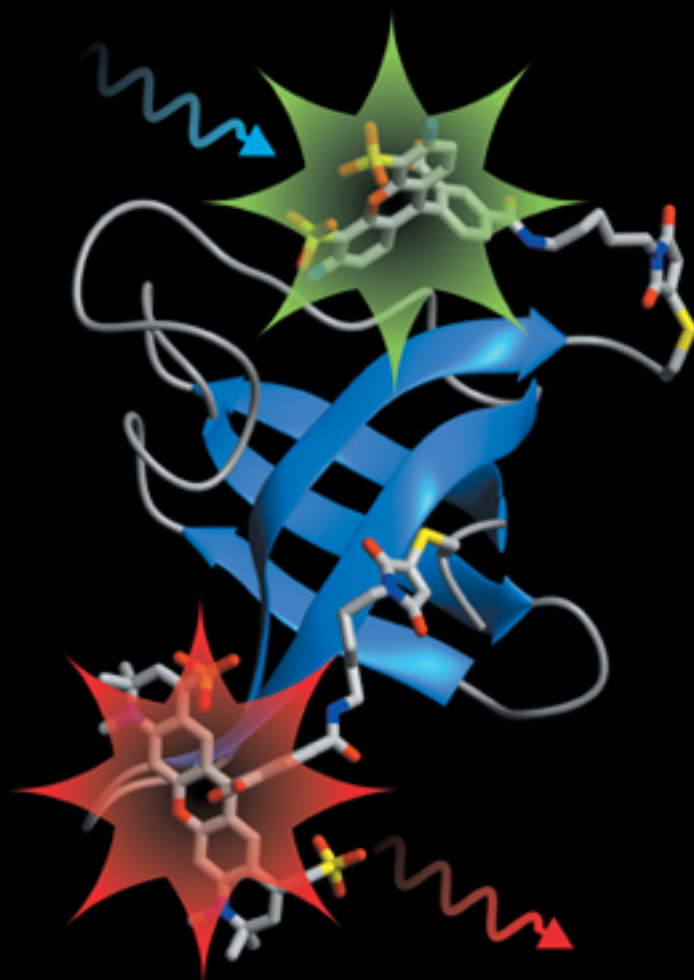


Molecular and Protein Engineering



Nohemi Altman

Sharlene Boswell

First Edition, 2012

ISBN 978-81-323-0957-4

© All rights reserved.

Published by:
Academic Studio
4735/22 Prakashdeep Bldg,
Ansari Road, Darya Ganj,
Delhi - 110002
Email: info@wtbooks.com

Table of Contents

Chapter 1 - Molecule

Chapter 2 - Molecular Modelling and Molecular Design Software

Chapter 3 - Molecular Graphics

Chapter 4 - Molecular Machine

Chapter 5 - Molecular Assembler

Chapter 6 - Molecular Orbital and Molecular Orbital Theory

Chapter 7 - Molecular Model

Chapter 8 - Protein Structure

Chapter 9 - Protein Folding

Chapter 10 - Protein Design & Fusion Protein

Chapter 11 - Directed Evolution

Chapter 12 - Protein Domain

Chapter 13 - Proteomics

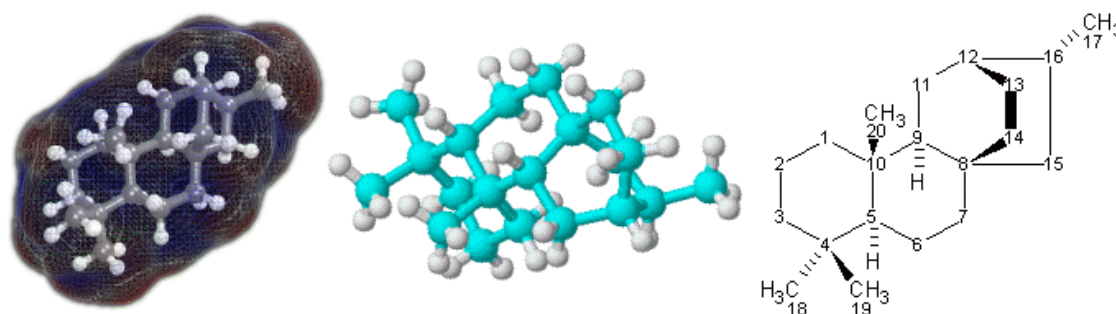
Chapter 14 - Proteinogenic Amino Acid

Chapter 15 - Protein

Chapter 16 - Cytoskeleton

Chapter 1

Molecule



3D (left and center) and 2D (right) representations of the terpenoid molecule atisane

A **molecule** is an electrically neutral group of at least two atoms held together by covalent chemical bonds. Molecules are distinguished from ions by their electrical charge. However, in quantum physics, organic chemistry, and biochemistry, the term *molecule* is often used less strictly and applied to polyatomic ions.

