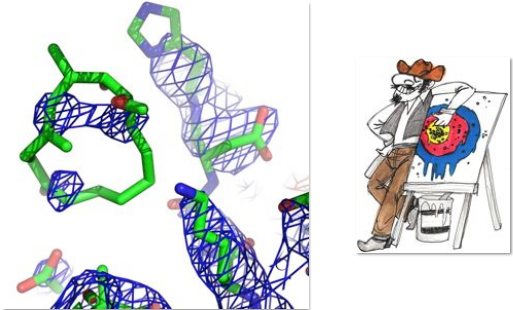
### Confirmation bias

- A scientific model is a hypothesis to be shot down
- We should be looking for **disconfirming** evidence
- But we often don't! We tend to look for **supporting** evidence
  - Reasonable expectation to find a ligand + Any old density blob in a reasonable ligand-binding site => Model the ligand!
  - Even if it isn't really there...
  - Conversely: we don't expect a ligand, so we model waters

PDBe

EMBL-EBI

### "Believing is seeing..."



PDBe

Retracted "ligand complex" published in *Nature*

EMBL-EBI

"A philosopher is a blind man in a dark room looking for a black cat that isn't there"

"A crystallographer is the man who finds it"

PDBe

Paraphrasing HL Mencken

EMBL-EBI

### Xtallography ≠ exact science

- Crystallographic models will contain errors
  - Crystallographers need to fix errors (if possible)
  - Users need to be aware of potentially problematic aspects of the model
  - Note: every crystallographer is also a user!
- Validation is important
  - Is the model as a whole reliable?
  - How about the bits that are of particular interest?
    - Active-site residues
    - Interface residues
    - Ligand, inhibitor, co-factor, ...

PDBe

EMBL-EBI

### Why don't people admit to their errors easily?

- To err is human
  - But so is denying that you erred
  - In some cases, "retraction battles" have raged for years
- Cognitive dissonance - discomfort caused by conflicting views of self
  - "I am an intelligent, hard-working scientist who makes good decisions"
  - "There is an error in my structure"
- How to resolve this discomfort?



PDBe

EMBL-EBI

### Cognitive dissonance – ways of coping

- (1) Self-justification/denial/passing the buck
  - "There's nothing wrong with it"
  - "It doesn't change the conclusions"
  - "Everybody makes those kinds of errors"
  - "It's really a matter of interpretation"
  - "It's probably low occupancy/high mobility"
  - "There is strain in the active site"
  - "It fits other data/my chemical intuition"
  - "It was my student's first structure"
  - "Legacy software changed the signs of  $\Delta F_{\text{atom}}$ "
- (2) Depression – no need for that!
- (3) Acceptance/reconciliation – the grown-up thing to do
  - "I made an error, I'll fix it and learn from it"
  - Still an intelligent, hard-working scientist!
  - Doing yourself and science a favour



Proceedings of the CCP4 Study Weekend. Accuracy and Reliability of Macromolecular Crystal Structures (1990)

PDBe

(David Eisenberg)

EMBL-EBI



## Cognitive dissonance in action

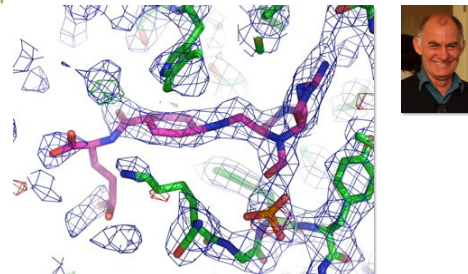
THE LIGAND N5G IN THIS ENTRY IS N5-IMINIUM PHOSPHATE. HOWEVER, THERE IS SOME DISCREPANCY IN THE GEOMETRY. THE GEOMETRY FOR N5G IS SUGGESTED BY THE REFINEMENT. THE CO-ORDINATES FIT WELL IN THE ELECTRON DENSITY MAP. THE MAP WAS GENERATED USING A DATASET COLLECTED AT 2.8 ANGSTROM RESOLUTION. THE DENSITY FOR THE LIGAND IS UNAMBIGUOUS AND THEREFORE THE GEOMETRIES ARE CORRECT AND ARE AS THEY WOULD BE IN A BIOLOGICAL MOLECULE, WHERE THE MICRO ENVIRONMENT HAS A PROFOUND INFLUENCE ON THE GEOMETRIES OF THE LIGAND.

- Single N-C bonds of 1.1 and 1.6Å
- Non-bonded C...C contact of 2.0Å
- PO<sub>3</sub> moiety separated by 2.7Å from O

EMBL-EBI

EMBL-EBI

## The experimental "evidence"



"Evidence that molecular-orbital theory breaks down in the presence of a protein crystallographer" (K. Henrick)

EMBL-EBI

EMBL-EBI



Mol	Type	Chain	Res	Link	Bond lengths			Bond angles		
					Counts	RMSZ	# Z  > 2	Counts	RMSZ	# Z  > 2
4	N5G	A	501	-	36,36,40	5.86	20 (55%)	49,50,57	8.25	27 (55%)
5	PO4	A	502	-	1,3,4	4.54	1 (100%)	0,3,6	0.00	-

EMBL-EBI

pdbe.org/valrep/3hy4

EMBL-EBI

## Errors and validation

- We need to take the drama out of the whole issue of errors and validation
- "When a friend makes a mistake, the friend remains a friend and the mistake remains a mistake" (S. Peres)
- Lao Tzu (more than 2500 years ago):  
A great nation is like a great man:  
When he makes a mistake, he realises it  
Having realised it, he admits it  
Having admitted it, he corrects it  
He considers those who point out his faults as his most benevolent teachers.



EMBL-EBI

EMBL-EBI

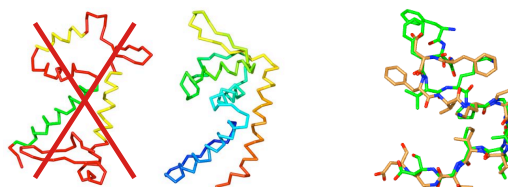
What kinds of errors do crystallographers make?

