

PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



... or a story about perceptions, expectations, naivety, macromolecular complexes and their complexity in bioinformatics and crystallography

Eugene Krissinel

Macromolecular Structure Database European Bioinformatics Institute, Genome Campus, Hinxton Cambridge CB10 1SD UK

http://www.ebi.ac.uk/msd

keb@ebi.ac.uk

2nd Annual CCP4 USA summer School and Workshop APS at ANL, Argonne, Illinois, 24 June - 2 July 2009.



http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html





Macromolecular Assemblies

- Complexes of protein, DNA/RNA chains and ligands, stable in native environment
- The way the chains assemble represents the [Protein] Quaternary Structure (PQS)
- Macromolecular assemblies are often the Biological Units, performing certain biochemical functions
- Biological significance of macromolecular assemblies is truly immense





European Bioinformatic





PQS is a difficult object for experimental studies

- Light / Neutron / X-ray / Small angle scatterings: mainly composition and multimeric state may be found. 3D shape may be guessed from mobility measurements.
- *Electron microscopy*: not a fantastic resolution and not applicable to all objects
- *NMR* is not good for big chains, even less so for protein assemblies.



Very few quaternary structures have been identified experimentally.







PQS are difficult to calculate

If we know the sequence

then we can calculate ...





10 - 90% Tertiary Structure (CASP 5), depending on method and target

Probably 0% Quaternary Structure. Docking of given number of given structures: 5 - 20% success (CAPRI 5)



PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



But PQS are assigned to many entries in PDB!







Most of those are **PROBABLE** Quaternary Structures.

The PDB "rules" are:

- 1. Depositor's say prevails.
- 2. Accept everything which passes formal validation checks.
- 3. No experimental evidence for PQS is required.
- 4. If a depositor does not know or does not care (60-80% of instances for PQS), the curator is to decide.
- 5. The curator may (EBI-MSD) or may not (Rutgers-RCSB, PDBj) use computing/modeling tools to assist the PQS annotation.







Crystallography is special

Because: A) crystal is made of assemblies







Crystallography is special

Because: B) there is no need to dock subunits - the docking is given by crystal structure



Macromolecular interfaces should be viewed as an additional important product of protein crystallography







PDB contains a wealth of experimental data on PQS

More than 80% of macromolecular structures are solved by means of X-ray diffraction on crystals.

It is reasonable to expect that PQS make construction blocks for the crystal.

An X-ray diffraction experiment produces atomic coordinates of the Asymmetric Unit (ASU), which is stored as a PDB file.

In general, neither ASU nor Unit Cell has any direct relation to PQS. The PQS may be made of

- a single ASU
- part of ASU
- several ASUseveral ASU parts

Crystal = translated Unit Cells





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



PDB entries and Biological Units



Biological unit 1P30 Homotrimer!





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



PDB entries and Biological Units



PDB entry 2TBV A trimer?





http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html



PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



PDB entries and Biological Units



PDB entry 1E94 ??????? 2 Biological Units in 1E94: A dodecamer and a hexamer!





In (very) simple words ...







A simple thing to do



PQS server @ EBI(Kim Henrick) Trends in Biochem. Sci. (1998) 23, 358PITA server @ EBI(Hannes Ponstingl) J. Appl. Cryst. (2003) 36, 1116







What is a significant interface?

Depends on the problem.

Protein functionality: the interface should be engaged in *any* sort of interaction, including transient short-living protein-ligand and protein-protein etc. associations. Obviously important properties:

• Affinity (comes from area, hydrophobicity, electrostatics, H-bond density etc.)

and properties that may be important for *reaction pathway and dynamics*:

- Aminoacid composition
- Geometrical complementarity
- Overall shape, compactness
- Charge distribution
- etc.

Stable macromolecular complexes, PQS: the interface should make a sound binding. Important properties:

- Sufficient free energy of binding
- something else?









A Common Identification Problem





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Real and superficial protein interfaces







Real and superficial protein interfaces

Most often used discrimination criteria - interface area.

