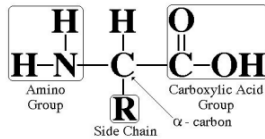
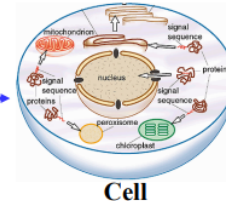
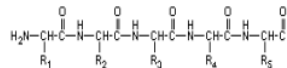


Protein Structure Modeling

Amino Acid Structure



AGCWY.....



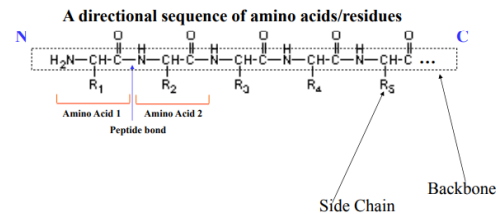
Protein Sequence – Primary Structure

- The first protein was sequenced by Frederick Sanger in 1953.
- Twice Nobel Laureate (1958, 1980) (other: Curie, Pauling, Bardeen).
- Determined the amino acid sequence of insulin and proved proteins have specific primary structure.



GIVEQCASVCSLYQLENYCN A chain (21 amino acids)
FVNQHLCGSHLVEALYLVCGERGFFYTPKA B chain (30 amino acids)

Protein Sequence



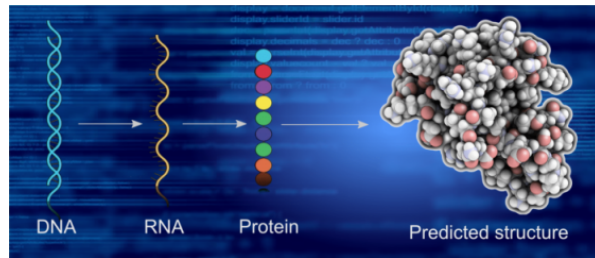
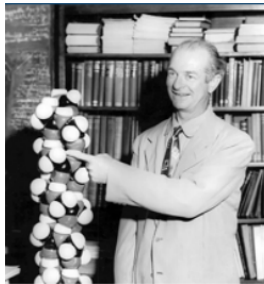
Amino Acids

Amino acid	Abbrev.	Side chain	Hydrophobic	Polar	Charged	Small	Tiny	Aromatic or Aliphatic	van der Waals volume	Codon	Occurrence in proteins (%)
Alanine	Ala, A	-CH ₃	X	-	-	X	X	-	67	GCU, GCC, GCA, GCG	7.8
Cysteine	Cys, C	-CH ₂ SH	X	-	-	X	-	-	86	UGU, UGC	1.9
Aspartate	Asp, D	-CH ₂ COOH	-	X	negative	X	-	-	91	GAU, GAC	5.3
Glutamate	Glu, E	-CH ₂ CH ₂ COOH	-	X	negative	-	-	-	109	GAA, GAG	6.3
Phenylalanine	Phe, F	-CH ₂ C ₆ H ₅	X	-	-	-	-	Aromatic	135	UUU, UUC	3.9
Glycine	Gly, G	-H	X	-	-	X	X	-	48	GGU, GGC, GGA, GGG	7.2
Histidine	His, H	-CH ₂ -C ₃ H ₃ N ₂	-	X	positive	-	-	Aromatic	118	CAU, CAC	2.3
Isoleucine	Ile, I	-CH(CH ₃)CH ₂ CH ₃	X	-	-	-	-	Aliphatic	124	AUU, AUC, AUA	6.3
Lysine	Lys, K	-(CH ₂) ₄ NH ₂	-	X	positive	-	-	-	135	AAA, AAG	5.9
Leucine	Leu, L	-CH ₂ CH(CH ₃) ₂	X	-	-	-	-	Aliphatic	124	UUA, UUG, CUU, CUC, CUA, CUG	9.1
Methionine	Met, M	-CH ₂ CH ₂ SCH ₃	X	-	-	-	-	-	124	AUG	2.3
Asparagine	Asn, N	-CH ₂ CONH ₂	-	X	-	X	-	-	96	AAU, AAC	4.3
Proline	Pro, P	-CH ₂ CH ₂ CH ₂ -	X	-	-	X	-	-	90	CCU, CCC, CCA, CCG	6.2
Glutamine	Gln, Q	-CH ₂ CH ₂ CONH ₂	-	X	-	-	-	-	114	CAA, CAG	4.2
Arginine	Arg, R	-(CH ₂) ₃ NH-C(NH ₂)NH ₂	-	X	positive	-	-	-	148	CGU, CGC, CGA, CGG, AGA, AGG	5.1
Serine	Ser, S	-CH ₂ OH	-	X	-	X	X	-	73	UCU, UCC, UCA, UCG, AGU, AGC	6.8
Threonine	Thr, T	-CH(OH)CH ₃	X	X	-	X	-	-	93	ACU, ACC, ACA, ACG	5.9
Valine	Val, V	-CH(CH ₃) ₂	X	-	-	X	-	Aliphatic	105	GUU, GUC, GUA, GUG	6.6
Tryptophan	Trp, W	-CH ₂ C ₈ H ₆ N	X	-	-	-	-	Aromatic	163	UGG	1.4
Tyrosine	Tyr, Y	-CH ₂ -C ₆ H ₄ OH	X	X	-	-	-	Aromatic	141	UAU, UAC	3.2

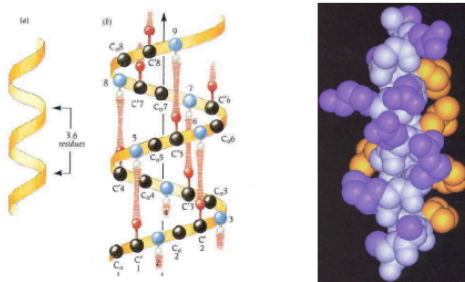
Hydrophilic

Protein Secondary Structure

- Determined by hydrogen bond patterns
- 3-Class categories: alpha-helix, beta-sheet, loop (or coil)
- First deduced by Linus Pauling et al.

