

Protein Structure Analysis

Iosif Vaisman

2023

Protein folding

DEVELOPMENTAL BIOLOGY SUPPLEMENT 2, 1-20 (1968)

I. SELF-ASSEMBLY OF MACROMOLECULAR STRUCTURES
Spontaneous Formation of the Three-Dimensional Structure of Proteins

CHRISTIAN B. ANFINSEN
Laboratory of Chemical Biology, National Institute of Arthritis and Metabolic Diseases, National Institutes of Health, Bethesda, Maryland

INTRODUCTION

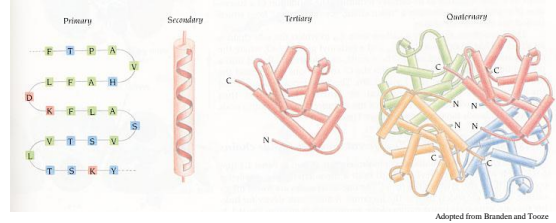
Our major consideration in this symposium will be the emergence of order during cellular differentiation and growth. The concept "emerging order" implies an organized, genetically complex process taking place over a reasonably extended stretch of time. In contrast, the restoration of linear genetic information in the form of three-dimensional protein structure results from a rapid and spontaneous interaction of amino acid side chains with each other, with the completed polypeptide backbone, and with the environment, without the necessity for additional genetic information (Anfinsen, 1967; Epstein *et al.*, 1963). The achievement of this unique geometry might be visualized as a rather hiter-skeleton process. An almost infinite number of sets of interactions are possible as an extended polypeptide chain coils upon itself (Fig. 1). If the process of folding involved even a small fraction of this number of conformational states, the specific folding of the chain could clearly require considerable time. It is probable that the rapidity of folding is made possible through the formation of one or more "nucleation sites" by side chain interactions that would predispose, during subsequent interactions, to the tertiary struc-



TABLE 1
 THE NUMBER OF WAYS IN WHICH 26 RESIDUES CAN COMBINE
 TO FORM 2-DIMENSIONAL SEQUENCES

Number of residues	Number of combinations
1	1
2	25
3	156
4	1000
5	9801
6	157480
7	1,679,616
8	28,247,524
9	384,901,655
10	6,471,921,615
11	117,493,100,125
12	2,182,874,848,125
13	39,905,328,861,250
14	713,006,805,781,250
15	12,918,233,362,207,500
16	239,926,782,523,207,500
17	4,328,207,752,828,500,000
18	80,344,803,607,718,500,000
19	1,428,743,281,281,500,000,000
20	26,308,967,727,727,187,500,000,000
21	471,113,076,128,799,684,406,250,000,000
22	8,506,308,068,630,689,684,406,250,000,000,000
23	153,271,322,282,742,627,090,125,000,000,000,000
24	2,755,641,927,348,122,396,748,125,000,000,000,000,000
25	49,494,144,047,272,318,555,67,800,000,000,000,000,000,000

Protein Structure Hierarchy



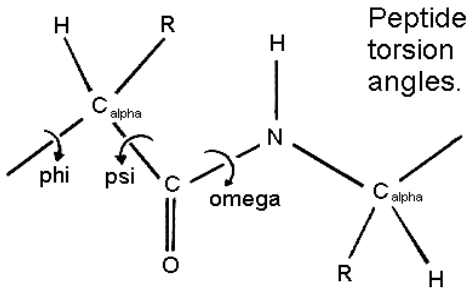
- Primary - the sequence of amino acid residues
- Secondary - ordered regions of primary sequence (helices, beta-sheets, turns)
- Tertiary - the three-dimensional fold of a protein subunit
- Quaternary - the arrangement of subunits in oligomers.

Anfinsen's Dogma

Three-dimensional structure of a protein is determined solely by its amino-acid sequence.

Native conformation of the protein is the global-minimum free energy conformation.