EMBL-EBI

EMBL-EBI

## Errors in protein structures

- Brändén & Jones (1990)
  - Mistracing an entire molecule or domain
  - Register errors
  - Local errors in the main chain



EMBL-EBI

Kleywegt, *Acta Cryst. D56*, 249 (2000)

EMBL-EBI

### Example of a tracing error

1PHY (1989, 2.4Å, *PNAS*)  
Entire molecule traced incorrectly

2PHY (1995, 1.4Å)

**PDBe** EMBL-EBI

### Example of a tracing error

1FZN (2000, 2.55Å, *Nature*)  
- One helix in register, two helices in place, rest wrong  
- 1FZN obsolete, but complex with DNA still in PDB (1FZP)

2FRH (2006, 2.6Å)

**PDBe** EMBL-EBI

### What are register errors?

- For a segment of a model, the assigned sequence is out-of-register with the actual density

**PDBe** EMBL-EBI

### Example of a register error

- 1CHR (light; 3.0Å, 1994, *Acta D*) vs. 2CHR (dark)

**PDBe** EMBL-EBI

### Example of a register error

1ZEN (green carbons), 1996, 2.5Å, *Structure*

1B57 (gold carbons), 1999, 2.0Å

```

1B57 (A)  ---SKIFDFVKPGVITGDDDVQKVFQ
.=ALIGN  | =ID  .. .....| | | | | | |
1ZEN (.)  SKI-PD-FVKPGVITGD-DVQKVFQ
    
```

Confirmed by iterative build-omit maps  
(Tom Terwilliger *et al.*, 2008)

**PDBe** EMBL-EBI

### Problems with ligands

**PDBe** EMBL-EBI

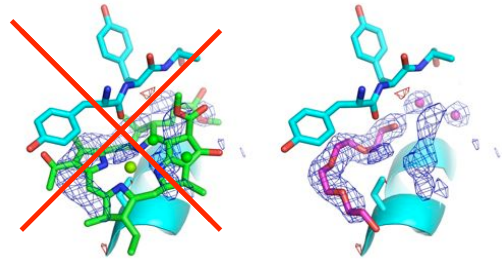
### Reasonable assumptions?

- Typical assumptions
  - We know what the ligand is
  - The modelled ligand was really there
  - We didn't miss anything important
  - The observed conformation is reliable
  - At high resolution we get all the answers
  - The H-bonding network is known
  - We can trust the waters
  - We are good chemists
  - (The complex structure is relevant for drug design)

EMBL-PDBe

EMBL-EBI

### A case of mistaken identity...



3OEG – bacteriochlorophyll-a

3VDI – PEG fragment and waters

EMBL-PDBe

Tronrud & Allen, *Photosynth. Res.* **112**, 71 (2012)

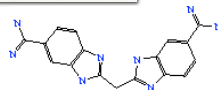
EMBL-EBI

### The ligand is really there?

**Structural Basis for BABIM Inhibition of Botulinum Neurotoxin Type B Protease** [*J. Am. Chem. Soc.* **2000**, *122*, 11268–11269]. Michael A. Hanson, Thorsten K. Oost, Chanokporn Sukonpan, Daniel H. Rich, Raymond C. Stevens\*

Page 11268. After a detailed analysis of the electron density maps for the structure of the inhibitor complex, we have concluded that the maps do not support the placement of the inhibitor as stated in the paper. Therefore, we are withdrawing the structural conclusions derived from PDB file 1FQH presented in the paper.

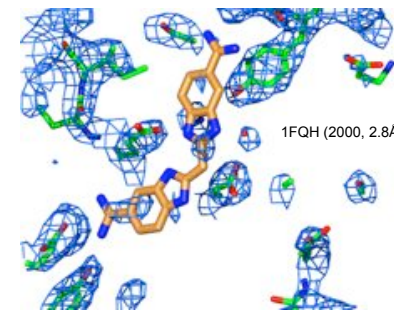
(*J. Amer. Chem. Soc.*, August 2002)



EMBL-PDBe

EMBL-EBI

### Dude, where's my density?

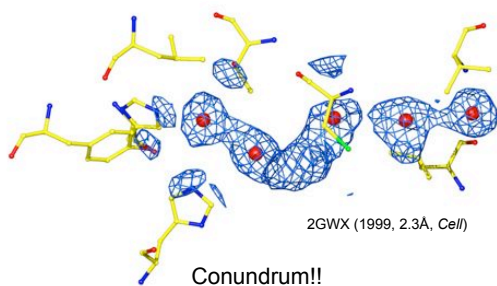


1FQH (2000, 2.8Å, JACS)

EMBL-PDBe

EMBL-EBI

### We didn't miss anything?



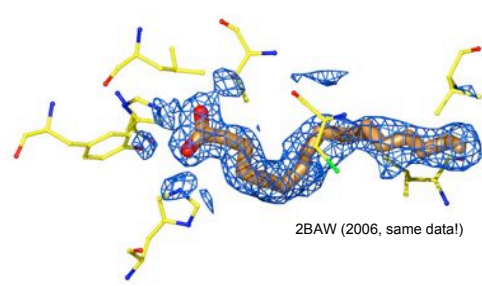
2GWX (1999, 2.3Å, Cell)

Conundrum!!

EMBL-PDBe

EMBL-EBI

### Oh, *that* ligand!





2BAW (2006, same data!)