A cut-off at 900 Å² gives about 80% success rate of discrimination between monomers and dimers.

Big proteins will be always sticky if this criteria is true ...







Real and superficial protein interfaces

Free energy gain of interface formation.

A cut-off at -8 kcal/M gives about 82% success rate of discrimination between monomers and dimers.

Can energy measure be uniform for all weights and shapes?







Is there a (good) measure of interface significance at all?







Real and superficial protein interfaces

- "No single parameter absolutely differentiates the interfaces from all other surface patches"
 - Jones, S. & Thornton, J.M. (1996) Principles of protein-protein interactions, *Proc. Natl. Acad. Sci. USA*, **93**, 13-20.
- Formation of N>2 -meric complexes is most probably a corporate process involving a set of interfaces. Therefore significance of an interface should not be detached from the context of macromolecular complex
- "…the type of complexes need to be taken into account when characterizing interfaces between them."
 - Jones, S. & Thornton, J.M., ibid.







Chemical stability of macromolecular complexes

- It is not properties of individual interfaces but rather chemical stability of complexes in general that really matters
- Macromolecular units will most likely associate into largest complexes that are still stable
- ✤ A complex is stable if its free energy of dissociation is positive:

$$\Delta G_{diss}^0 = -\Delta G_{\rm int} - T\Delta S > 0$$



European Bioinformati







Chemical stability of macromolecular complexes

$$\Delta G_{diss} = -\Delta G_{int} - T\Delta S > 0$$





http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html











Solvation free energy

Eisenberg, D. & McLachlan, A.D. (1986) *Nature* 319, 199-203.

Atomic solvation parameters

Atom's accessible surface area

 $\Delta G_{sol}(A) = \sum_{k} \Delta \sigma_k \left(a_k - a_k^r \right)$

Atom's accessible surface area in the reference (unfolded) state



In this model, binding energy is function of individual interfaces.



European Bioinformatics In

http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html





Entropy of macromolecules in solutions



$$S_{trans}(m) \approx c_t + \frac{3R}{2}\log(m)$$

$$S_{rot}(\hat{I},\sigma_S) \approx c_r + \frac{R}{2}\log(I_1I_2I_3/\sigma_S^2)$$

Murray C.W. and Verdonk M.L. (2002) J. Comput.-Aided Mol. Design 16, 741-753.

$$S_{surf}(a) \approx Fa$$

 C_t, C_r and F are semiempirical parameters





=(n-1)



Entropy of dissociation

$$\Delta S = \sum_{i=1}^{n} S(A_i) - S(A_1, A_2 \dots A_n)$$

 $\frac{R}{2}\log$

Mass of i-th subunit

k-th principal moment of inertia of i-th subunit

$$\left(\frac{1}{A_n}\right) + Fa_{buried}$$

Fitted parameter

By its very nature, entropy of dissociation is function of protein complex rather than that of individual interfaces.



http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html

Fitted parameter





Entropy of dissociation

- Drives thermodynamic systems towards most disordered (dissolved) state.
- Makes bias towards less symmetric states. However, in practice, this is overweighed by binding energy, which is normally maximal in most symmetric states.
- Responsible for complex "instability".













PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.









We know how it looks, can we see it?

We now know (or we think that we know) how to evaluate chemical stability of protein complexes:

$$\Delta G_{diss}^0 = -\Delta G_{\rm int} - T\Delta S > 0$$

which depends on interface properties as well as on the geometry of protein complex and dissociated subunits.

How to find stable complexes in a crystal?



