

Fig. 1. Frequency distributions. (A) pI of successfully crystallized proteins. (B) Reported pH of crystallization for proteins. (C) pI of successfully crystallized protein–nucleic acid complexes. (D) Reported pH of crystallization for protein–nucleic acid complexes.

pH screening in protein crystallization and the corresponding reduction in material requirements could lead to cost savings of millions of US\$ for structural genomics projects using high-throughput crystallographic structure determination.

SYSTEMS AND METHODS

We have used the SEQRES records of 9596 PDB entries comprising a non-redundant protein data set (Kantardjieff and Rupp, 2003), which contain the sequence of the entire expressed construct, including any tags, fusions or linkers, to calculate the pI using the pK_a values of Bjellqvist *et al.* (1993), and we have treated complexes of proteins and nucleic acids (469 entries) as a separate group. The frequency distribution for pI of proteins is bimodal (Fig. 1A), with highest frequencies (modes) at approximately pH 5.7 and 9.0, similar to the pI distribution seen for proteins encoded by sequenced genomes (see, e.g. Baisnee *et al.*, 2001; Urquhart *et al.*, 1998;

Adams *et al.*, 2003). The frequency distribution for the reported crystallization pH of proteins is unimodal, with mean = 6.7, median = 6.9 and mode = 7.5 (Fig. 1B). For the complexes, we observe a similar bimodal distribution of pI, with modes at 6.1 and 9.5 (Fig. 1C), and a unimodal distribution of crystallization pH, with mean = 6.6, median = 6.5 and mode = 6.5 (Fig. 1D). A similar distribution of crystallization pH has been observed from successful crystallizations of proteins resulting from unbiased random screening experiments in a structural genomics initiative (Rupp, 2003b).

We find that while there is no statistically significant direct correlation between the pI of a crystallized protein and the pH of crystallization, there is a good correlation ($R^2 = 0.62$) between the pI of a crystallized protein and the difference between the pH of crystallization and pI (Fig. 2). The delta (pH – pI) histograms for acidic and basic proteins are shown in Figure 3. It is apparent that acidic proteins crystallize with highest likelihood ~0–2.5 pH units above their pI, whereas basic proteins preferably crystallize ~0.5–3 pH units below