EMBL-PDBe

EMBL-EBI

### Small-molecule anomalies

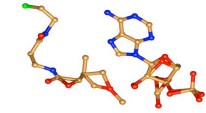
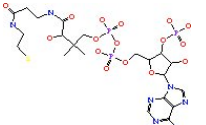

- 3-Phenylpropylamine in 1TNK, 1994, 1.8Å, *Nature Struct. Biol.*
- Aromatic carbon in between planar (0°) and pyramidal (35°) ... 17°

**PDBe** EMBL-EBI

### Oops-a-daisy!

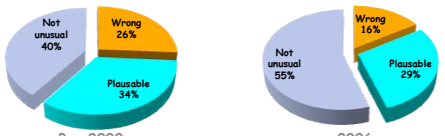
- COA = coenzyme A
- 2.25Å, R 0.25/0.28, *Mol. Cell*
- Deposited 2003
- Non-bonded contacts as close as 0.54Å
- Bond lengths up to 6.7Å
- Bond angles as low as 18°
- Impropers of 160°

**PDBe** EMBL-EBI

### Validation of PDB ligand structures by CCDC

- 16% of PDB entries deposited in 2006 had ligand geometries that were almost certainly in significant error (*in-house analysis using Relibase+/Mogul*)
- The good news - for structures before 2000 the figure was 26%



Time Period	Not unusual	Wrong	Plausible
Pre 2000	40%	26%	34%
2006	55%	16%	29%

(Jana Hennemann & John Liebeschuetz)

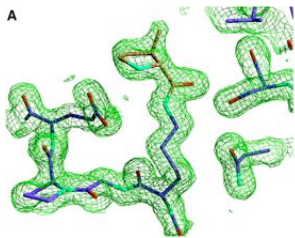
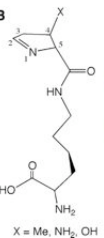
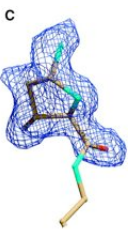
**PDBe** Liebeschuetz et al., *J. Comput. Aid. Mol. Des.* 26, 169 (2012) EMBL-EBI

### High resolution reveals all?

- Even at very high resolution there are sources of subjectivity and ambiguity
  - How to model temperature factors?
  - Is a blob of density a water or not?
  - How to model alternative conformations?
  - How to interpret density of unknown entities?
  - How to tell C/N/O apart?

**PDBe** EMBL-EBI

### The 22<sup>nd</sup> amino acid @ 1.55Å

Sodium chloride  
(Hao et al., 2002; PDB entries 1L2Q and 1L2R)

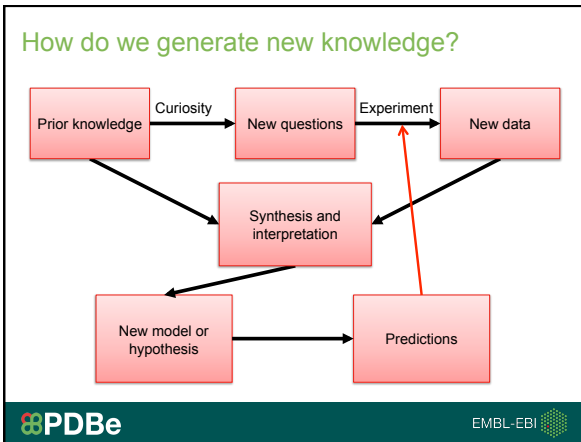
Ammonium sulfate

X = Me, NH<sub>2</sub>, OH

**PDBe** EMBL-EBI

### The *what* of validation

**PDBe** EMBL-EBI



### Errors affect measurements

- Random errors (noise)
  - Affect precision
  - Usually normally distributed
  - Reduce by increasing nr of observations
- Systematic errors (bias)
  - Affect accuracy
  - Incomplete knowledge or inadequate design
  - Reproducible
- Gross errors (bloopers)
  - Incorrect assumptions, undetected mistakes or malfunctions
  - Sometimes detectable as outliers

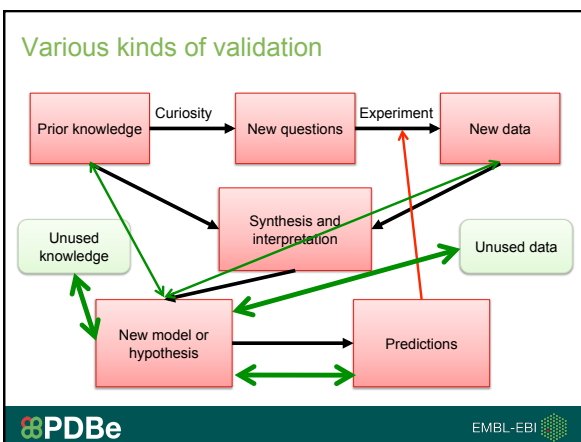
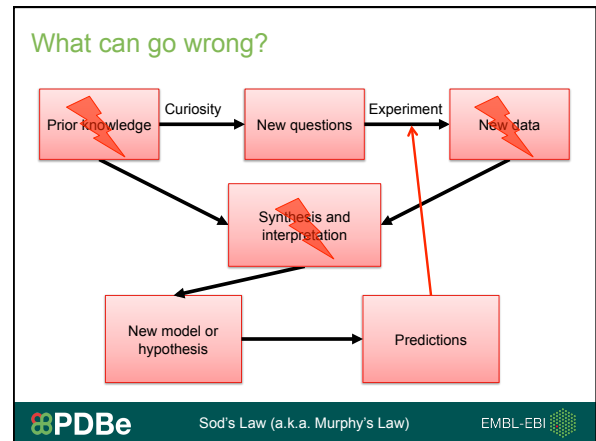
EMBL-EBI

### Errors affect measurements

- How tall is Gerard?
- 200 203 202 203 202
- 201 203
- Random error
- Systematic error
- Gross error

Anisotropic model of Gerard

EMBL-EBI



This model of hypothesis validation is entirely general for experimental sciences

How does it apply to protein crystallography?

EMBL-EBI



### The *how* of validation

### What is a good model?

- A good model makes SENSE in all respects!