In the kinetic theory of gases, the term *molecule* is often used for any gaseous particle regardless of its composition. According to this definition noble gas atoms are considered molecules despite the fact that they are composed of a single non-bonded atom.

A molecule may consist of atoms of a single chemical element, as with oxygen (O_2), or of different elements, as with water (H_2O). Atoms and complexes connected by non-covalent bonds such as hydrogen bonds or ionic bonds are generally not considered single molecules.

Molecules as components of matter are common in organic substances (and therefore biochemistry). They also make up most of the oceans and atmosphere. A large number of familiar solid substances, however, including most of the minerals that make up the crust, mantle, and core of the Earth itself, contain many chemical bonds, but are *not* made of identifiable molecules. No typical molecule can be defined for ionic crystals (salts) and covalent crystals (network solids), although these are often composed of repeating unit cells that extend either in a plane (such as in graphene) or three-dimensionally (such as in

diamond or sodium chloride). The theme of repeated unit-cellular-structure also holds for most condensed phases with metallic bonding. In glasses (solids that exist in a vitreous disordered state), atoms may also be held together by chemical bonds without any definable molecule, but also without any of the regularity of repeating units that characterises crystals.

Molecular science

The science of molecules is called *molecular chemistry* or *molecular physics*, depending on the focus. Molecular chemistry deals with the laws governing the interaction between molecules that results in the formation and breakage of chemical bonds, while molecular physics deals with the laws governing their structure and properties. In practice, however, this distinction is vague. In molecular sciences, a molecule consists of a stable system (bound state) comprising two or more atoms. Polyatomic ions may sometimes be usefully thought of as electrically charged molecules. The term *unstable molecule* is used for very reactive species, i.e., short-lived assemblies (resonances) of electrons and nuclei, such as radicals, molecular ions, Rydberg molecules, transition states, van der Waals complexes, or systems of colliding atoms as in Bose-Einstein condensate

History and etymology

According to Merriam-Webster and the Online Etymology Dictionary, the word "molecule" derives from the Latin "moles" or small unit of mass.

- **Molecule** (1794) – "extremely minute particle," from Fr. *molécule* (1678), from Mod.L. *molecula*, dim. of L. *moles* "mass, barrier". A vague meaning at first; the vogue for the word (used until late 18th century only in Latin form) can be traced to the philosophy of Descartes.

Although the existence of molecules has been accepted by many chemists since the early 19th century as a result of Dalton's laws of Definite and Multiple Proportions (1803–1808) and Avogadro's law (1811), there was some resistance among positivists and physicists such as Mach, Boltzmann, Maxwell, and Gibbs, who saw molecules merely as convenient mathematical constructs. The work of Perrin on Brownian motion (1911) is considered to be the final proof of the existence of molecules.

The definition of the molecule has evolved as knowledge of the structure of molecules has increased. Earlier definitions were less precise, defining molecules as the smallest particles of pure chemical substances that still retain their composition and chemical properties. This definition often breaks down since many substances in ordinary experience, such as rocks, salts, and metals, are composed of large networks of chemically bonded atoms or ions, but are not made of discrete molecules.

Molecular size

Most molecules are far too small to be seen with the naked eye, but there are exceptions. DNA, a macromolecule, can reach macroscopic sizes, as can molecules of many polymers. The smallest molecule is the diatomic hydrogen (H_2), with a length of 0.74 Å. Molecules commonly used as building blocks for organic synthesis have a dimension of a few Å to several dozen Å. Single molecules cannot usually be observed by light (as noted above), but small molecules and even the outlines of individual atoms may be traced in some circumstances by use of an atomic force microscope. Some of the largest molecules are macromolecules or supermolecules.

Radius

Effective molecular radius is the size a molecule displays in solution. The table of permselectivity for different substances contains examples.