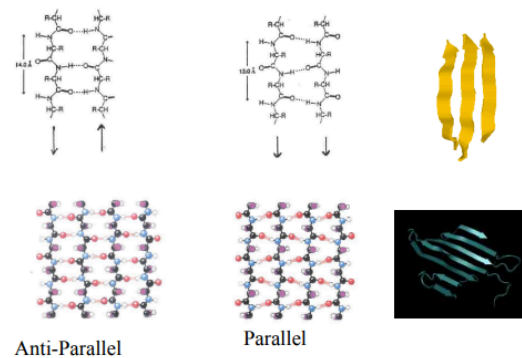


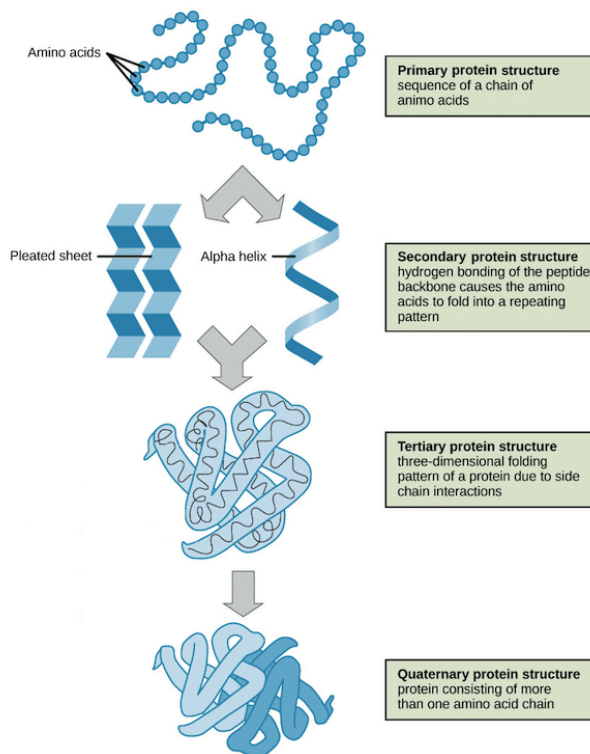
Alpha-Helix



Jurnak, 2003

Beta-Sheet





Protein Extraction Methods:

Proteins are essential macromolecules that perform a variety of functions in the body, like DNA replication, catalyzing reactions, and providing structural support. They are studied in three main ways:

1. **In Vivo:** Studying proteins within the organism to understand how they interact.
2. **In Vitro:** Studying purified proteins in controlled lab settings to avoid interference from other factors.
3. **In Silico:** Using computer simulations to study proteins, saving time and resources.

Proteins are classified into types like **extracellular matrix proteins** (e.g., elastin, collagen) and **globular proteins** (e.g., enzymes, antibodies). Purifying proteins from other cellular components is crucial for research, whether for large-scale production (e.g., insulin) or analysis of small protein amounts.

Uses of Isolated Proteins:

Protein extraction is widely applied in both research and industry. The purification process requires multiple steps and detection methods, including absorbance, spectrometry, and antibody-based techniques.

In **clinical applications**, isolated proteins can help diagnose diseases like diabetes or be used in treatments (e.g., collagen in skincare). In **research**, purified proteins enable various studies:

- **Immunoprecipitation (IP):** Isolates proteins using antibodies.
- **Proteomics:** Studies the entire set of proteins in an organism.
- **Enzyme Assays:** Measure enzyme activity, including different experiment types like relaxation or transient kinetics.
- **Western Blot:** Detects specific proteins in a sample.
- **Gel Electrophoresis:** Separates proteins by size and charge.
- **Biomarkers:** Used to track biological processes like disease or treatment effects.

These techniques are crucial for advancing both clinical diagnostics and scientific understanding.

What is Protein Sequencing?

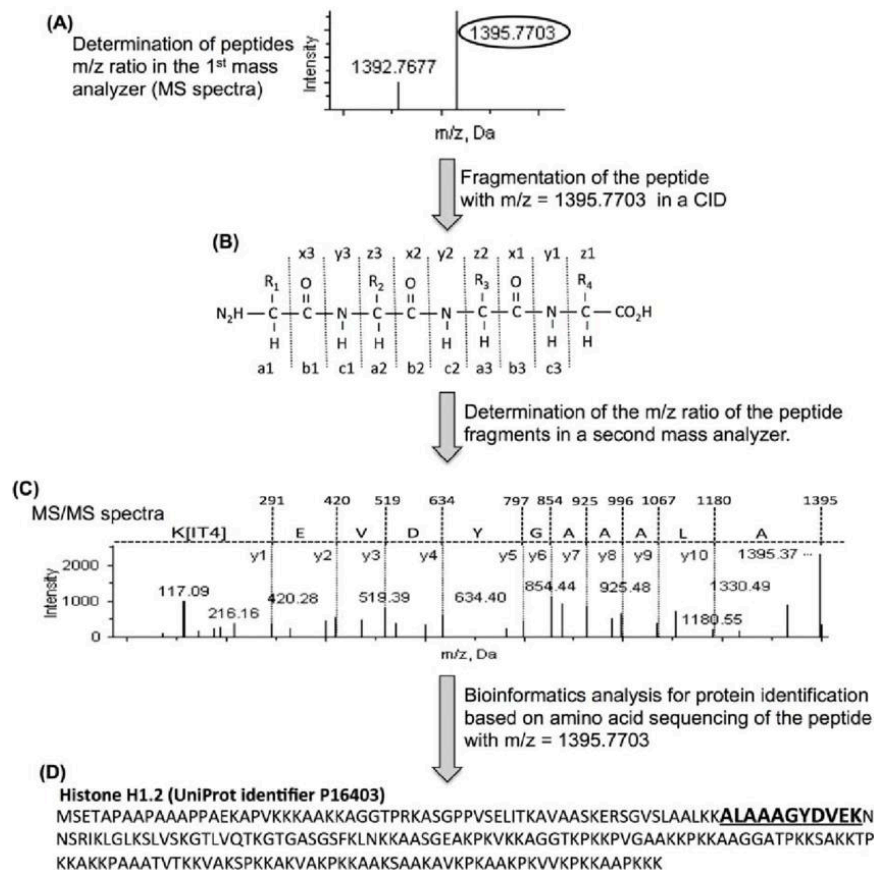
Protein sequencing is a technique used to determine the exact order of amino acids in a protein. This sequence is known as the protein's primary structure, and understanding it helps us know how a protein functions in the body.