### Various kinds of crystal structure validation

Flowchart illustrating the process of crystal structure validation:

- Prior knowledge** (with *Curiosity*) leads to **Unused knowledge**.
- Unused knowledge** leads to **New model or hypothesis**.
- A central green box lists validation criteria: **Geometry, Stereo-chemistry, Close contacts, Sequence, Chemical structure, Biosynthetic pathways, ...**

### Various kinds of crystal structure validation

Flowchart illustrating the cycle of crystal structure validation:

- Prior knowledge** leads to **Unused knowledge**.
- Unused knowledge** leads to **New model or hypothesis**.
- New model or hypothesis** leads to **Predictions**.
- Predictions** leads to **Experiment**.
- Experiment** leads to **New data**.
- New data** leads to **Unused data**.
- Unused data** leads back to **Unused knowledge**.
- A central green box lists validation criteria: **R-value, Real-space fit, B-values,  $k_{sol}$ , ...**

### Various kinds of crystal structure validation

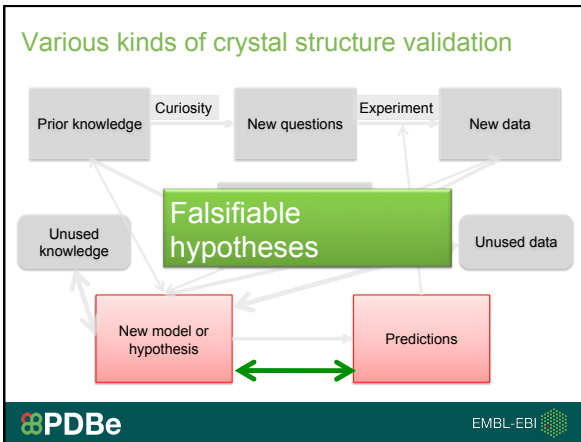
Flowchart illustrating the cycle of crystal structure validation:

- Prior knowledge** leads to **Unused knowledge**.
- Unused knowledge** leads to **New model or hypothesis**.
- New model or hypothesis** leads to **Predictions**.
- Predictions** leads to **Experiment**.
- Experiment** leads to **New data**.
- New data** leads to **Unused data**.
- Unused data** leads back to **Unused knowledge**.
- A central green box lists validation criteria:  **$R_{free}$ , Binding data, Mutant data, Conserved residues, Heavy-atom sites, SAXS envelope, ...**

### Various kinds of crystal structure validation

Flowchart illustrating the cycle of crystal structure validation:

- Prior knowledge** leads to **Unused knowledge**.
- Unused knowledge** leads to **New model or hypothesis**.
- New model or hypothesis** leads to **Predictions**.
- Predictions** leads to **Experiment**.
- Experiment** leads to **New data**.
- New data** leads to **Unused data**.
- Unused data** leads back to **Unused knowledge**.
- A central green box lists validation criteria: **Ramachandran, Rotamers, Environments, ...**



### Validation in a nutshell

- Compare your model to the experimental data and to the prior knowledge. It should:
  - Reproduce knowledge/information/data used in the construction of the model
    - R, RMSD bond lengths, chirality, ...
  - Predict knowledge/information/data not used in the construction of the model
    - $R_{free}$ , Ramachandran plot, packing quality, ...
  - Global and local
  - Model alone, data alone, fit of model and data
  - ... and if your model fails to do this, there had better be a plausible explanation!

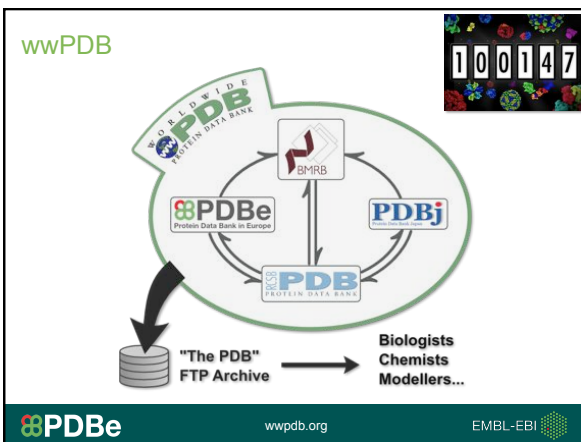
**PDBe** EMBL-EBI

### What is "the PDB" doing about validation?

**PDBe** EMBL-EBI

### What is "the PDB"?

**PDBe** EMBL-EBI




### wwPDB partnership

- Collaborate on "data in"
  - Policy issues
  - Weekly releases
  - Validation standards
  - Format specifications
  - Chemical Component Dictionary
  - Deposition and annotation procedures
  - Archive quality and remediation
  - Journal interactions
  - Community interactions
- Friendly competition on "data out"
  - Serving PDB data with added-value
  - PDB-based services
  - Other services, resources and activities

**PDBe** wwpdb.org

### Validation addresses important questions


- Entry-specific validation (quality control)
  - Is this model ready for archiving and publication?
  - Is this model a faithful, reliable and complete interpretation of the experimental data?
  - Are there any obvious errors/problems?
  - Are the conclusions drawn in the paper justified by the data?
  - Is this model suitable for my application?
- Archive-wide validation (comparative)
  - Is this model a better interpretation of the data?
  - What is the best model for this molecule/complex to answer my research question?
  - Which models should I select/omit when mining the PDB?



**PDBe** EMBL-EBI

### Validation by wwPDB - advantages

- Applies community-agreed methods uniformly
- Improves the quality and consistency of the PDB archive
- Supports editors and referees
- Helps users assess if an entry is suitable
- Helps users compare related entries
- Enables identification of outliers when mining the PDB
- Stimulates adoption of better protocols by the community



**PDBe** EMBL-EBI

### The future of validation


WORLDWIDE **PDB** PROTEIN DATA BANK

- wwPDB X-ray Validation Task Force

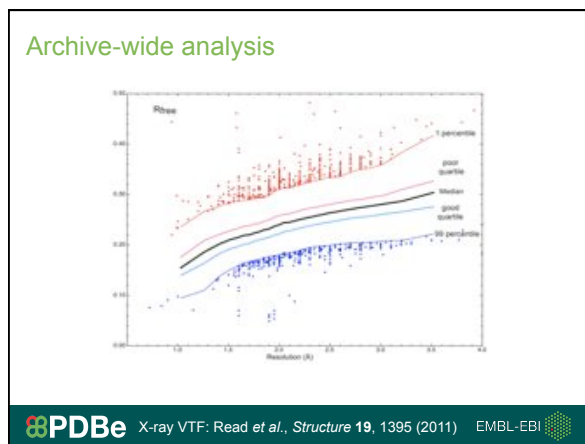
Structure **Ways & Means** Cell PRESS

#### A New Generation of Crystallographic Validation Tools for the Protein Data Bank

Randy J. Read,<sup>1\*</sup> Paul D. Adams,<sup>2</sup> W. Bryan Arendall, III,<sup>2</sup> Axel T. Brunger,<sup>4</sup> Paul Emsley,<sup>5</sup> Robbie P. Joosten,<sup>6,7</sup> Gerard J. Kleywegt,<sup>8,9</sup> Eugene B. Krissinel,<sup>10,11</sup> Thomas Luttmann,<sup>12</sup> Zbyszek Otwinowski,<sup>13</sup> Anastassis Perrakis,<sup>7</sup> Jane S. Richardson,<sup>3</sup> William H. Sheffler,<sup>14</sup> Janet L. Smith,<sup>15</sup> Ian J. Tickle,<sup>16</sup> Gert Vriend,<sup>8</sup> and Peter H. Zwart<sup>17</sup>