Molecular formula

A compound's empirical formula is the **simplest** integer ratio of the chemical elements that constitute it. For example, water is always composed of a 2:1 ratio of hydrogen to oxygen atoms, and ethyl alcohol or ethanol is always composed of carbon, hydrogen, and oxygen in a 2:6:1 ratio. However, this does not determine the kind of molecule uniquely – dimethyl ether has the same ratios as ethanol, for instance. Molecules with the same atoms in different arrangements are called isomers. Also carbohydrates, for example, have the same ratio (carbon:hydrogen:oxygen = 1:2:1) (and thus the same empirical formula) but different total numbers of atoms in the molecule.

The molecular formula reflects the exact number of atoms that compose the molecule and so characterizes different molecules. However different isomers can have the same atomic composition while being different molecules.

The empirical formula is often the same as the molecular formula but not always. For example the molecule acetylene has molecular formula C_2H_2 , but the simplest integer ratio of elements is CH.

The molecular mass can be calculated from the chemical formula and is expressed in conventional atomic mass units equal to 1/12 of the mass of a neutral carbon-12 (^{12}C isotope) atom. For network solids, the term formula unit is used in stoichiometric calculations.

Molecular geometry

Molecules have fixed equilibrium geometries—bond lengths and angles—about which they continuously oscillate through vibrational and rotational motions. A pure substance is composed of molecules with the same average geometrical structure. The chemical formula and the structure of a molecule are the two important factors that determine its

properties, particularly its reactivity. Isomers share a chemical formula but normally have very different properties because of their different structures. Stereoisomers, a particular type of isomers, may have very similar physico-chemical properties and at the same time different biochemical activities.

Molecular spectroscopy

Molecular spectroscopy deals with the response (spectrum) of molecules interacting with probing signals of known energy (or frequency, according to Planck's formula). Molecules have quantized energy levels that can be analyzed by detecting the molecule's energy exchange through absorbance or emission. Spectroscopy does not generally refer to diffraction studies where particles such as neutrons, electrons, or high energy X-rays interact with a regular arrangement of molecules (as in a crystal).

Theoretical aspects

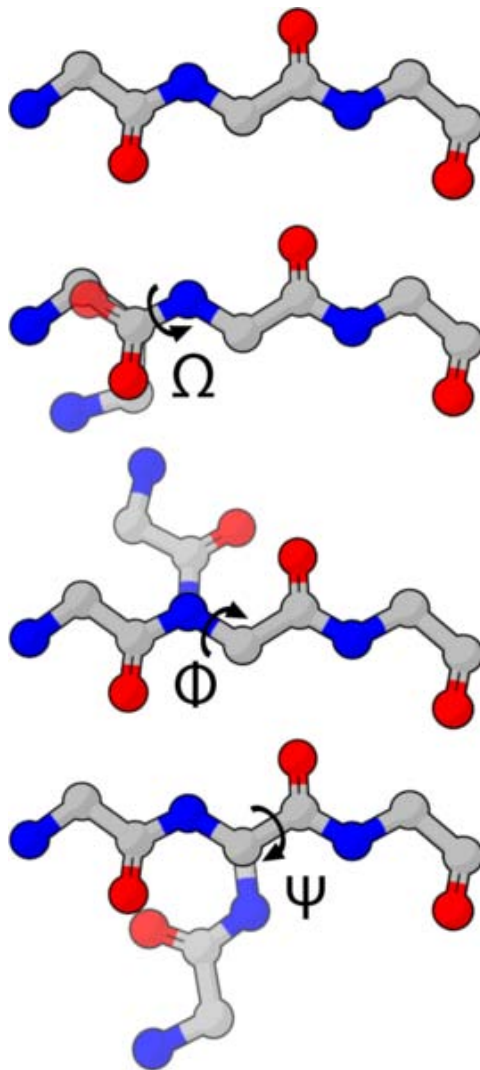
The study of molecules by molecular physics and theoretical chemistry is largely based on quantum mechanics and is essential for the understanding of the chemical bond. The simplest of molecules is the hydrogen molecule-ion, H_2^+ , and the simplest of all the chemical bonds is the one-electron bond. H_2^+ is composed of two positively charged protons and one negatively charged electron, which means that the Schrödinger equation for the system can be solved more easily due to the lack of electron–electron repulsion. With the development of fast digital computers, approximate solutions for more complicated molecules became possible and are one of the main aspects of computational chemistry.