Significance of Protein Sequencing:

1. **Decoding Genetic Information:** It translates DNA's genetic code into proteins, revealing the exact sequence of amino acids in a protein.
2. **Unveiling Protein Function:** The sequence affects how a protein folds and interacts with other molecules, helping scientists understand its role.
3. **Biotechnology:** It aids in designing drugs, enzymes, and other proteins for medical and industrial uses.
4. **Personalized Medicine:** By sequencing proteins, doctors can identify genetic mutations linked to diseases and create personalized treatments.
5. **Structural Biology:** Protein sequencing is essential for understanding the three-dimensional shape of proteins, which is key to drug design.
6. **Proteomics Advancements:** It helps in studying all proteins in an organism, uncovering complex biological processes.

Methods of Protein Sequencing:

- **Edman Degradation:** This older method removes and identifies amino acids from the protein's starting point, but requires large samples.
- **Mass Spectrometry:** It analyzes protein fragments to determine their sequence and is useful for complex mixtures. Mass Spectrometry (MS) measures the mass-to-charge ratio of ions to identify and quantify molecules.



- Next-Generation Sequencing (NGS): This modern technique uses mRNA to indirectly infer protein sequences, offering high throughput.

Applications of Protein Sequencing:

- Drug Development: Helps create targeted therapies by identifying protein structures linked to diseases.
- Structural Biology: Reveals a protein's 3D shape, crucial for designing drugs that target specific protein structures.
- Proteomics Research: Identifies proteins in biological samples, providing insights into diseases like cancer and neurological disorders.
- Biotechnology: Protein sequencing is key in designing and producing therapeutic proteins and enzymes.
- Personalized Medicine: It helps tailor treatments by identifying genetic variations in proteins.

Challenges and Solutions:

- **Sample Preparation:** Extracting and purifying proteins can be difficult, but new techniques are making it easier.
- **Data Analysis:** The large amount of data requires advanced software to interpret protein sequences accurately.
- **Cost:** Protein sequencing can be expensive, but efforts are ongoing to reduce costs and improve efficiency.

Technological Advancements:

- **High-Throughput Sequencing:** Allows simultaneous analysis of many proteins, speeding up research.
- **Mass Spectrometry Improvements:** Enhanced sensitivity helps detect low-abundance proteins and modifications.
- **Hybrid Approaches:** Combining different sequencing methods improves accuracy and efficiency.

Anfinsen's Folding Experiment

- Structure is uniquely determined by protein sequence
- Protein function is determined by protein structure



Protein Structure Determination

- X-ray crystallography
- Nuclear Magnetic Resonance (NMR) Spectroscopy
- X-ray: any size, accurate (1-3 Angstrom (10^{-10} m)), sometime hard to grow crystal
- NMR: small to medium size, moderate accuracy, structure in solution

Historically, determining protein structures was a slow and laborious process, but advances in technology, including automation and better equipment, have sped up this process. The **Protein Data Bank now holds over 206,000 protein structures, and technological improvements have led to a "resolution revolution,"** especially in cryo-EM.

X-ray crystallography is a key method for determining protein structures, offering detailed data on atomic arrangement. It works by passing X-rays through a rotating protein crystal and analyzing the diffracted rays. The technique provides high resolution but requires high-quality crystals and large amounts of protein, often produced through recombinant methods. It is particularly effective for rigid proteins but less so for flexible ones.

Electron Microscopy (EM) and Cryo-EM: These methods are ideal for studying large macromolecules and cellular structures, avoiding the need for protein crystallization. Cryo-EM, performed at very low temperatures, provides high-resolution images and works with small amounts of protein, reducing radiation damage.

Nuclear Magnetic Resonance (NMR) Spectroscopy: This technique uses radiofrequency waves to analyze protein atoms. It requires larger quantities of stable protein at room temperature and works best for small proteins, offering high resolution, especially for flexible proteins.

Small-Angle X-Ray Scattering (SAXS) and Neutron Scattering: These methods are useful for studying protein structures in solution when high resolution isn't necessary, allowing better control over experimental conditions.

Homology Modeling: This technique creates a 3D protein model based on a known similar protein, relying on sequence similarities.

Partial Structural Study Methods: These include ultracentrifugation, mass spectrometry, and fluorescence spectrometry, often used alongside other techniques to gain further insights into protein structure.

<https://pmc.ncbi.nlm.nih.gov/articles/PMC2705668/>

<https://sci-hub.se/https://doi.org/10.1002/9780470015902.a0002716.pub2>

In the 1930s, William Astbury studied diffraction from biological fibers, while Dorothy Crowfoot (Hodgkin) and J.D. Bernal explored crystals of macromolecules. Over 20 years later, in 1958, John Kendrew and his team solved the first crystal structure of myoglobin from sperm whale muscle. This was followed by Max Perutz's work on haemoglobin (1962) and David Phillips on lysozyme (1965).

Meanwhile, in the 1930s, Rabi and colleagues demonstrated nuclear resonance (NMR) by applying electromagnetic radiation to molecular beams. Though NMR was theoretically possible in solids and liquids, early attempts failed due to low sensitivity and long relaxation times of the nuclei. Advances in electronics led to the first NMR spectra in condensed phases by groups led by Bloch, Pound, Purcell, and others in the mid-1940s. The chemical shift effect was observed by Knight (1949), Proctor and Yu (1950-1951), and Dickinson (1950), and spin-spin coupling was discovered by Hahn and Maxwell (1951).

Fourier transform (FT) techniques introduced by Ernst and Anderson in 1966 paved the way for two-dimensional NMR experiments, transforming NMR's role in biological systems. The first proton NMR spectrum for a protein was recorded in 1957 (Saunders et al.), and since then, NMR methods have advanced to enable routine assignments of proton resonances in proteins up to 50 kDa. One key discovery was the Overhauser Effect (OE) by Overhauser in 1953, which improved signal-to-noise ratios and led to the development of NOESY experiments used in protein structure determination. In 1985, Wüthrich and colleagues reported the first complete NMR structure of a globular protein (Williamson et al., 1985).

As of January 2012, the Protein Data Bank (PDB) contained about 80,000 entries, with 70,000 from X-ray crystallography and 10,000 from NMR spectroscopy. NMR is primarily used for smaller proteins (under 50 kDa), while X-ray crystallography is used for larger proteins (over 35 kDa).