**PDBe** EMBL-EBI



### Percentile scores

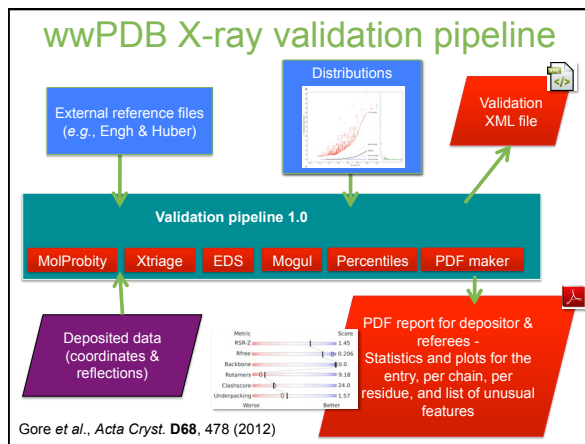
Metric	Score
Rfree	0.256
RSR-Z	0.123
Backbone	0.056
Rotamers	0.025
Clashscore	17.3
Underpacking	1.3
RNA puckers	0.031

Worse | Better

Legend: | Absolute Percentile, 0 Relative Percentile

[More Details](#)

**PDBe** EMBL-EBI



## What does it mean for a crystallographer?

- There are three uses of the validation pipeline
  - At deposition time
    - Not all checks can be run, e.g. some sequence and ligand checks
    - Report for depositor
  - At annotation time
    - Complete validation report, also suitable for editors/referees
  - Independently of deposition
    - Anonymous web-based server to use on models not (yet) in the PDB
    - Not all checks can be done
    - Will be developed once the production pipeline is up and running
    - Will not be available as a stand-alone software package

**PDBe** EMBL-EBI

## Validation reports

wwPDB X-ray Structure Validation Summary Report

Feb 6, 2015 - 12:00 PM GMT

PDB ID : 1CBS  
 Title : CRYSTAL STRUCTURE OF CELLULAR RETINOIC-ACID-BINDING PROTEINS 1 AND 2 IN COMPLEX WITH ALL-TRANS-RETINOIC ACID AND A SYNTHETIC RETINOID

Authors : Klywerg, G.J.; Berglin, T.; Jones, T.A.  
 Deposited on : 1996-09-28  
 Resolution : 1.80 Å (reported)

This is a wwPDB validation summary report for a publicly released PDB entry. We welcome your comments at validation@wwpdb.org. A user guide is available at <http://www.pdb.org/Validation/ValidationPDFViewer.html>

The following versions of software and data (see references) were used in the production of this report:

MultiSub : 4.05b-87  
 MolProbity : 1.1.7 November 2013  
 XrayDiffraction : 1.1.1  
 EDS : 1.0.0.10.0.0  
 Phenix : 1.8.5-1019  
 CCTF : 4.3.0 (2013)  
 Unit geometry (symm) : Eng & Baker (2011)  
 Unit geometry (DNA, RNA) : Parkison et al. (1996)  
 Validation Pipeline (wwPDB-V) : 1.0.0.10.0.0

**PDBe** [pdb.org/valrep/1cbs](http://pdb.org/valrep/1cbs) EMBL-EBI

## Validation reports

### 1 Overall quality at a glance

The reported resolution of this entry is 1.80 Å. Percentile scores (ranging between 0-100) for global validation metrics of the entry are shown in the following graphic. The table shows the number of entries on which the scores are based.

Metric	Percentile Score	Value
Ramachandran outliers	4.2%	1
Clash geometry	4.2%	1
Sidechain outliers	2.8%	1
RSR-Z outliers	2.8%	1

The table below summarizes the geometric issues observed across the polymer(s) chains and their fit to the electron density. The red, orange, yellow and green squares on the lower bar indicate the fraction of residues that contain outliers for  $\sigma > 2$ , 1 and 0 types of geometric quality criteria. The upper red bar (where present) indicates the fraction of residues that have poor fit to the electron density.

Mol. Chain	Length	Quality of chains
1 A	127	100%

The following table lists non-polymeric compounds, carbohydrate monomers and non-standard residues in protein, DNA, RNA chains that are outliers for geometric or electron-density fit criteria.

Mol. Chain	Type	Chain	Res	Chemistry	Geometry	Electron density
1	2	BDG	A	100	-	X
1	2	BDG	A	101	-	X
1	2	BDG	A	102	-	X

**PDBe** EMBL-EBI

## Validation reports

### 2 Entry composition

There are 3 unique types of molecules in this entry. The entry contains 1213 atoms, of which 9 are hydrogen and 0 are deuterium.

In the table below, the ZeroOcc column contains the number of atoms modelled with zero occupancy, the AltConf column contains the number of residues with at least one atom in alternate conformation and the Trace column contains the number of residues modelled with at most 2 atoms.

Mol. Chain	Residues	Atoms	ZeroOcc	AltConf	Trace
1 A	127	1094	0	0	0

Molecule 1 is a protein called CELLULAR RETINOIC ACID BINDING PROTEIN TYPE 1.

Molecule 2 is RETINOIC ACID (ligand code: REX) (Formula: C<sub>20</sub>H<sub>28</sub>O<sub>2</sub>).

Molecule 3 is water.