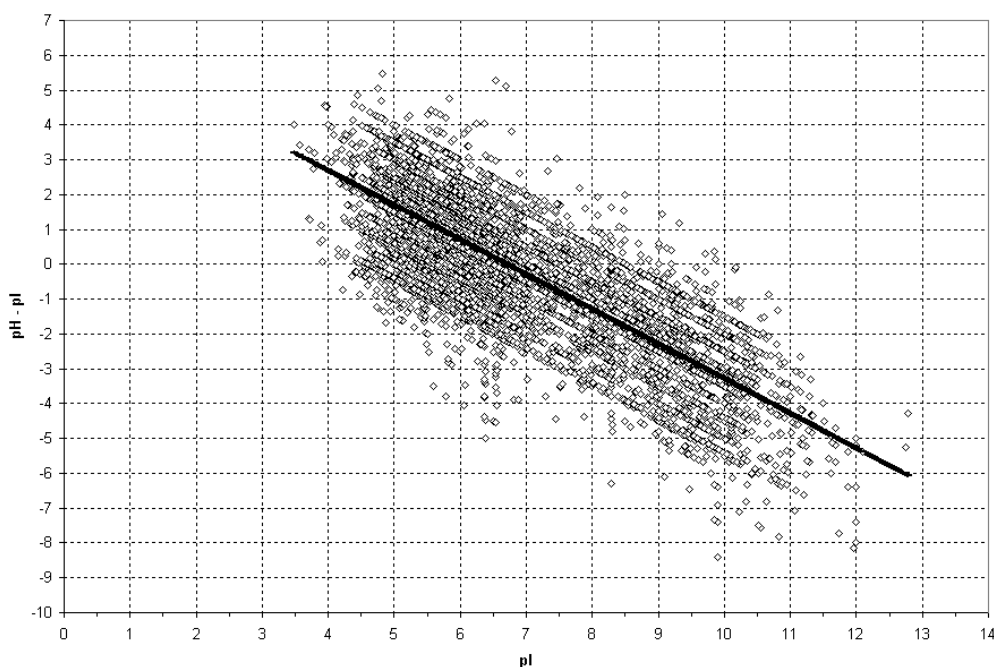


Fig. 2. Correlation between pI and pH. Correlation between calculated pI of successfully crystallized protein and difference between reported crystallization pH and pI. $R^2 = 0.62$, P -value $< 10^{-7}$. Not shown, protein–nucleic acid complexes ($R^2 = 0.77$, P -value $< 10^{-7}$).

their pI. Extreme values of pH do not contribute significantly to successful crystallization for most proteins, except for those that have unusually high or low pI values. For nucleic acid-bound proteins (data not shown), the correlation is also strong ($R^2 = 0.77$), with similar tendencies for optimal pH of crystallization, ~ 0 – 2 pH units above the pI for acidic proteins, ~ 2 – 4 pH units below the pI for basic proteins. We have not accounted for the pI of DNA (pH ~ 4), however, which generally lacks functional groups that change ionization state near physiological pH (Peek and Williams, 2001). Although conditions for crystallizing DNA–protein complexes have been shown to be similar to protein-only crystallization conditions, we do not use this last correlation for predictive purposes due to the above-mentioned uncertainties, as well as the limited number of data points.

IMPLEMENTATION

To demonstrate the utility of our analysis, we have implemented a prototype pH range calculator, CrysPred (<http://www-structure.llnl.gov/cryspred/>). The purpose of this small server-based applet is to show how prior information can be used to optimize efficiency of initial crystallization screening in HTPX. Effective initial crystallization screening aims to identify with the highest overall efficiency (least material, supplies and resources and thus cost) the proteins that are most likely to yield useful or suitable crystals and structures. The purpose of efficient initial screening is not to find

conditions for each and every protein but to focus resources (scale-up, Se-Met incorporation, etc.) on those proteins that have the highest probability to yield structures with the least effort (a.k.a. ‘the first cut’, ‘cherry picking’, etc.).

CrysPred accepts as input the amino acid sequence of the protein moiety to be crystallized, including the sequence of any tags, linkers or fusions, if present, and the number of crystallization experiments to be attempted. The program returns the calculated pI for the protein as well as a histogram showing the ‘delta’ bins (pH–pI) for successfully crystallized proteins with similar pI, grouped in clusters of two pH units. A table is provided also, showing the delta bin frequency expressed as a percentage of the pI cluster, the population of experiments (equal distribution) for a random screen, the recommended population of experiments based on the ‘delta’ prior information, and a suggested range of pH for the specified experiments (Fig. 4). Finally, CrysPred estimates the expected efficiency increase compared to pH screening with equally populated bins of each pH over the selected range. Depending on the shape of the corresponding frequency distribution and the extent of the pH sampling range, the total savings of material is predicted typically to be between 30 and 50%.

The values from CrysPred can be easily imported into any customizable screen generator that allows us to define the frequency of occurrence for selected pH ranges [e.g. CrysTool; Segelke and Rupp (1998) and Segelke (2001)]. The pH frequency distribution data are available for download from the CrysPred site to allow a custom implementation if desired.