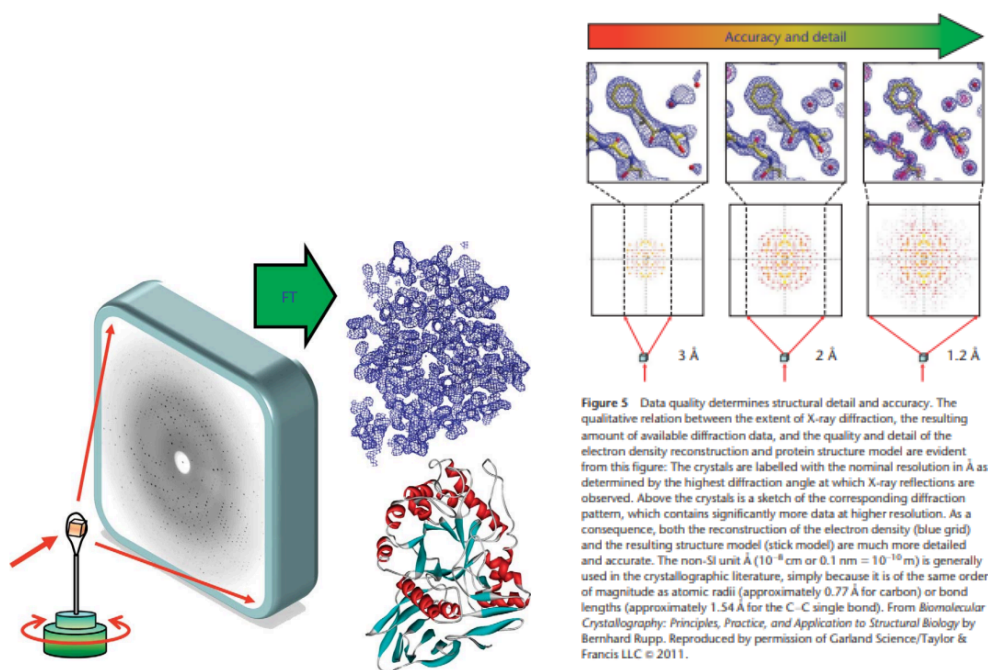
X-ray Diffraction (XRD) and NMR Spectroscopy:

X-ray diffraction (XRD) and nuclear magnetic resonance (NMR) spectroscopy use different physical processes to determine the 3D structure of macromolecules. Both techniques involve electromagnetic radiation, but with different energy levels.

In XRD, high-energy X-rays are scattered by the electrons in a protein crystal (about 100–10 microns in size). This scattering creates a pattern of diffraction spots, which are recorded and used to map the electron density of the molecules in the crystal. This map helps build an atomic model of the structure. However, the phases of the diffraction spots are lost during detection, so separate methods are needed to recover them for accurate reconstruction.

To solve this phase problem, two main methods are used: Transforming Data into a Model:

1. **Molecular Replacement:** If a similar protein structure is already known, it can be used to estimate the initial phases. This method works quickly, but the initial structure may have biases from the model used. Around 75% of X-ray structures in the Protein Data Bank (PDB) are solved this way.



2. **Experimental Phasing:** If no similar structure is available, new phases must be determined through experiments. This often involves introducing heavy atoms into the crystal to create differences in the diffraction pattern (called isomorphous replacement or anomalous scattering). These differences help determine the phases. For example, replacing methionine with seleno-methionine in proteins can provide a useful anomalous signal for phasing. Around 25% of structures in the PDB are solved using experimental phasing.

After obtaining the initial phases, density modification techniques are used to refine the electron density map and improve the protein model.

Synchrotrons, which provide high-intensity X-ray beams, are commonly used for data collection, and their adjustable wavelengths are especially helpful for anomalous data collection.

X-ray Structure Determination:

In X-ray crystallography, a crystal is placed on a rotating goniostat and exposed to an intense X-ray beam (5–20 keV). The crystal is rotated in small increments, and diffraction images are captured for each position. These images don't directly show the molecule but provide data that can be used to reconstruct its structure.

Using a mathematical method called Fourier transform (FT), along with phase data for each diffraction spot, the electron density of the molecules in the crystal is reconstructed. Finally, an atomic model is built based on this electron density map.

Nuclear Magnetic Resonance (NMR) Overview:

NMR is based on the magnetic properties of atomic nuclei. Nuclei with nonzero spin have a magnetic dipole that moves in response to an external magnetic field. When a sample is placed in a strong magnetic field and exposed to radiofrequency (RF) radiation, the nuclei absorb energy at specific frequencies. The absorption depends on the magnetic properties of the nuclei.

In biomolecular NMR, the focus is on simple nuclei like ^1H , ^{15}N , ^{13}C , and ^{31}P , because they produce clear spectra. NMR works by applying RF pulses that disturb the nuclear spin states. The response to these pulses (a signal called the Free Induction Decay, or FID) is collected, and then Fourier transformed into an NMR spectrum.

The NMR spectrum reveals details about the chemical environment of nuclei (via "chemical shifts") and the interactions between nearby nuclei (called "coupling constants"). These interactions provide insights into the structure of the molecule.

For complex molecules like proteins, one-dimensional NMR spectra often show overlapping peaks. Therefore, multidimensional NMR techniques are used to spread the data across multiple dimensions (e.g., 2D NMR or 3D NMR) to separate the peaks and make interpretation easier.

NMR Structure Determination:

NMR is different from X-ray crystallography because it can study molecules in solution without requiring them to be in a perfect crystal form. NMR provides information not only on the structure but also on the dynamics of the molecule.

For proteins, NMR requires isotopic labeling (like ^{15}N or ^{13}C) to improve resolution, especially for larger proteins. The chemical shifts and coupling constants help assign specific resonances to atoms in the molecule. NOESY experiments are then used to determine distances between protons, further refining the structure.

Resonance Assignment:

Assigning the correct resonance to each proton is one of the most time-consuming parts of NMR structure determination. Special experiments, like CBCA(CO)NH, help assign the backbone atoms (such as the N, $\text{C}\alpha$, and $\text{C}\beta$ atoms) in the protein chain. These experiments build a sequential assignment, linking each residue to the next one. For larger proteins, additional experiments are used to overcome relaxation effects and improve the assignment process.