When trying to define rigorously whether an arrangement of atoms is "sufficiently stable" to be considered a molecule, IUPAC suggests that it "must correspond to a depression on the potential energy surface that is deep enough to confine at least one vibrational state". This definition does not depend on the nature of the interaction between the atoms, but only on the strength of the interaction. In fact, it includes weakly bound species that would not traditionally be considered molecules, such as the helium dimer, He_2 , which has one vibrational bound state and is so loosely bound that it is only likely to be observed at very low temperatures.

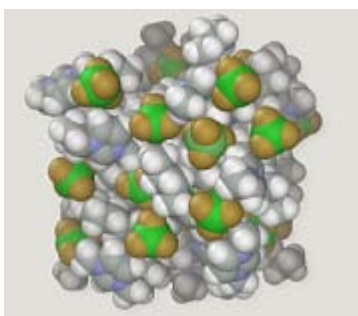
Chapter 2

Molecular Modelling and Molecular Design Software

Molecular modelling



The backbone dihedral angles are included in the molecular model of a protein.



Modelling of ionic liquid

Molecular modelling encompasses all theoretical methods and computational techniques used to model or mimic the behaviour of molecules. The techniques are used in the fields of computational chemistry, computational biology and materials science for studying molecular systems ranging from small chemical systems to large biological molecules and material assemblies. The simplest calculations can be performed by hand, but inevitably computers are required to perform molecular modelling of any reasonably sized system. The common feature of molecular modelling techniques is the atomistic level description of the molecular systems; the lowest level of information is individual atoms (or a small group of atoms). This is in contrast to quantum chemistry (also known as electronic structure calculations) where electrons are considered explicitly. The benefit of molecular modelling is that it reduces the complexity of the system, allowing many more particles (atoms) to be considered during simulations.

Molecular mechanics

Molecular mechanics is one aspect of molecular modelling, as it refers to the use of classical mechanics/Newtonian mechanics to describe the physical basis behind the models. Molecular models typically describe atoms (nucleus and electrons collectively) as point charges with an associated mass. The interactions between neighbouring atoms are described by spring-like interactions (representing chemical bonds) and van der Waals forces. The Lennard-Jones potential is commonly used to describe van der Waals forces. The electrostatic interactions are computed based on Coulomb's law. Atoms are assigned coordinates in Cartesian space or in internal coordinates, and can also be assigned velocities in dynamical simulations. The atomic velocities are related to the temperature of the system, a macroscopic quantity. The collective mathematical expression is known as a potential function and is related to the system internal energy (U), a thermodynamic quantity equal to the sum of potential and kinetic energies. Methods which minimize the potential energy are known as energy minimization techniques (e.g., steepest descent and conjugate gradient), while methods that model the behaviour of the system with propagation of time are known as molecular dynamics.

$$E = E_{\text{bonds}} + E_{\text{angle}} + E_{\text{dihedral}} + E_{\text{non-bonded}}$$
$$E_{\text{non-bonded}} = E_{\text{electrostatic}} + E_{\text{van der Waals}}$$

This function, referred to as a potential function, computes the molecular potential energy as a sum of energy terms that describe the deviation of bond lengths, bond angles and torsion angles away from equilibrium values, plus terms for non-bonded pairs of atoms describing van der Waals and electrostatic interactions. The set of parameters consisting of equilibrium bond lengths, bond angles, partial charge values, force constants and van der Waals parameters are collectively known as a force field. Different implementations of molecular mechanics use different mathematical expressions and different parameters for the potential function. The common force fields in use today have been developed by using high level quantum calculations and/or fitting to experimental data. The technique known as energy minimization is used to find positions of zero gradient for all atoms, in other words, a local energy minimum. Lower energy states are more stable and are commonly investigated because of their role in chemical and biological processes. A molecular dynamics simulation, on the other hand, computes the behaviour of a system as a function of time. It involves solving Newton's laws of motion, principally the second law, $\mathbf{F} = m\mathbf{a}$. Integration of Newton's laws of motion, using different integration algorithms, leads to atomic trajectories in space and time. The force on an atom is defined as the negative gradient of the potential energy function. The energy minimization technique is useful for obtaining a static picture for comparing between states of similar systems, while molecular dynamics provides information about the dynamic processes with the intrinsic inclusion of temperature effects.