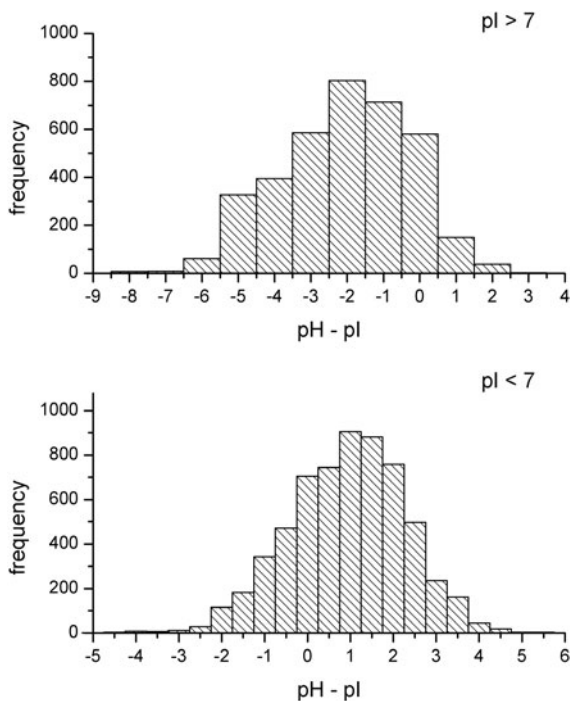


Fig. 3. Delta histograms for successfully crystallized proteins. Top panel shows frequency distribution of the difference between crystallization pH and pI of successfully crystallized basic proteins. Bottom panel shows this frequency distribution for acidic proteins. It is clear that basic proteins have a tendency to crystallize 0.5–3 pH units below their pI, whereas acidic proteins prefer to crystallize 0–2.5 pH units above their pI. Similar tendencies are observed for protein–nucleic acid complexes, although this is shifted 2–4 units below the pI for basic proteins.

DISCUSSION

Methods for choosing protein crystallization conditions have largely been empirical, based on knowledge of what has worked in the past (McPherson, 1982). More recently, random screening methods have been developed (Segelke and Rupp, 1998; Segelke, 2001), and it is anticipated that statistical analysis will provide predictive frameworks that increase the probability of producing high-quality crystals. Because pH is one of the few consistently reported crystallization parameters in the PDB, we have completed such a statistical analysis and implemented into a predictive framework called CrysPred the significant relationship between calculated pI of successfully crystallized proteins and the reported pH at which they were crystallized.

Crystallization is a special case of phase separation from a thermodynamically metastable solution under the control of kinetic parameters (Rupp, 2003b). While control over kinetic parameters such as nucleation or growth rates is rather difficult to achieve, attractive interaction among molecules as a thermodynamically necessary—but not sufficient—condition for crystallization can be discussed on the basis of

thermodynamic excess properties, in particular their manifestation in the second virial coefficient, B_{22} , as determined by static light scattering and osmotic pressure measurements.

More than fifty years ago, Zimm (1946) examined theoretically the osmotic second virial coefficient of proteins, B_{22} . At the molecular level, B_{22} reflects the nature of protein–protein interactions, which involve van der Waals attractions, electrostatic repulsions, non-centrosymmetric dipole interactions, hydrophobic interactions, hydrogen bonding and ion-bridge mechanisms. More negative values of B_{22} are indicative of more attractive interactions. Protein solubility is affected by solvent and additives, which alter protein size and surface characteristics (Farnum and Zukoski, 1999). Quantitative links between the second virial coefficient and solubility have suggested that large classes of globular proteins will exhibit similar solubility with the same normalized B_{22} (Fine *et al.*, 1996; Rosenbaum *et al.*, 1996; Rosenbaum and Zukoski, 1996). A number of groups (Farnum and Zukoski, 1999; George and Wilson, 1994; George *et al.*, 1997; Bonnete *et al.*, 1999, 2001; Beretta *et al.*, 2000; Tardieu *et al.*, 2001) have shown that for proteins under conditions where they were crystallized, the second virial coefficient is negative, falling in a narrow range termed the ‘crystallization slot’ (George and Wilson, 1994), and it is well documented that protein crystallization occurs in or close to attractive regimes (Tardieu *et al.*, 2001). Tardieu *et al.* (2001) have recommended that to crystallize soluble proteins (starting from a monodisperse solution), one should start far from precipitation and gently adjust repulsive interactions toward more attractive ones. However, although interactions tend to be attractive near the pI, in accord with the van der Waals potential, van der Waals forces are considerable only for small compact proteins (Tardieu *et al.*, 2001).