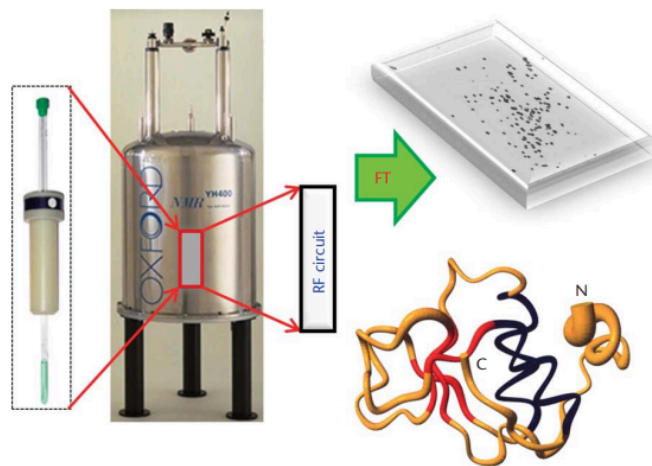


Figure 2 Experimental basis of NMR structure determination. Biomolecules in solution at close to physiological conditions are inserted into a magnet. The radio frequency circuit detects the time domain signal corresponding response of the nuclear spins to resonance. This analogue time domain signal detected by the circuit is amplified and digitised prior to Fourier transformation into the spectral domain. A combination of 2D and 3D experiments are generally collected, processed and analysed to obtain NMR restraint parameter that are sensitive to determine both local structural relations and events (through chemical shifts and coupling constants) as well as the global fold (via NOEs) of a protein. Three-dimensional structural models generated by NMR methods also carry additional information on residue-specific dynamic motion.

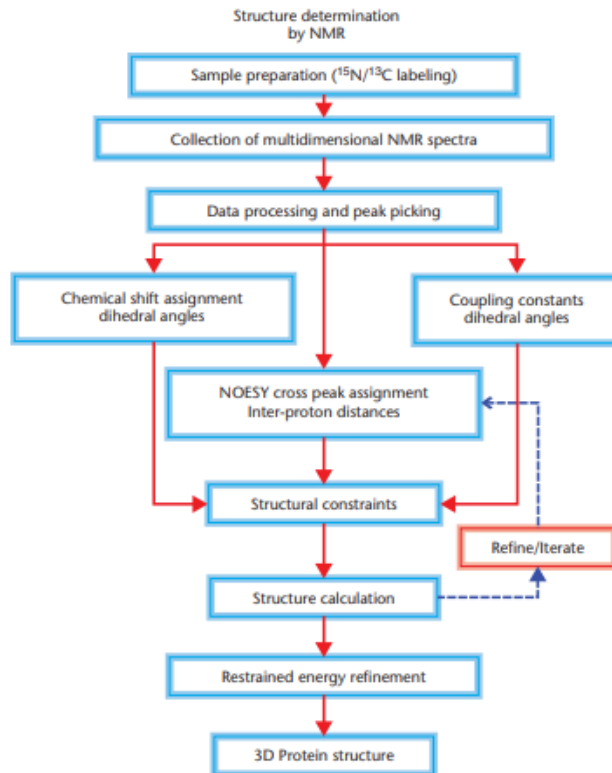
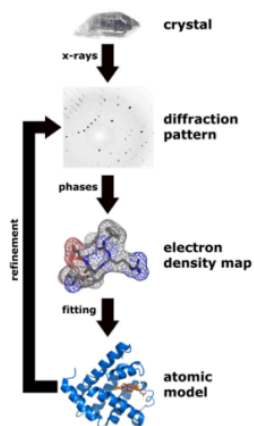


Figure 3 Flowchart describing the major steps in NMR based structure determination. Four principal elements are combined in the NMR method for protein structure determination: (1) multidimensional NMR experiments that provide the data for all the following steps; (2) sequence-specific resonance assignment – matching each proton in the protein to respective peaks in the spectra; (3) the Nuclear Overhauser Effect (NOE) data – providing inter-proton distances to identify the global fold of the protein; (4) computational tools such as distance geometry (restraint) based approach for the structural interpretation of the NMR data and the evaluation of the resulting molecular structures and Each of these elements is critically important in obtaining a good quality NMR structure.

Distance Geometry Restraints:

Once resonance assignments are complete, distance constraints are needed to model the protein structure. NOESY experiments provide data about the proximity of protons, giving distance restraints. These experiments are especially useful for measuring distances between neighboring atoms (less than ~0.6 nm). The information from these experiments, along with known bond lengths and angles, is used to generate a 3D model of the protein structure.



Pacific Northwest National Laboratory's high magnetic field (800 MHz 18.8 T) NMR spectrometer being loaded with a sample.

Storage in Protein Data Bank

Welcome to the RCSB PDB

The RCSB PDB provides a variety of tools and resources for studying the structures of biological macromolecules and their relationships to sequence, function, and disease.

The RCSB is a member of the [wwPDB](#) whose mission is to ensure that the PDB archive remains an international resource with uniform data.

This site offers tools for browsing, searching, and reporting that utilize the data resulting from ongoing efforts to create a more consistent and comprehensive archive.

Information about computational resources can be found [here](#).

A narrated tutorial illustrates how to search, navigate, browse, generate reports and visualize structures using the new site. (This requires the [Java plug-in](#).)

Comments? [info@rcsb.org](#)

Materials of the Month: RNA-Proteins

How would you make a protein-coding machine that would be safe to use inside a cell? Optimize proteins for function and improve their small and efficient effect on proteins and other cells. This would mean work in a cell. The cell has to be able to read the code in the cell and translate it into a protein. This is the challenge.