Levinthal paradox



3 conformations per residue is a very conservative estimate

Complexity of protein structure (Levinthal paradox)

100 residue protein
 3 conformations per residue

number of distinct conformations:
 $3^{100} \approx 10^{48}$

sampling time $\approx 10^{30}$ years

Complexity

P (Polynomial)

complexity class of decision problems for which execution time of a computation is no more than a polynomial function of the problem size

NP (Nondeterministic Polynomial)

complexity class of decision problems for which answers can be checked by an algorithm whose run time is polynomial in the size of the input

Protein Folding Problem

Given: **sequence**

Find: **structure**

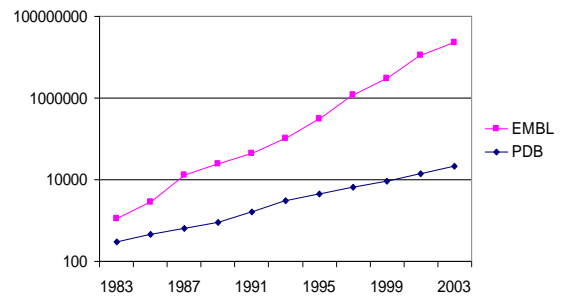
The problem is NP-complete

Protein Folding Problem

Problem for us, not for proteins.
They just fold...

(Ken Dill)

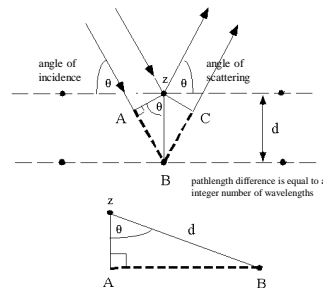
Dynamics of Database Growth



Protein Structure Determination

- X-ray crystallography
- NMR spectroscopy
- Neutron diffraction
- Electron microscopy
- Atomic force microscopy

X-ray crystallography

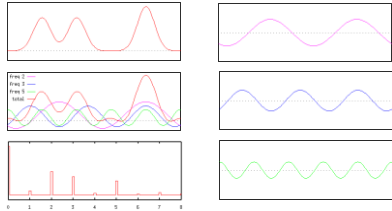


Bragg's Law

$$n\lambda = 2d \sin\theta$$

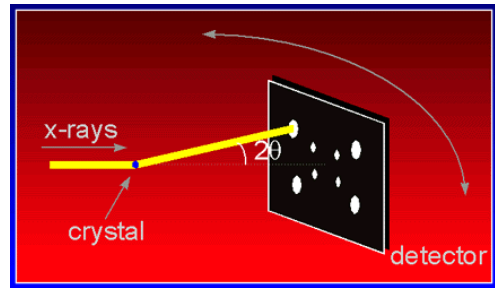
X-ray crystallography

Phase determination: MIR and MAD
(Multiple Isomorphous Replacement and Multiwavelength Anomalous Diffraction)

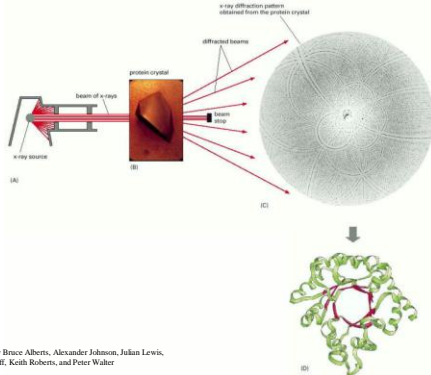


Fourier Transforms

X-ray crystallography

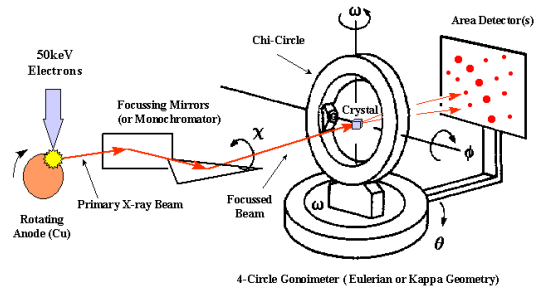


X-ray crystallography



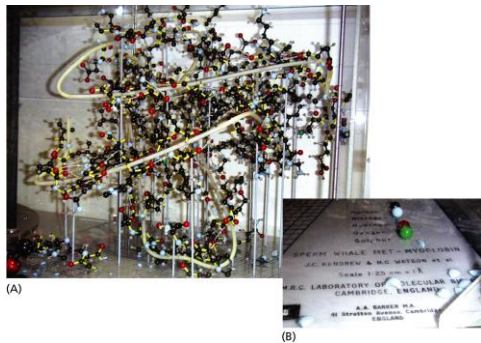
© 2002 by Bruce Alberts, Alexander Johnson, Julian Lewis, Martin Raff, Keith Roberts, and Peter Walter