A number of studies on protein solutions and crystals (Haas *et al.*, 1999; Haas and Drenth, 1999; Ruppert *et al.*, 2001; Sear, 2002) have shown that protein–protein interactions can be described by a sum of surface contacts between proteins, but that the mutual arrangement of proteins requires some anisotropy (Beretta *et al.*, 2000; Haas *et al.*, 1999; Haas and Drenth, 2000) or complementarity (molecular recognition) (Neal *et al.*, 1998). Neal *et al.* (1998) have shown that as pH values approach pI, and charge and repulsive interactions are decreased, B_{22} becomes more negative at constant values of ionic strength. The magnitude of repulsive interactions and appearance of attractive interactions depend on the spatial distribution of charges and not simply on the global net charge of the protein, although accounting for short-range effects due to hydrogen bonding and solvation is not straightforward. Whereas changing the pH to approach the pI reduces the overall protein charge and decreases longer range electrostatic repulsion, Debye–Hückel screening of repulsive charge interactions may be exploited by searching for crystals under conditions of pH away from the pI (Juarez-Martinez *et al.*, 2001). B_{22} (and the possibility to crystallize) is determined

Estimated pI = 6.15

Table for cutoff excluding bins with expected success rates below 0.1%

pH-pI bin	: -8.0	-7.0	-6.0	-5.0	-4.0	-3.0	-2.0	-1.0	0.0	1.0	2.0	3.0	4.0	5.0	6.0
Expected %	: 0.0	0.0	0.0	0.0	0.0	0.3	1.1	6.2	14.6	25.6	28.9	19.0	3.8	0.5	0.0
Population of 288 experiments in 9 bins :															
equal pop.	: 0	0	0	0	0	32	32	32	32	32	32	32	32	32	288
suggested	: 0	0	0	0	0	1	3	17	41	73	83	54	10	1	0
Expected relative hit rates															
equal pop.	: 0.0	0.0	0.0	0.0	0.0	0.1	0.4	2.0	4.7	8.2	9.3	6.1	1.2	0.2	0.0
suggested	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.1	6.1	18.9	24.1	10.4	0.4	0.0
pH	---	---	---	---	---	3.1	4.1	5.1	6.1	7.1	8.1	9.1	10.1	11.1	---
Experiments:	---	---	---	---	---	1	3	17	41	73	83	54	10	1	---

Expected efficiency increase compared to pH screening with equally populated bins: 91%

Table for cutoff excluding bins with expected success rates below 1.0%

pH-pI bin	: -8.0	-7.0	-6.0	-5.0	-4.0	-3.0	-2.0	-1.0	0.0	1.0	2.0	3.0	4.0	5.0	6.0
Expected %	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.1	6.2	14.6	25.6	28.9	19.0	3.8	0.0
Population of 288 experiments in 7 bins :															
equal pop.	: 0	0	0	0	0	0	41	41	41	41	41	41	41	0	287
suggested	: 0	0	0	0	0	0	3	17	42	74	84	55	10	0	0
Expected relative hit rates															
equal pop.	: 0.0	0.0	0.0	0.0	0.0	0.0	0.5	2.5	6.0	10.5	11.9	7.8	1.5	0.0	0.0
suggested	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.1	6.2	19.0	24.3	10.5	0.4	0.0	0.0
pH	---	---	---	---	---	---	4.1	5.1	6.1	7.1	8.1	9.1	10.1	---	---
Experiments:	---	---	---	---	---	---	3	17	42	74	84	55	10	---	---

Expected efficiency increase compared to pH screening with equally populated bins: 51%