Search database

Crystal structure of putative lipase from the G-D-L family from *Bacillus* sp. ATCC 25411

Authors: Joint Center for Structural Genomics (JCSG)

Primary: Joint Center for Structural Genomics (JCSG)

History: Deposited: 2004-02-19, Released: 2004-03-16

Representative images: Type: JCSG DIFFRACTION, Data: [JCSG]

Parameters: Resolution: 2.01 Å, R-value: 0.175 (R_{free}: 0.218), Space Group: P₂, 2₁, 2₁

EMR Data: Chain ID: 1, 2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12, 13, 14, 15, 16, 17, 18, 19, 20, 21, 22, 23, 24, 25, 26, 27, 28, 29, 30, 31, 32, 33, 34, 35, 36, 37, 38, 39, 40, 41, 42, 43, 44, 45, 46, 47, 48, 49, 50, 51, 52, 53, 54, 55, 56, 57, 58, 59, 60, 61, 62, 63, 64, 65, 66, 67, 68, 69, 70, 71, 72, 73, 74, 75, 76, 77, 78, 79, 80, 81, 82, 83, 84, 85, 86, 87, 88, 89, 90, 91, 92, 93, 94, 95, 96, 97, 98, 99, 100, 101, 102, 103, 104, 105, 106, 107, 108, 109, 110, 111, 112, 113, 114, 115, 116, 117, 118, 119, 120, 121, 122, 123, 124, 125, 126, 127, 128, 129, 130, 131, 132, 133, 134, 135, 136, 137, 138, 139, 140, 141, 142, 143, 144, 145, 146, 147, 148, 149, 150, 151, 152, 153, 154, 155, 156, 157, 158, 159, 160, 161, 162, 163, 164, 165, 166, 167, 168, 169, 170, 171, 172, 173, 174, 175, 176, 177, 178, 179, 180, 181, 182, 183, 184, 185, 186, 187, 188, 189, 190, 191, 192, 193, 194, 195, 196, 197, 198, 199, 200, 201, 202, 203, 204, 205, 206, 207, 208, 209, 210, 211, 212, 213, 214, 215, 216, 217, 218, 219, 220, 221, 222, 223, 224, 225, 226, 227, 228, 229, 230, 231, 232, 233, 234, 235, 236, 237, 238, 239, 240, 241, 242, 243, 244, 245, 246, 247, 248, 249, 250, 251, 252, 253, 254, 255, 256, 257, 258, 259, 260, 261, 262, 263, 264, 265, 266, 267, 268, 269, 270, 271, 272, 273, 274, 275, 276, 277, 278, 279, 280, 281, 282, 283, 284, 285, 286, 287, 288, 289, 290, 291, 292, 293, 294, 295, 296, 297, 298, 299, 300, 301, 302, 303, 304, 305, 306, 307, 308, 309, 310, 311, 312, 313, 314, 315, 316, 317, 318, 319, 320, 321, 322, 323, 324, 325, 326, 327, 328, 329, 330, 331, 332, 333, 334, 335, 336, 337, 338, 339, 340, 341, 342, 343, 344, 345, 346, 347, 348, 349, 350, 351, 352, 353, 354, 355, 356, 357, 358, 359, 360, 361, 362, 363, 364, 365, 366, 367, 368, 369, 370, 371, 372, 373, 374, 375, 376, 377, 378, 379, 380, 381, 382, 383, 384, 385, 386, 387, 388, 389, 390, 391, 392, 393, 394, 395, 396, 397, 398, 399, 400, 401, 402, 403, 404, 405, 406, 407, 408, 409, 410, 411, 412, 413, 414, 415, 416, 417, 418, 419, 420, 421, 422, 423, 424, 425, 426, 427, 428, 429, 430, 431, 432, 433, 434, 435, 436, 437, 438, 439, 440, 441, 442, 443, 444, 445, 446, 447, 448, 449, 450, 451, 452, 453, 454, 455, 456, 457, 458, 459, 460, 461, 462, 463, 464, 465, 466, 467, 468, 469, 470, 471, 472, 473, 474, 475, 476, 477, 478, 479, 480, 481, 482, 483, 484, 485, 486, 487, 488, 489, 490, 491, 492, 493, 494, 495, 496, 497, 498, 499, 500, 501, 502, 503, 504, 505, 506, 507, 508, 509, 510, 511, 512, 513, 514, 515, 516, 517, 518, 519, 520, 521, 522, 523, 524, 525, 526, 527, 528, 529, 530, 531, 532, 533, 534, 535, 536, 537, 538, 539, 540, 541, 542, 543, 544, 545, 546, 547, 548, 549, 550, 551, 552, 553, 554, 555, 556, 557, 558, 559, 560, 561, 562, 563, 564, 565, 566, 567, 568, 569, 570, 571, 572, 573, 574, 575, 576, 577, 578, 579, 580, 581, 582, 583, 584, 585, 586, 587, 588, 589, 590, 591, 592, 593, 594, 595, 596, 597, 598, 599, 600, 601, 602, 603, 604, 605, 606, 607, 608, 609, 610, 611, 612, 613, 614, 615, 616, 617, 618, 619, 620, 621, 622, 623, 624, 625, 626, 627, 628, 629, 630, 631, 632, 633, 634, 635, 636, 637, 638, 639, 640, 641, 642, 643, 644, 645, 646, 647, 648, 649, 650, 651, 652, 653, 654, 655, 656, 657, 658, 659, 660, 661, 662, 663, 664, 665, 666, 667, 668, 669, 670, 671, 672, 673, 674, 675, 676, 677, 678, 679, 680, 681, 682, 683, 684, 685, 686, 687, 688, 689, 690, 691, 692, 693, 694, 695, 696, 697, 698, 699, 700, 701, 702, 703, 704, 705, 706, 707, 708, 709, 710, 711, 712, 713, 714, 715, 716, 717, 718, 719, 720, 721, 722, 723, 724, 725, 726, 727, 728, 729, 730, 731, 732, 733, 734, 735, 