European Bioinformatic





Enumerating assemblies in crystal

- crystal is represented as a periodic graph with monomeric chains as vertices and interfaces as edges
- each set of assemblies is identified by engaged interface types
- all assemblies may be enumerated by a backtracking scheme engaging all possible combinations of different interface types



Example: crystal with 3 interface types

sembly set	Engag interface	led types	Assembly set in	Engaç nterface	jed types
1	000		5	100	- dimer N3
2	001	- dimer N1	6	101	
3	010	- dimer N2	7	110	
4	011		8	111	







Detection of Biological Units in Crystals: Method Summary

- 1. Build periodic graph of the crystal
- 2. Enumerate all possibly stable assemblies
- 3. Evaluate assemblies for chemical stability
- 4. Leave only sets of stable assemblies in the list and range them by chances to be a biological unit :
 - Larger assemblies take preference
 - Single-assembly solutions take preference
 - Otherwise, assemblies with higher ΔG_{diss} take preference







Classification of protein assemblies

Assembly classification on the benchmark set of 218 protein structures published in

Ponstingl, H., Kabir, T. and Thornton, J. (2003) Automatic inference of protein quaternary structures from crystals. J. Appl. Cryst. 36, 1116-1122.

	1mer	2mer	3mer	4mer	6mer	Other	Sum	Correct
1mer	49	3	0	1	1	1	55	89%
2mer	3	71+11	0	2+1	0	0	76+1 <u>2</u>	93%
3mer	1	0	22	0	1	0	24	92%
4mer	2	2+1	0	26+ <mark>6</mark>	0	1	31+7	84%
6mer	0	0	0	0+1	10+ <mark>2</mark>	0	10+ <mark>3</mark>	92%
196+22 <=> 196 homomers and 22 heteromers					Total:	196+22	90%	

kcal/mol

kcal/mol

kcal/mol

 $T \cdot F = 0.57 \cdot 10^{-3} \text{ kcal/(mol*Å^2)}$

E_{hb}

= 0.51

 $E^{nv} = 0.21$ $T^{sb} \cdot C = 11.7$

Fitted parameters:

- 1. Free energy of a H-bond :
- 2. Free energy of a salt bridge :
- 3. Constant entropy term :
- 4. Surface entropy factor :

Classification error in ΔG_{diss}^0 : ± 5 kcal/mol







Classification of protein-DNA complexes

Assembly classification on the benchmark set of 212 protein – DNA complexes published in

Luscombe, N.M., Austin, S.E., Berman H.M. and Thornton, J.M. (2000) An overview of the structures of protein-DNA complexes. Genome Biol. 1, 1-37.

	2mer	3mer	4mer	5mer	6mer	10mer	Other	Sum	Correct
2mer	1	0	0	0	0	0	0	1	100%
3mer	6	96	0	0	1	0	2	105	91%
4mer	0	2	83	0	0	0	0	85	98%
5mer	0	0	2	3	0	0	0	5	60%
6mer	1	0	0	0	13	0	1	15	87%
10mer	0	0	0	0	0	1	0	1	100%
							Total:	212	93%









Free energy distribution of misclassifications



http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html



PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1QEX

BACTERIOPHAGE T4 GENE PRODUCT 9 (GP9), THE TRIGGER OF TAIL CONTRACTION AND THE LONG TAIL FIBERS CONNECTOR



Predicted: homohexamer

Dissociates into 2 trimers $\Delta G_{diss}^0 \approx 106 \text{ kcal/mol}$



Biolgical unit: homotrimer

Dissociates into 3 monomers $\Delta G_{diss}^0 \approx 90 \text{ kcal/mol}$





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1QEX

BACTERIOPHAGE T4 GENE PRODUCT 9 (GP9), THE TRIGGER OF TAIL CONTRACTION AND THE LONG TAIL FIBERS CONNECTOR



Rossmann M.G., Mesyanzhinov V.V., Arisaka F and Leiman P.G. (2004) *The bacteriophage T4 DNA injection machine*. Curr. Opinion Struct. Biol. **14**:171-180.



http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html



PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1QEX

BACTERIOPHAGE T4 GENE PRODUCT 9 (GP9), THE TRIGGER OF TAIL CONTRACTION AND THE LONG TAIL FIBERS CONNECTOR





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1D3U

TATA-BINDING PROTEIN / TRANSCRIPTION FACTOR



Predicted: octamer

Dissociates into 2 tetramers $\Delta G_{diss}^0 \approx 20 \text{ kcal/mol}$

Functional unit: tetramer





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1CRX

CRE RECOMBINASE / DNA COMPLEX REACTION INTERMEDIATE



Predicted: dodecamer

Dissociates into 2 hexamers $\Delta G_{diss}^0 \approx 28 \text{ kcal/mol}$

Functional unit: trimer



http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html



PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1CRX

CRE RECOMBINASE / DNA COMPLEX REACTION INTERMEDIATE



Guo F., Gopaul D.N. and van Duyne G.D. (1997)

Structure of Cre recombinase complexed with DNA in a sitespecific recombination synapse.