Variables

Molecules can be modelled either in vacuum or in the presence of a solvent such as water. Simulations of systems in vacuum are referred to as *gas-phase* simulations, while those that include the presence of solvent molecules are referred to as *explicit solvent* simulations. In another type of simulation, the effect of solvent is estimated using an empirical mathematical expression; these are known as *implicit solvation* simulations.

Applications

Molecular modelling methods are now routinely used to investigate the structure, dynamics, surface properties and thermodynamics of inorganic, biological and polymeric systems. The types of biological activity that have been investigated using molecular modelling include protein folding, enzyme catalysis, protein stability, conformational changes associated with biomolecular function, and molecular recognition of proteins, DNA, and membrane complexes.

Molecular design software

Molecular design software is a software for molecular modeling, distinctive property of which is the presence of the special support for developing the molecular models.

In contrast to the usual molecular modeling programs such as the molecular dynamics and quantum chemistry programs, such software directly supports the aspects related to the construction of molecular models:

- Molecular graphics
- interactive molecular drawing and conformational editing
- building of polymeric molecules, crystals and solvated systems
- partial charges development
- geometry optimization
- support for the different aspects of Force Field development
- *etc.*

Comparative table of packages covering the major aspects of molecular design

3D - Molecular Graphics, **Mouse** - drawing molecule by mouse, **Poly** - polymer building, **DNA** - Nucleic acid building, **Pept** - Peptide building, **Cryst** - crystal building, **Solv** - solvent addition, **Q** - partial charges, **Dock** - docking, **Min** - optimization, **MM** - Molecular mechanics, **QM** - Quantum mechanics. **FF** - Support for Force Field development. **QSAR** - 2D, 3D and Group QSAR.

	3D	Mouse	Poly	DNA	Pept	Cryst	Solv	Q	Dock	Min	MM	QM	FF	QSAR	Homepage	Comments
AMBER				+	+			+		+	+			+	ambermd.org	Classical molecular modeling program
ArgusLab	+	+			+			++		+	+				Planaria Software	A molecular modeling, graphics, and drug design program
Ascalaph Designer	+	+	+	+	+	+	+	+		+	+	+	+		Agile Molecule	common molecular modeling suite
Avogadro (software)	+	+			+	+		+		+	+			+	OpenMolecules.net	Extensible, free, open source molecular editor

BALL / BALLView	+	+			+	+			+	+	+	+	+	ball-project.org	open source, C++ molecular modelling and visualization tool with scripting interface
BOSS									+	+	+	+	+	Yale University	OPLS inventor
DOCK									+	+				University of California	DOCK algorithm
Firefly (PC GAMES S)									+	+			+	Moscow State University	ab initio and DFT computational chemistry program
FoldX									+	+	+		+	CRG	A force field for energy calculations and protein design
Maestro	+	+	+	+	+	+	+	+	+	+	+	+	+	Schrodinger	A molecular modeling, visualization, and drug design program
Materials Studio	+	+	+			+	+	+	+	+	+	+	+	Accelrys	software environment
MedeA	+		+			+				+	+	+		Materials Design	software environment for inorganic materials science
MOE	+	+	+	+	+	+	+	+	+	+	+	+	+	Chemical Computing Group	Molecular Operating Environment
VLifeMDS	+	+							+	+	+	+	+	VLife Sciences Technologies Pvt. Ltd.	VLife Molecular Design Suite
NAB				+						+	+			Rutgers University	molecular manipulation language for nucleic acids
PCMODEL	+		+	+	+				+	+	+	+		Serena Software	common molecular modeling tool
SPARTAN	+	+		+	+			+	+	+	+	+		Wavefunction	molecular modeling tool with molecular mechanics and quantum chemical engines
StruMM3D	+	+	+	+	+	+	+	+	+	+				Exorga Software	molecular modeling tool

Chapter 3

Molecular Graphics

Molecular graphics (MG) is the discipline and philosophy of studying molecules and their properties through graphical representation. IUPAC limits the definition to representations on a "graphical display device". Ever since Dalton's atoms and Kekule's benzene, there has been a rich history of hand-drawn atoms and molecules, and these representations have had an important influence on modern molecular graphics. Here we concentrate on the use of computers to create molecular graphics. Note, however, that many molecular graphics programs and systems have close coupling between the graphics and editing commands or calculations such as in molecular modelling.