736, 737, 738, 739, 740, 741, 742, 743, 744, 745, 746, 747, 748, 749, 750, 751, 752, 753, 754, 755, 756, 757, 758, 759, 760, 761, 762, 763, 764, 765, 766, 767, 768, 769, 770, 771, 772, 773, 774, 775, 776, 777, 778, 779, 780, 781, 782, 783, 784, 785, 786, 787, 788, 789, 790, 791, 792, 793, 794, 795, 796, 797, 798, 799, 800, 801, 802, 803, 804, 805, 806, 807, 808, 809, 810, 811, 812, 813, 814, 815, 816, 817, 818, 819, 820, 821, 822, 823, 824, 825, 826, 827, 828, 829, 830, 831, 832, 833, 834, 835, 836, 837, 838, 839, 840, 841, 842, 843, 844, 845, 846, 847, 848, 849, 850, 851, 852, 853, 854, 855, 856, 857, 858, 859, 860, 861, 862, 863, 864, 865, 866, 867, 868, 869, 870, 871, 872, 873, 874, 875, 876, 877, 878, 879, 880, 881, 882, 883, 884, 885, 886, 887, 888, 889, 890, 891, 892, 893, 894, 895, 896, 897, 898, 899, 900, 901, 902, 903, 904, 905, 906, 907, 908, 909, 910, 911, 912, 913, 914, 915, 916, 917, 918, 919, 920, 921, 922, 923, 924, 925, 926, 927, 928, 929, 930, 931, 932, 933, 934, 935, 936, 937, 938, 939, 940, 941, 942, 943, 944, 945, 946, 947, 948, 949, 950, 951, 952, 953, 954, 955, 956, 957, 958, 959, 960, 961, 962, 963, 964, 965, 966, 967, 968, 969, 970, 971, 972, 973, 974, 975, 976, 977, 978, 979, 980, 981, 982, 983, 984, 985, 986, 987, 988, 989, 990, 991, 992, 993, 994, 995, 996, 997, 998, 999, 1000, 1001, 1002, 1003, 1004, 1005, 1006, 1007, 1008, 1009, 1010, 1011, 1012, 1013, 1014, 1015, 1016, 1017, 1018, 1019, 1020, 1021, 1022, 1023, 1024, 1025, 1026, 1027, 1028, 1029, 1030, 1031, 1032, 1033, 1034, 1035, 1036, 1037, 1038, 1039, 1040, 1041, 1042, 1043, 1044, 1045, 1046, 1047, 1048, 1049, 1050, 1051, 1052, 1053, 1054, 1055, 1056, 1057, 1058, 1059, 1060, 1061, 1062, 1063, 1064, 1065, 1066, 1067, 1068, 1069, 1070, 1071, 1072, 1073, 1074, 1075, 1076, 1077, 1078, 1079, 1080, 1081, 1082, 1083, 1084, 1085, 1086, 1087, 1088, 1089, 1090, 1091, 1092, 1093, 1094, 1095, 1096, 1097, 1098, 1099, 1100, 1101, 1102, 1103, 1104, 1105, 1106, 1107, 1108, 1109, 1110, 1111, 1112, 1113, 1114, 1115, 1116, 1117, 1118, 1119, 1120, 1121, 1122, 1123, 1124, 1125, 1126, 1127, 1128, 1129, 1130, 1131, 1132, 1133, 1134, 1135, 1136, 1137, 1138, 1139, 1140, 1141, 1142, 1143, 1144, 1145, 1146, 1147, 1148, 1149, 1150, 1151, 1152, 1153, 1154, 1155, 1156, 1157, 1158, 1159, 1160, 1161, 1162, 1163, 1164, 1165, 1166, 1167, 1168, 1169, 1170, 1171, 1172, 1173, 1174, 1175, 1176, 1177, 1178, 1179, 1180, 1181, 1182, 1183, 1184, 1185, 1186, 1187, 1188, 1189, 1190, 1191, 1192, 1193, 1194, 1195, 1196, 1197, 1198, 1199, 1200, 1201, 1202, 1203, 1204, 1205, 1206, 1207, 1208, 1209, 1210, 1211, 1212, 1213, 1214, 1215, 1216, 1217, 1218, 1219, 1220, 1221, 1222, 1223, 1224, 1225, 1226, 1227, 1228, 1229, 1230, 1231, 1232, 1233, 1234, 1235, 1236, 1237, 1238, 1239, 1240, 1241, 1242, 1243, 1244, 1245, 1246, 1247, 1248, 1249, 1250, 1251, 1252, 1253, 1254, 1255, 1256, 1257, 1258, 1259, 1260, 1261, 1262, 1263, 1264, 1265, 1266, 1267, 1268, 1269, 1270, 1271, 1272, 1273, 1274, 1275, 1276, 1277, 1278, 1279, 1280, 1281, 1282, 1283, 1284, 1285, 1286, 1287, 1288, 1289, 1290, 1291, 1292, 1293, 1294, 1295, 1296, 1297, 1298, 1299, 1300, 1301, 1302, 1303, 1304, 1305, 1306, 1307, 1308, 1309, 1310, 1311, 1312, 1313, 1314, 1315, 1316, 1317, 1318, 1319, 1320, 1321, 1322, 1323, 1324, 1325, 1326, 1327, 1328, 1329, 1330, 1331, 1332, 1333, 1334, 1335, 1336, 1337, 1338, 1339, 1340, 1341, 1342, 1343, 1344, 1345, 1346, 1347, 1348, 1349, 1350, 1351, 1352, 1353, 1354, 1355, 1356, 1357, 1358, 1359, 1360, 1361, 1362, 1363, 1364, 1365, 1366, 1367, 1368, 1369, 1370, 1371, 1372, 1373, 1374, 1375, 1376, 1377, 1378, 1379, 1380, 1381, 1382, 1383, 1384, 1385, 1386, 1387, 1388, 1389, 1390, 1391, 1392, 1393, 1394, 1395, 1396, 1397, 1398, 1399, 1400, 1401, 1402, 1403, 1404, 1405, 1406, 1407, 1408, 1409, 1410, 1411, 1412, 1413, 1414, 1415, 1416, 1417, 1418, 1419, 1420, 1421, 1422, 1423, 1424, 1425, 1426, 1427, 1428, 1429, 1430, 1431, 1432, 1433, 1434, 1435, 1436, 1437, 1438, 1439, 1440, 1441, 1442, 1443, 1444, 1445, 1446, 1447, 1448, 1449, 1450, 1451, 