Mol. Chain	Residues	Atoms	ZeroOcc	AltConf
2 A	1	22	0	0

Mol. Chain	Residues	Atoms	ZeroOcc	AltConf
3 A	100	198	0	0

**PDBe** EMBL-EBI

## Validation reports

### Residue quality

- One plot per polymer
- Coloured by number of types of geometric outliers
- Grey if not modelled
- Red dots: poor density (RSR-Z > 2, as in EDS)

Molecule 1: MEMBRANE COPPER AMINE OXIDASE  
 Chain A: [Residue plot showing outliers and density]

Molecule 1: CELLULAR RETINOIC ACID BINDING PROTEIN TYPE II  
 Chain A: [Residue plot showing outliers and density]

Molecule 1: Aerobic glycerol-3-phosphate dehydrogenase  
 Chain A: [Residue plot showing outliers and density]

**PDBe** EMBL-EBI

## Validation reports

### 4 Data and refinement statistics

- "Table 1"
- Xtriage

Property	Value	Source
Space group	P 21 21 21	Depositor
Cell constants	45.65 Å 47.56 Å 77.61 Å	Depositor
a, b, c, α, β, γ	90.00° 90.00° 90.00°	Depositor
Resolution (Å)	8.00 - 1.80	Depositor
% Data completeness (in resolution range)	90.2 (3.00-1.80)	EDS
R <sub>int</sub>	90.5 (14.95-1.80)	EDS
R <sub>merge</sub>	(Not available)	Depositor
<math>\langle I/\sigma(I) \rangle > 1</math>	3.77 (at 1.79 Å)	Xtriage
Refinement program	X-PLOR	Depositor
R, R <sub>free</sub>	0.200 / 0.237	Depositor
R <sub>int</sub> test set	0.184 / 0.189	DCC
Wilson B-factor (Å <sup>2</sup> )	14.8	Xtriage
Anisotropy	0.134	Xtriage
Bulk solvent E <sub>sol</sub> (v(A <sup>3</sup> ), B <sub>sol</sub> (Å <sup>3</sup> ))	0.41, 58.9	EDS
Estimated twinning fraction	0.027 for k,h,l	Xtriage
L-test for twinning <sup>2</sup>	<math>\langle  L  \rangle = 0.51 < L^2 = 0.36</math>	Xtriage
Outliers	0 of 13678 reflections	Xtriage
F <sub>o</sub> -F <sub>c</sub> correlation	0.95	EDS
Total number of atoms	1213	wwPDB-V <sup>1</sup>
Average B, all atoms (Å <sup>2</sup> )	16.0	wwPDB-V <sup>1</sup>

Xtriage's analysis on translational NCS is as follows: The largest off-origin peak in the Patterson function is 9.26% of the height of the origin peak. No significant pseudotranslation is detected.

**PDBe** EMBL-EBI

## Validation reports

- Model quality
  - Bond lengths and angles
  - Torsion angles (Ramachandran, rotamers)
  - Clashes
  - Separately for standard residues, non-standard residues, ligands, carbohydrates
  - Generally: information about distribution, outlier stats, percentile scores, list of up to 5 (worst) outliers (full reports contain all outliers)

**5.3 Torsion angles**

**5.3.1 Protein backbone**

In the following table, the Percentile column shows the percent Ramachandran outliers of the chain as a percentile score with respect to all X-ray entries followed by that with respect to entries of similar resolution. The Analyzed column shows the number of residues for which the backbone conformation was analyzed, and the total number of residues.

Mol. Chain	Analyzed	Favoured	Allowed	Outliers	Percentile
1 A	125/127 (98%)	122 (98%)	3 (2%)	0	100

There are no Ramachandran outliers to report.

**5.3.2 Protein sidechains**

In the following table, the Percentile column shows the percent sidechain outliers of the chain as a percentile score with respect to all X-ray entries followed by that with respect to entries of similar resolution. The Analyzed column shows the number of residues for which the sidechain conformation was analyzed, and the total number of residues.

Mol. Chain	Analyzed	Rotameric	Outliers	Percentile
1 A	122/122 (100%)	120 (98%)	3 (2%)	60

All (2) residues with a non-rotameric sidechain are listed below:

Mol. Chain	Res	Type
1 A	11	ASN
1 A	27	ASP

## Validation reports

- Geometry validation of ligands and non-standard entities
  - Mogul (CCDC)
- wwPDB will get CSD coordinates for new and existing compounds (if they are available, of course)

**5.6 Ligand geometry**

Ligand is modelled in this entry.

In the following table, the Counts column lists the number of bonds (or angles) for which Mogul statistics could be generated, the number of bonds (or angles) that are observed in the model and the number of bonds (or angles) that are defined in the chemical component dictionary. The Link column lists molecule types, if any, to which the group is linked. The Z score for a bond length (or angle) is the number of standard deviations the observed value is removed from the expected value. A bond length (or angle) with  $|Z| > 2$  is considered an outlier worth inspection. RMSZ is the root-mean-square of all Z scores of the bond lengths (or angles).

Mol. Type	Chain	Res	Link	Counts	RMSZ	$\mu$ (Z) > 3	Counts	RMSZ	$\mu$ (Z) > 2	
2	REA	A	200	-	19,22/22	1.06	1 (5%)	26,30/30	1.00	2 (7%)

All (1) bond length outliers are listed below:

Mol. Chain	Res	Type	Atom1	Z	Observed(A)	Model(X)
2 A	200	REA	CH1-O	2.18	1.56	1.33

All (2) bond angle outliers are listed below:

Mol. Chain	Res	Type	Atom1	Z	Observed(A)	Model(X)
2 A	200	REA	CH1-CH2-O	-2.20	121.50	127.20
2 A	200	REA	CH1-CH2-O	2.22	126.93	124.64

There are no chirality outliers.  
There are no torsion outliers.  
There are no ring outliers.

## Validation reports

- Model/data fit proteins, DNA, RNA
  - RSR and RSR-Z (EDS)
- Ligands etc.
  - RSR and LLDF

**6 Fit of model and data**

**6.1 Protein, DNA and RNA chains**

In the following table, the column labelled '#RSRZ-Z' contains the number (and percentage) of RSRZ outliers, followed by percent RSRZ outliers for the chain as percentile scores relative to all X-ray entries and entries of similar resolution. The OWAB column contains the minimum, median, 90th percentile and maximum values of the occupancy-weighted average B-factor per residue. The column labelled 'Q<0.9' lists the number of (and percentage) of residues with an average occupancy less than 0.9.