Relation to molecular models

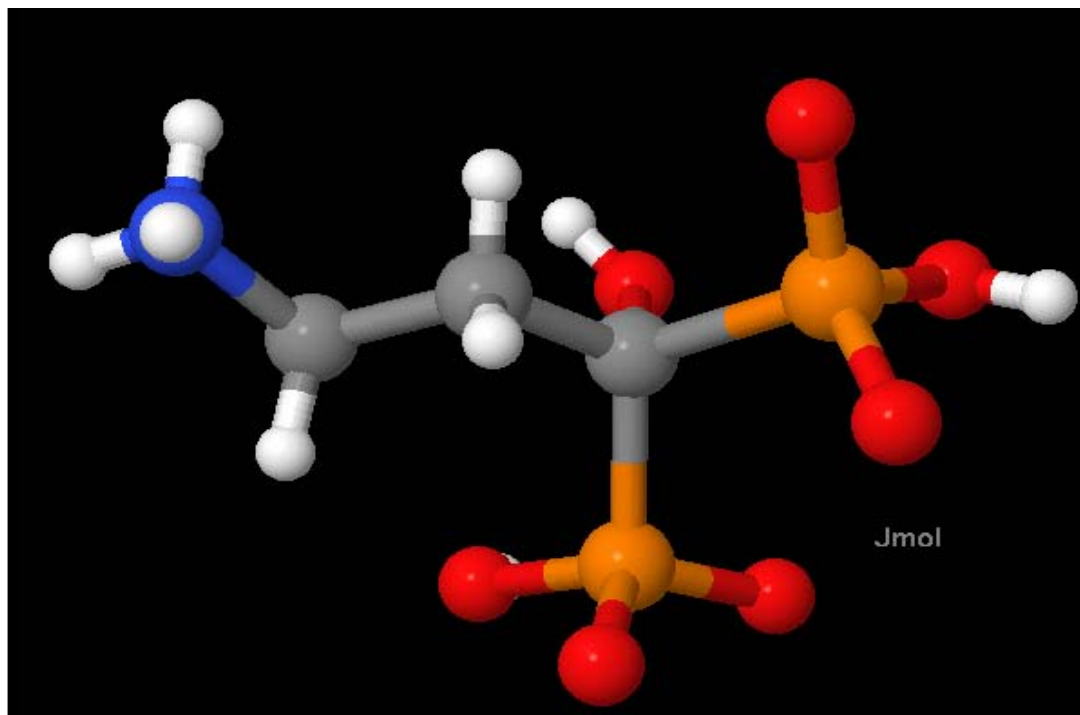


Fig. 1. Key: Hydrogen = white, carbon = grey, nitrogen = blue, oxygen = red, and phosphorus = orange.

There has been a long tradition of creating molecular models from physical materials. Perhaps the best known is Crick and Watson's model of DNA built from rods and planar sheets, but the most widely used approach is to represent all atoms and bonds explicitly using the "ball and stick" approach. This can demonstrate a wide range of properties, such as shape, relative size, and flexibility. Many chemistry courses expect that students will have access to ball and stick models. One goal of mainstream molecular graphics has been to represent the "ball and stick" model as realistically as possible and to couple this with calculations of molecular properties.

Figure 1 shows a small molecule ($\text{NH}_3\text{CH}_2\text{CH}_2\text{C}(\text{OH})(\text{PO}_3\text{H})(\text{PO}_3\text{H})^-$), as drawn by the Jmol program. It is important to realize that the colors and shapes are purely a convention, as individual atoms are not colored, nor do they have hard surfaces. Bonds between atoms are also not rod-shaped.

Comparison of physical models with molecular graphics

Physical models and computer models have partially complementary strengths and weaknesses. Physical models can be used by those without access to a computer and now can be made cheaply out of plastic materials. Their tactile and visual aspects cannot be easily reproduced by computers (although haptic devices have occasionally been built). On a computer screen, the flexibility of molecules is also difficult to appreciate; illustrating the pseudorotation of cyclohexane is a good example of the value of mechanical models.

However, it is difficult to build large physical molecules, and all-atom physical models of even simple proteins could take weeks or months to build. Moreover, physical models are not robust and they decay over time. Molecular graphics is particularly valuable for representing global and local properties of molecules, such as electrostatic potential. Graphics can also be animated to represent molecular processes and chemical reactions, a feat that is not easy to reproduce physically.