1452, 1453, 1454, 1455, 1456, 1457, 1458, 1459, 1460, 1461, 1462, 1463, 1464, 1465, 1466, 1467, 1468, 1469, 1470, 1471, 1472, 1473, 1474, 1475, 1476, 1477, 1478, 1479, 1480, 1481, 1482, 1483, 1484, 1485, 1486, 1487, 1488, 1489, 1490, 1491, 1492, 1493, 1494, 1495, 1496, 1497, 1498, 1499, 1500, 1501, 1502, 1503, 1504, 1505, 1506, 1507, 1508, 1509, 1510, 1511, 1512, 1513, 1514, 1515, 1516, 1517, 1518, 1519, 1520, 1521, 1522, 1523, 1524, 1525, 1526, 1527, 1528, 1529, 1530, 1531, 1532, 1533, 1534, 1535, 1536, 1537, 1538, 1539, 1540, 1541, 1542, 1543, 1544, 1545, 1546, 1547, 1548, 1549, 1550, 1551, 1552, 1553, 1554, 1555, 1556, 1557, 1558, 1559, 1560, 1561, 1562, 1563, 1564, 1565, 1566, 1567, 1568, 1569, 1570, 1571, 1572, 1573, 1574, 1575, 1576, 1577, 1578, 1579, 1580, 1581, 1582, 1583, 1584, 1585, 1586, 1587, 1588, 1589, 1590, 1591, 1592, 1593, 1594, 1595, 1596, 1597, 1598, 1599, 1600, 1601, 1602, 1603, 1604, 1605, 1606, 1607, 1608, 1609, 1610, 1611, 1612, 1613, 1614, 1615, 1616, 1617, 1618, 1619, 1620, 1621, 1622, 1623, 1624, 1625, 1626, 1627, 1628, 1629, 1630, 1631, 1632, 1633, 1634, 1635, 1636, 1637, 1638, 1639, 1640, 1641, 1642, 1643, 1644, 1645, 1646, 1647, 1648, 1649, 1650, 1651, 1652, 1653, 1654, 1655, 1656, 1657, 1658, 1659, 1660, 1661, 1662, 1663, 1664, 1665, 1666, 1667, 1668, 1669, 1670, 1671, 1672, 1673, 1674, 1675, 1676, 1677, 1678, 1679, 1680, 1681, 1682, 1683, 1684, 1685, 1686, 1687, 1688, 1689, 1690, 1691, 1692, 1693, 1694, 1695, 1696, 1697, 1698, 1699, 1700, 1701, 1702, 1703, 1704, 1705, 1706, 1707, 1708, 1709, 1710, 1711, 1712, 1713, 1714, 1715, 1716, 1717, 1718, 1719, 1720, 1721, 1722, 1723, 1724, 1725, 1726, 1727, 1728, 1729, 1730, 1731, 1732, 1733, 1734, 1735, 1736, 1737, 1738, 1739, 1740, 1741, 1742, 1743, 1744, 1745, 1746, 1747, 1748, 1749, 1750, 1751, 1752, 1753, 1754, 1755, 1756, 1757, 1758, 1759, 1760, 1761, 1762, 1763, 1764, 1765, 1766, 1767, 1768, 1769, 1770, 1771, 1772, 1773, 1774, 1775, 1776, 1777, 1778, 1779, 1780, 1781, 1782, 1783, 1784, 1785, 1786, 1787, 1788, 1789, 1790, 1791, 1792, 1793, 1794, 1795, 1796, 1797, 1798, 1799, 1800, 1801, 1802, 1803, 1804, 1805, 1806, 1807, 1808, 1809, 1810, 1811, 1812, 1813, 1814, 1815, 1816, 1817, 1818, 1819, 1820, 1821, 1822, 1823, 1824, 1825, 1826, 1827, 1828, 1829, 1830, 1831, 1832, 1833, 1834, 1835, 1836, 1837, 1838, 1839, 1840, 1841, 1842, 1843, 1844, 1845, 1846, 1847, 1848, 1849, 1850, 1851, 1852, 1853, 1854, 1855, 1856, 1857, 1858, 1859, 1860, 1861, 1862, 1863, 1864, 1865, 1866, 1867, 1868, 1869, 1870, 1871, 1872, 1873, 1874, 1875, 1876, 1877, 1878, 1879, 1880, 1881, 1882, 1883, 1884, 1885, 1886, 1887, 1888, 1889, 1890, 1891, 1892, 1893, 1894, 1895, 1896, 1897, 1898, 1899, 1900, 1901, 1902, 1903, 1904, 1905, 1906, 1907, 1908, 1909, 1910, 1911, 1912, 1913, 1914, 1915, 1916, 1917, 1918, 1919, 1920, 1921, 1922, 1923, 1924, 1925, 1926, 1927, 1928, 1929, 1930, 1931, 1932, 1933, 1934, 1935, 1936, 1937, 1938, 1939, 1940, 1941, 1942, 1943, 1944, 1945, 1946, 1947, 1948, 1949, 1950, 1951, 1952, 1953, 1954, 1955, 1956, 1957, 1958, 1959, 1960, 1961, 1962, 1963, 1964, 1965, 1966, 1967, 1968, 1969, 1970, 1971, 1972, 1973, 1974, 1975, 1976, 1977, 1978, 1979, 1980, 1981, 1982, 1983, 1984, 1985, 1986, 1987, 1988, 1989, 1990, 1991, 1992, 1993, 1994, 1995, 1996, 1997, 1998, 1999, 2000, 2001, 2002, 2003, 2004, 2005, 2006, 2007, 2008, 2009, 2010, 2011, 2012, 2013, 2014, 2015, 2016, 2017, 2018, 2019, 2020, 2021, 2022, 2023, 2024, 2025, 2026, 2027, 2028, 2029, 2030, 2031, 2032, 2033, 2034, 2035, 2036, 2037, 2038, 2039, 2040, 2041, 2042, 2043, 2044, 2045, 2046, 2047, 2048, 2049, 2050, 2051, 2052, 2053, 2054, 2055, 2056, 2057, 2058, 2059, 2060, 2061, 2062, 2063, 2064, 2065, 2066, 2067, 2068, 2069, 2070, 2071, 2072, 2073, 2074, 2075, 2076, 2077, 2078, 2079, 2080, 2081, 2082, 2083, 2084, 2085, 2086, 2087, 2088, 2089, 2090, 2091, 2092, 2093, 2094, 2095, 2096, 2097, 2098, 2099, 2100, 2101, 2102, 2103, 2104, 2105, 2106, 2107, 2108, 2109, 2110, 2111, 2112, 2113, 2114, 2115, 2116, 2117, 2118, 2119, 2120, 2121, 2122, 212