Table for cutoff excluding bins with expected success rates below 2.0%

pH-pI bin	: -8.0	-7.0	-6.0	-5.0	-4.0	-3.0	-2.0	-1.0	0.0	1.0	2.0	3.0	4.0	5.0	6.0
Expected %	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	6.2	14.6	25.6	28.9	19.0	3.8	0.0
Population of 288 experiments in 6 bins :															
equal pop.	: 0	0	0	0	0	0	0	48	48	48	48	48	48	0	288
suggested	: 0	0	0	0	0	0	0	18	42	75	85	55	11	0	0
Expected relative hit rates															
equal pop.	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	3.0	7.0	12.3	13.9	9.1	1.8	0.0	0.0
suggested	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.1	6.2	19.2	24.6	10.6	0.4	0.0	0.0
pH	---	---	---	---	---	---	---	5.1	6.1	7.1	8.1	9.1	10.1	---	---
Experiments:	---	---	---	---	---	---	---	18	42	75	85	55	11	---	---

Expected efficiency increase compared to pH screening with equally populated bins: 32%

Table for cutoff excluding bins with expected success rates below 5.0%

pH-pI bin	: -8.0	-7.0	-6.0	-5.0	-4.0	-3.0	-2.0	-1.0	0.0	1.0	2.0	3.0	4.0	5.0	6.0
Expected %	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	0.0	6.2	14.6	25.6	28.9	19.0	0.0	0.0
Population of 288 experiments in 5 bins :															
equal pop.	: 0	0	0	0	0	0	0	57	57	57	57	57	0	0	285
suggested	: 0	0	0	0	0	0	0	18	44	78	88	58	0	0	0
Expected relative hit rates															
equal pop.	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	3.5	8.3	14.6	16.5	10.8	0.0	0.0	0.0
suggested	: 0.0	0.0	0.0	0.0	0.0	0.0	0.0	1.2	6.5	20.0	25.6	11.0	0.0	0.0	0.0
pH	---	---	---	---	---	---	---	5.1	6.1	7.1	8.1	9.1	---	---	---
Experiments:	---	---	---	---	---	---	---	18	44	78	88	58	---	---	---

Expected efficiency increase compared to pH screening with equally populated bins: 19%

5.5 - 7.5

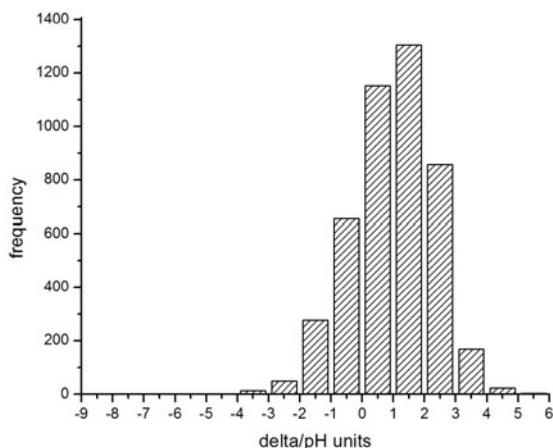


Fig. 4. Sample CrysPred output. Shown are the calculated pI for the protein and a histogram of the 'delta' bins (pH – pI) for successfully crystallized proteins with similar pI, grouped in clusters of two pH units. Table reports the delta bin frequency expressed as a percentage of the pI cluster, the equal population of experiments for a random screen, the recommended population of experiments based on the pH prior information, and a suggested range of pH for the specified experiments. Expected efficiency increase compared with pH screening with equally populated bins of each pH over the selected range is also predicted (in this example, between 19 and 92%).

largely by relatively few attractive interactions, the molecular configurations of which are influenced by pH and ionic strength.

Thus, while buffering at a pH equal to or very near the pI value of a protein offers a reasonable probability of yielding crystals, this pH is not necessarily that value with the highest probability. The 'knowledge' occasionally perpetuated at protein crystallization workshops and by unreviewed publications that a protein has the best chance of crystallizing at a pH very near its solubility minimum, pI, is not reflected statistically in the PDB data. We have found one commercially available crystallization screen that recommends empirically, 'The high efficiency of this kit can be further improved by pre-determining the isoelectric point (pI) of the subject macromolecule, followed by screening within a range at or near that value (within 2–3 pH units of the pI)' (Harris and McPherson, 1995). Our statistical analysis suggests optimal pH ranges for crystallization screening, and to improve efficiency of any crystallization screen, we recommend that the pI of the protein moiety to be crystallized be used to design an optimized pH distribution for incorporation into screening experiments.