X-ray crystallography



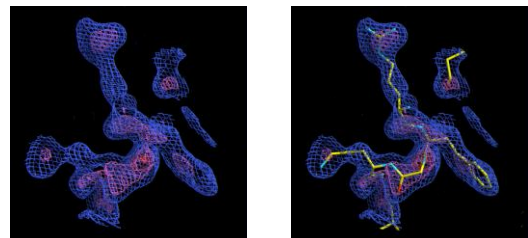
4-Circle Goniometer (Eulerian or Kappa Geometry)

X-ray crystallography



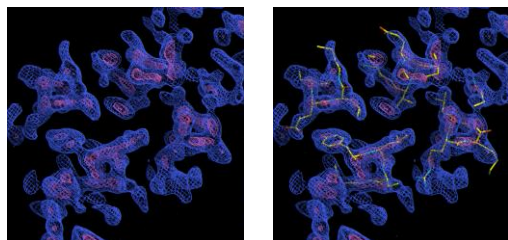
Adapted from Zeebii, Baum, 2008

X-ray crystallography



Electron density map created from multi-wavelength data (Arg)

X-ray crystallography



Experimental electron density map and model fitting
(apoE four helix bundle)

X-ray crystallography

Confidence in structural features of proteins determined by X-ray crystallography

(These are rough estimates, and depend strongly on the quality of the data.)

Structural feature	Resolution				
	5 Å	3 Å	2.5 Å	2.0 Å	1.5 Å
Chain tracing	—	Fair	Good	Good	Good
Secondary structure	Helices fair	Fair	Good	Good	Good
Sidechain conformations	—	—	Fair	Good	Good
Orientation of peptide planes	—	—	Fair	Good	Good
Protein hydrogen atoms visible	—	—	—	—	Good

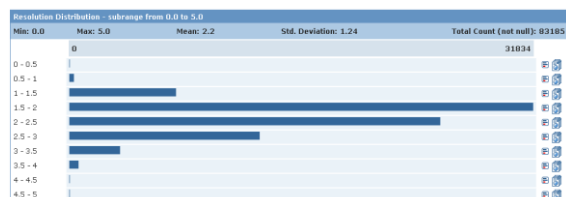
wwPDB statistics

Year	Total Depositions	Deposited To			Processed By		
		RCSB PDB	PDBj	PDBe	RCSB PDB	PDBj	PDBe
2001	3287	2673	118	496	2408	383	496
2002	3565	2769	289	507	2401	657	507
2003	4830	3488	673	669	3135	1026	669
2004	5508	3796	900	812	3082	1614	812
2005	6678	4507	1166	1005	3563	2110	1005
2006	7282	5145	1052	1085	4252	1945	1085
2007	8130	5399	1603	1128	4703	2299	1128
2008	7073	5452	648	973	4106	1994	973
2009	8300	6715	527	1058	5069	2173	1058
2010	8878	6912	593	1373	5464	2041	1373
2011	9250	7172	582	1496	5938	1816	1496
2012	9972	7695	601	1676	6408	1888	1676
2013	10566	8031	749	1786	6652	2128	1786
2014	10364	8178	501	1685	6040	1779	2545
2015	8070	6880	49	1141	3692	1411	2968
TOTAL	114736	87257	10061	17418	69210	25422	20105

PDB statistics

Exp.Method	Proteins	Nucleic Acids	Protein/NA Complexes	Other	Total
X-RAY	93956	1668	4692	4	100320
NMR	9751	1130	227	8	11116
ELECTRON MICROSCOPY	619	29	204	0	852
HYBRID	76	3	2	1	82
other	168	4	6	13	191
Total	104570	2834	5131	26	112561

PDB resolutions



PDB redundancy

Method	Description	# of Clusters
blast	95% identity (one chain)	17575
blast	90% identity (one chain)	16853
blast	70% identity (one chain)	15114
blast	50% identity (one chain)	12886
blast	40% identity (one chain)	11218
blast	30% identity (one chain)	9294

PDB ambiguities

Table 1 The number of PDB structures retrieved by ambiguous chemical component codes

Code	Name	Number of PDB structures ^a
SUL	Sulfate anion	156 (3.6%)
SO4	Sulfate ion	4083 (96.4%)
SUL and SO4	Sulfate anion and sulfate ion	1 (0.03%)
NET	Tetraethylammonium ion	9 (90%)
E4N	Tetraethylammonium ion	1 (10%)
MMC	Methyl mercury ion	8 (66.66%)
HGC	Methyl mercury ion	4 (33.33%)

^aPercentages of the total number of structures with the chemical component are shown in brackets. Search carried out August 2006.