Nature 389:40-46.







Example of misclassification: 1TON

TONIN

Predicted: dimer

Dissociates at $\Delta G_{diss}^0 \approx 37 \text{ kcal/mol}$

Biological unit: monomer

Apparent dimerization is an artefact due to the presence of Zn⁺² ions added to the buffer to aid crystallization. Removal Zn from the file results in $\Delta G_{diss}^0 \approx 3 \text{ kcal/mol}$

Fujinaga M., James M.N.G. (1997) *Rat submaxillary* gland serine protease, tonin structure solution and refinement at 1.8 Å resolution. J.Mol.Biol. **195**:373-396.





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1YWK







Structural homologue 1XRU: RMSD ≈ 0.9 Å Seq.Id $\approx 50\%$ Homohexameric with $\Delta G_{diss} \approx$ 9.3 kcal/mol

Predicted: homohexameric $\Delta G_{diss} \approx$ 4.4 kcal/mol dissociating into 3 dimers

http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html







Choice of ASU



http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html



PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Example of misclassification: 1YWK







Structural homologue 1XRU: RMSD ≈ 0.9 Å Seq.Id $\approx 50\%$ Homohexameric with $\Delta G_{diss} \approx 9.3 \text{ kcal/mol}$

Predicted: homohexameric $\Delta G_{diss} \approx$ 4.4 kcal/mol dissociating into 3 dimers

6 units in ASU







Reasons for the misclassification of Biological Units

- Crystal packing may introduce interactions that are not engaged in native environment
- Definition of Biological Unit is based on physiological function and may be to a certain degree subjective
- Interaction artifacts are often due to the addition of binding agents in order to aid crystallization

- Theoretical models of macromolecular complexation are simplified
- Experimental data are of a limited accuracy
- No explicit account for concentrations, temperature, pH and ionic strength

50%

50%



European Bioinformati



PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.



Comparison with current PDB annotation







Web-server PISA http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html

A new MSD-EBI tool for working with *P*rotein *Interfaces*, *Surfaces* and *Assemblies*







Session 564-RI-ACR map





PISA, or a story about macromolecular complexes in bioinformatics and crystallography 2nd Annual CCP4 USA Summer School and Workshop, APS at ANL, Illinois, 24 June – 2 July 2009.









Conclusions

- Chemical-thermodynamical models for protein complex stability allow one to recover biological units from protein crystallography data at 80-90% success rate
- Considerable part of misclassifications is due to the difference of experimental and native environments and artificial interactions induced by crystal packing
- Functional significance of protein interfaces cannot be reliably inferred only from their properties. Due to entropy contribution and entangled interactions, interface function is also subject to protein complex composition and geometry.
- Protein interface and assembly analysis software (PISA) is available, please use it







Acknowledgements

Kim Henrick European Bioinformatics Institute

Mark Shenderovich Structural Bioinformatics Inc.

Hannes Ponstingl Sanger Centre

Sergei Strelkov University of Leuven

MSD & PDB teams EBI & Rutgers

CCP4 Daresbury-York-Oxford-Cambridge

~4000 PISA users Worldwide

Biotechnology and Biological Sciences Research Council (BBSRC) UK General introduction, job assignment and gentle push

Helpful discussion

Sharing the expertise and benchmark data

"Mystery" of bacteriophage T4

Everyday use of PISA, examples, verification and feedback

Encouragement and publicity

Using PISA and feedback

Research grant No. 721/B19544



European Bioinformatics

http://www.ebi.ac.uk/msd-srv/prot_int/pistart.html