Mol. Chain	Analyzed	#RSRZ-Z	#RSRZ-Z	OWAB(AV)	Q<0.9
1 A	671/735 (91%)	-0.22	1 (0%)	63	43, 64, 79, 87
1 B	671/735 (91%)	-0.24	2 (0%)	65	41, 60, 75, 81
1 C	671/735 (91%)	-0.26	1 (0%)	64	43, 64, 79, 88
1 D	671/735 (91%)	-0.23	2 (0%)	63	41, 60, 75, 81
All All	2684/2940 (91%)	-0.26	6 (0%)	62	41, 62, 77, 88

The worst 5 of 6 RSRZ outliers are listed below:

Mol. Chain	Res	Type	RSRZ
1 C	161	GLY	3.6
1 B	79	LEU	3.2
1 D	203	HIS	2.5
1 A	109	PRO	2.1
1 D	286	VAL	2.1

## Public X-ray Validation Reports

pdbe.org - rcsb.org - pdj.org EMBL-EBI

## Beta site at PDB

http://wwwdev.ebi.ac.uk/pdbe/entry/pdb/1cbs EMBL-EBI

## Other methods?

... crystallography. "These days, being a crystallographer is not good enough," says Gerard Kleywegt, a structural biologist at the European Bioinformatics Institute in Hinxton, who heads the European annex of the PDB.

Hybrid methods take an 'everything but the kitchen sink' approach to structural biology, incorporating many different techniques. Some can offer a dynamic view of a molecular machine in motion; for example, fluorescence resonance energy transfer measures the distance and interactions between proteins. Others, such as cryo-electron microscopy, can deliver near-atomic detail of entire complexes without the need to crystallize them. Computer programs

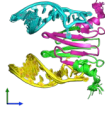

*"These days, being a crystallographer is not good enough."*

Nature 514, 416 (2014) EMBL-EBI



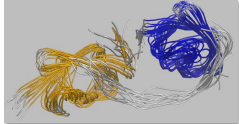
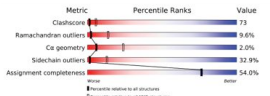
### Other Methods?

- Model validation using same criteria as X-ray
  - MolProbity, Mogul
- Some special model-related issues per technique
  - X-ray: alternative conformations
  - NMR: ensemble of models; well-defined regions
  - 3DEM: clashes of rigid-body fitted models; difference in species of model and sample sequence
- Data quality and model/data-fit assessment will be different for each technique


### NMR Validation

- NMR VTF recommendations published
- Global quality scores reported for “well-defined residues” only
  - As averages over the ensemble
  - Mediodi model only

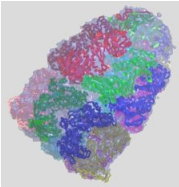

Metric	Whole archive (# Entries)	NMR archive (# Entries)
Clashscore	99129	10081
Ramachandran outliers	96105	8982
Cα geometry	96347	8988
Sidechain outliers	96047	8965
Assignment completeness	1540	1532

- Molecule 1: Small ubiquitin-related modifier 3
- Chain A: [Sequence]
- Molecule 2: Small ubiquitin-related modifier 3
- Chain B: [Sequence]
- Molecule 3: ES ubiquitin-protein ligase RNF4
- Chain C: [Sequence]



### 3DEM Validation

- Model validation
  - Clashes?
  - Taxonomy?
  - Homology models?
  - Non-atomistic models?
  - Ca-only models?
  - Rigid-body vs. flexible fitting vs. de novo modelling?
- Data and map validation
  - Per technique and resolution regime
  - Tilt-pair analysis; handedness; projections vs. raw data
  - Map + model
  - Depending on resolution regime and model-building method?


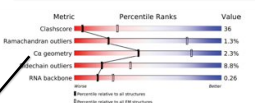



### EM Validation Reports

- Metrics relevant for EM models
- Define “Table 1” for EM

4 Experimental information

Property	Value	Source
Redundant method	Not provided	Depositor
Imposed symmetry	1	Depositor
Number of images	2600	Depositor
Resolution information method	2DCC at 0.143 cut-off	Depositor
CTF correction method	Each particle	Depositor
Microscope	CRYOEM	Depositor
Voltage (kV)	300	Depositor
Electron dose (e <sup>-</sup> /Å <sup>2</sup> )	15	Depositor
Maximum defocus (nm)	1000	Depositor
Minimum defocus (nm)	5000	Depositor
Magnification	50000	Depositor
Image distance	Kodak SO 163 film	Depositor

Metric	Whole archive (# Entries)	EM structures (# Entries)
Clashscore	99129	735
Ramachandran outliers	96105	539
Cα geometry	96347	682
Sidechain outliers	96047	526
RNA backbone	2607	214

wwPDB EM Map/Model Validation Report

Aug 22, 2014 - 03:59 PM BST


EMBI ID: EMB-3024

Title: Structures of the mammalian ubiquitin cycle by 400 subunit effector: a ubiquitin-specific protease reorganization

Authors: Hoshino, T., Ghoshal, J., Robinson, E., Goshima, J., Banerji, D.J.F., Miele, T., Lee, J., Hoshino, F., Tang, C.-S., Nathan, K.R., Sankaranarayanan, S.V., Dreyer, C.H.T.

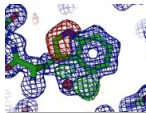
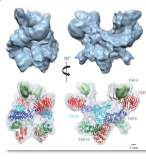

Deposited on: 2014-08-21

Resolution: 3.50 Å (reported)




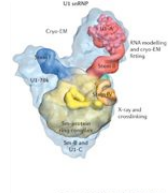
### Validation by wwPDB

- By no means the end of the story!
  - Room for extension and improvement
    - Ligands, nucleic acids, carbohydrates, NCS, spacegroup errors, ...
  - wwPDB ligand-validation workshop in 2015
  - X-ray
    - Re-convolve X-ray VTF in 2015 to evaluate and update recommendations
  - NMR
    - Further development in progress
  - EM
    - Rudimentary at present, lots more work needed
  - All methods: annual re-compute of distributions
  - User feedback welcome at [validation@mail.wwpdb.org](mailto:validation@mail.wwpdb.org)






### “Other other” methods

- SAS – wwPDB task force (2012, 2014)
- Hybrid methods – wwPDB task force (2014)
  - For example: solid-state NMR + EM + SAXS + solution NMR + homology modelling ...
- Questions
  - What to archive and where?
  - What to accept?
  - What requirements for deposition?
  - How to validate?
  - What to do with non-atomistic models?
  - What to do with homology models?