ACKNOWLEDGEMENTS

LLNL is operated by University of California for the US DOE under contract W-7405-ENG-48. This work was funded by NIH P50 GM62410 (TB Structural Genomics) center grant and the Robert Welch Foundation. The W.M. Keck Foundation Center for Molecular Structure is a core facility of the California State University Program for Education and Research in Biotechnology.

REFERENCES

- Adams, M.W., Dailey, H.A., De Lucas, L.J., Luo, M., Prestegard, J.H., Rose, J.P., Wang, B.C. (2003) The Southeast Collaboratory for Structural Genomics: a high-throughput gene to structure factory. *Acc. Chem. Res.*, **36**, 191–198.
- Baisnee, P.F., Baldi, P., Brunok, S. and Pedersen, A.G. (2001) Flexibility of the genetic code with respect to DNA structure. *Bioinformatics*, **17**, 237–248.
- Beretta, S., Chirico, G. and Baldini, G. (2000) Short-range interactions of globular proteins at high ionic strengths. *Macromolecules*, **33**, 8663–8670.
- Bjellqvist, B., Hughes, G.J., Pasquali, C., Paquet, N., Ravier, F., Sanchez, J.C., Frutiger, S. and Hochstrasser, D. (1993) The focusing positions of polypeptides in immobilized pH gradients can be predicted from their amino acid sequences. *Electrophoresis*, **14**, 1023–1031.
- Bodenstaff, E.R., Hoedemaeker, F.J., Kuil, M.E., de Vrind, H.P. and Abrahams, J.P. (2002) The prospects of protein nanocrystallography. *Acta Crystallogr. D Biol. Crystallogr.*, **58**, 1901–1906.
- Bonnete, F., Finet, S. and Tardieu, A. (1999) Second virial coefficient: variations with lysozyme crystallization conditions. *J. Cryst. Growth*, **196**, 403–414.
- Bonnete, F., Vivares, D., Robert, C. and Colloch, N. (2001). Interactions in solution and crystallization of *Asperigillus flavus* urate oxidase. *J. Cryst. Growth*, **232**, 330–339.
- Farnum, M. and Zukoski, C. (1990) Effect of glycerol on the interactions and solubility of bovine pancreatic trypsin inhibitor. *Biophys. J.*, **76**, 2716–2726.
- Fine, B.M., Lomakin, A., Ogun, O.O. and Benedek, G.B. (1996) The concentration dependence of the diffusion coefficient for bovine pancreatic trypsin inhibitor: a dynamic light scattering study of a small protein. *Biopolymers*, **28**, 2001–2024.
- George, A. and Wilson, W.W. (1994) Predicting protein crystallization from a dilute solution property. *Acta Crystallogr. D*, **50**, 361–365.
- George, A., Chiang, Y., Guo, B., Arabshahi, A., Cai, Z. and Wilson, W.E. (1997) Second virial coefficient as predictor in protein growth. *Methods Enzymol.*, **276**, 100–110.
- Gilliland, G.L., Tung, M. and Ladner, J.E. (2002) The Biological Macromolecule Crystallization Database: crystallization procedures and strategies. *Acta Crystallogr. D Biol. Crystallogr.*, **58**, 916–920.
- Haas, C. and Drenth, J. (1999) Understanding protein crystallization on the basis of the phase diagram. *J. Cryst. Growth*, **196**, 388–394.
- Haas, C. and Drenth, J. (2000) The interface between a protein crystal and an aqueous solution and its effects on nucleation and crystal growth. *J. Phys. Chem. B*, **104**, 368–377.
- Haas, C., Drenth, J. and Wilson, W.W. (1999) Relation between the solubility of proteins in aqueous solutions and the second virial coefficient of the solution. *J. Phys. Chem. B*, **103**, 2808–2811.
- Hansen, C.L., Skordalakes, E., Berger, J.M. and Quake, S.R. (2002) A robust and scalable microfluidic metering method that allows protein crystal growth by free interface diffusion. *Proc. Natl Acad. Sci., USA*, **99**, 16531–16536.
- Harris, L.J. and McPherson, A. (1995) Crystallization of intact monoclonal antibodies. *Prot. Struct. Funct. Genet.*, **23**, 285–289.
- Juarez-Martinez, G., Garza, C., Castillo, R. and Moreno, A. (2001) A dynamic light scattering investigation of the nucleation and growth of thaumatin crystals. *J. Cryst. Growth*, **232**, 119–131.
- Kantardjiev, K.A. and Rupp, B. (2003) Matthews coefficient probabilities: improved estimates for unit cell contents of proteins, DNA, and protein–nucleic acid complex crystals. *Protein Sci.*, **12**, 1865–1871.
- McPherson, A. (1982) *Preparation and Analysis of Protein Crystals*. Wiley, New York.
- Neal, B.L., Asthagiri, D. and Lenhoff, A.M. (1998) Molecular origins of osmotic second virial coefficients of proteins. *Biophys. J.*, **75**, 2469–2477.
- Peek, M.E. and Williams, L.D. (2001) X-ray crystallography of DNA–drug complexes. *Methods Enzymol.*, **340**, 282–290.
- Rosenbaum, D.F. and Zukoski, C.F. (1996) Protein interactions and crystallization. *J. Cryst. Growth*, **169**, 752–758.
- Rosenbaum, D.F., Zamora, P.C. and Zukoski, C.F. (1996) Phase behavior of small attractive colloidal particles. *Phys. Rev. Lett.*, **76**, 150–153.
- Rupp, B. (2003a) High-throughput crystallography at an affordable cost: the TB Structural Genomics Consortium Crystallization Facility. *Acc. Chem. Res.*, **36**, 173–181.
- Rupp, B. (2003b) Maximum-likelihood crystallization. *J. Struct. Biol.*, **142**, 162–169.