Copyright © 2008 Nature Publishing Group  
Nature Reviews Molecular Cell Biology



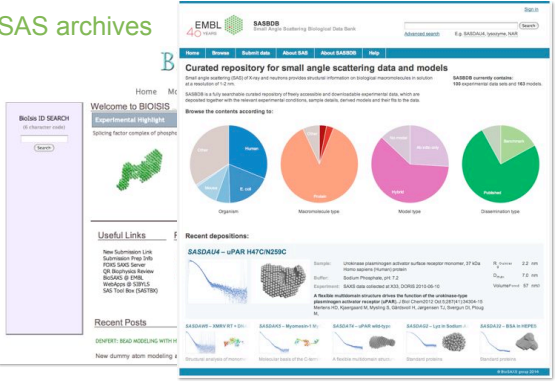
### SAS Task Force recommendations

- Need repository for SAXS and SANS data
- Need dictionary (data model) for SAXS and SANS
- Shape/bead and atomistic models should be archived (somewhere, somehow)
- Validation criteria need to be defined
- Archive of non-atomistic models from hybrid data
- What should (not) be in the PDB?



Trewhella *et al.*, *Structure* 21, 875 (2013)

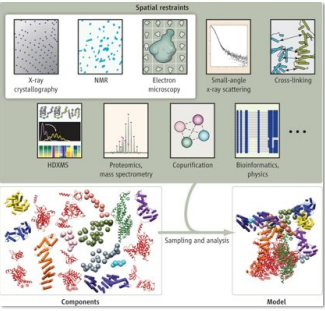
### SAS archives



bioisis.net – sasbdb.org

### Hybrid Methods

- Task Force met in October 2014
- Representatives of existing task forces, other methods, integrative modellers, and wwPDB
- Questions about what to archive where, what data and meta-data, how to validate



### wwPDB Hybrid Methods Task Force

EMBL-EBI, Hinxton, 6-7 October, 2014

#### Data bank struggles as protein imaging ups its game

Hybrid methods to solve structures of molecular machines create a storage headache.



Nature 514, 416 (2014)

### Key outcomes of discussion

- Be as inclusive as possible in collecting data from many different experimental methods
- Accommodate many types of structural representations
- Create a federated system to collect/curate data
- Use a common interface to collect data
- wwPDB should play a leadership role
- Whitepaper to describe vision

### What have we learned?

### Why do/did things sometimes go horribly wrong in X-ray?

- Blind optimism/naïveté/ignorance
  - Belief in (wrong) numbers and in "magic" refinement programs
- Inappropriate (use of) modelling/refinement methods
  - Fitting too many parameters
- No/inappropriate quality control/validation
- "Believing is seeing"
- Large influx of non-experts

EMBL-EBI

EMBL-EBI

### Of course, none of this should be news or surprising...

Hendrickson (CCP4 Proc., 1980) - "That which is not restricted will take its liberties"

Knight *et al.* (CCP4 Proc., 1990) - "None of this evidence is dependent on a refined model and instead makes use of known facts about proteins in general and the S subunit of RuBisCO in particular"



EMBL-EBI

EMBL-EBI

1990

Protein crystallography is a highly competitive field where the same or similar structures are being worked on in a number of different laboratories. Although this may result in an urge to publish quickly and prematurely, it is the responsibility of crystallographers to check their preliminary models carefully before rushing into print. Journals must insist on the publication of enough data for crystallographers to convince the reader that they have a correct structure, and the readership should be sophisticated enough to judge the quality of the data. We strongly object to publication of structural work where authors supply a minimum amount of detail in the form of a cartoon. □

Carl-Ivar Brändén and T. Alwyn Jones are in the Department of Molecular Biology, Uppsala Biomedical Center, PO Box 590, S-751 24 Uppsala, Sweden.



EMBL-EBI

Brändén & Jones, *Nature* 353, 687 (1990)

EMBL-EBI

### Lessons

- Have we learned anything from 25 years of errors?
  - Education is important
    - Avoid blind optimism, naïveté, belief in "magic" programs
    - Don't be afraid to ask a colleague's help or opinion
  - Use restraint and restraints when modelling
    - Consider the ratio of observations and parameters
    - Consider the information content of your data
  - **Null-hypothesis: everything is normal!**
    - Trans-peptides, bond lengths/angles, rotamers, NCS, ...
    - Unless your data shouts at you otherwise, or you have reliable prior knowledge

EMBL-EBI

EMBL-EBI

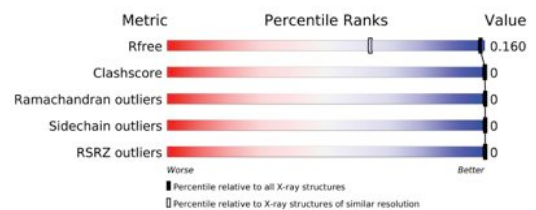
### Lessons

- Have we learned anything from 25 years of errors?
  - Use (lots of) validation tools throughout, not just when you deposit
    - Or worse, rely on wwPDB annotators to tell you what's dodgy about your model...
  - Be your own fiercest critic!
    - Avoid confirmation bias - try to shoot down your own models and hypotheses
    - How will you deal with cognitive dissonance?

EMBL-EBI

EMBL-EBI

### What you would like your plots to look like...



- Molecule 1: Potassium channel toxin ShK

Chain A:

There are no outlier residues recorded for this chain.

EMBL-EBI

pdbe.org/valrep/4lfq

EMBL-EBI