- Ruppert,S., Sandler,S.I. and Lenhoff,A.M. (2001) Correlation between the osmotic second virial coefficient and the solubility of proteins. *Biotechnol. Prog.*, **17**, 182–187.
- Santarsiero,B.D., Yegian,D.T., Lee,C.C., Spraggon,G., Gu,J., Scheibe,D., Uber,D.C., Cornell,E.W., Nordmeyer,R.A., Kolbe,W.S. *et al.* (2002) An approach to rapid protein crystallization using nanodroplets. *J. Appl. Crystallogr.*, **35**, 278–281.
- Sear,R.P. (2002) Distribution of the second virial coefficients of globular proteins. *Europhys. Lett.*, **60**, 938–944.
- Segelke,B.W. (2001) Efficiency analysis of sampling protocols used in protein crystallization screening. *J. Cryst. Growth*, **232**, 553–562.
- Segelke,B. and Rupp,B. (1998) Beyond the sparse matrix screen: a web service for randomly generating crystallization experiments. *Annual Meeting of the American Crystallographic Association*, Arlington, Virginia, USA, July 18–23.
- South East Collaboratory for Structural Genomics, Calculated pI values in *C.elegans* genome.
- Tardieu,A., Finet,S. and Bonnete,F. (2001) Structure of the macromolecular solutions that generate crystals. *J. Cryst. Growth*, **232**, 1–9.
- Urquhart,B.L., Cordwell,S.J. and Humphrey-Smith,I. (1998) Comparison of predicted and observed properties of proteins encoded in the genome of *Mycobacterium tuberculosis* H37Rv. *Biochem. Biophys. Resd. Commun.*, **253**, 70–79.
- Zimm,B.H. (1946) Applications of the methods of molecular distribution to solutions of large molecules. *J. Phys. Chem.*, **14**, 164–179.