History

Initially the rendering was on early Cathode ray tube screens or through plotters drawing on paper. Molecular structures have always been an attractive choice for developing new computer graphics tools, since the input data are easy to create and the results are usually highly appealing. The first example of MG was a display of a protein molecule (Project MAC, 1966) by Cyrus Levinthal and Robert Langridge. Among the milestones in high-performance MG was the work of Nelson Max in "realistic" rendering of macromolecules using reflecting spheres.

By about 1980 many laboratories both in academia and industry had recognized the power of the computer to analyse and predict the properties of molecules, especially in materials science and the pharmaceutical industry. The discipline was often called "molecular graphics" and in 1982 a group of academics and industrialists in the UK set up the Molecular Graphics Society (MGS). Initially much of the technology concentrated

either on high-performance 3D graphics, including interactive rotation or 3D rendering of atoms as spheres (sometimes with radiosity). During the 1980s a number of programs for calculating molecular properties (such as molecular dynamics and quantum mechanics) became available and the term "molecular graphics" often included these. As a result the MGS has now changed its name to the Molecular Graphics and Modelling Society (MGMS).

The requirements of macromolecular crystallography also drove MG because the traditional techniques of physical model-building could not scale. Alwyn Jones' FRODO program (and later "O") were developed to overlay the molecular electron density determined from X-ray crystallography and the hypothetical molecular structure.

Art, science and technology in molecular graphics

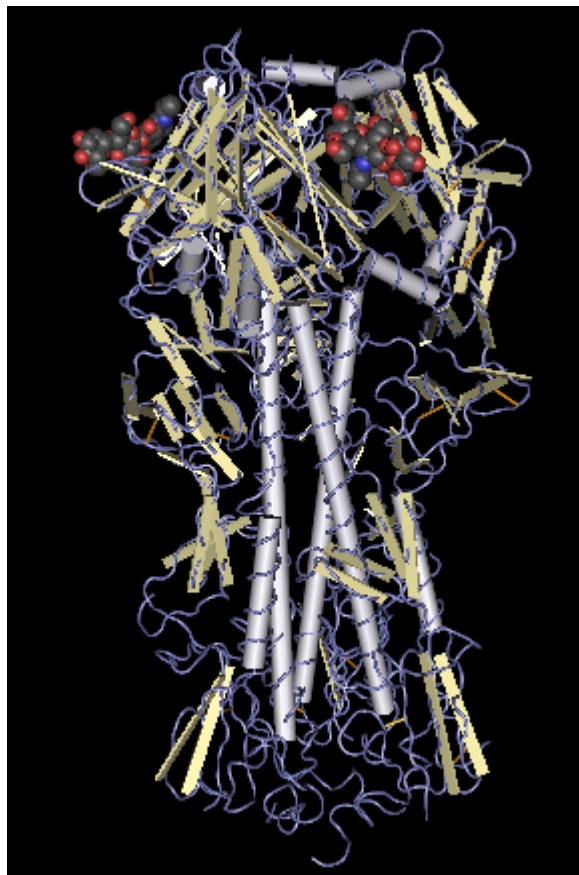


Fig. 2. Image of hemagglutinin with alpha helices depicted as cylinders and the rest of the chain as silver coils. The individual protein atoms (several thousand) have been hidden. All of the non-hydrogen atoms in the two ligands (presumably sialic acid) have been shown near the top of the diagram. Key: Carbon = grey, oxygen = red, nitrogen = blue.

Both computer technology and graphic arts have contributed to molecular graphics. The development of structural biology in the 1950s led to a requirement to display molecules with thousands of atoms. The existing computer technology was limited in power, and in

any case a naive depiction of all atoms left viewers overwhelmed. Most systems therefore used conventions where information was implicit or stylistic. Two vectors meeting at a point implied an atom or (in macromolecules) a complete residue (10-20 atoms).

The macromolecular approach was popularized by Dickerson and Geis' presentation of proteins and the graphic work of Jane Richardson through high-quality hand-drawn diagrams such as the "ribbon" representation. In this they strove to capture the intrinsic 'meaning' of the molecule. This search for the "messages in the molecule" has always accompanied the increasing power of computer graphics processing. Typically the depiction would concentrate on specific areas of the molecule (such as the active site) and this might have different colours or more detail in the number of explicit atoms or the type of depiction (e.g., spheres for atoms).

In some cases the limitations of technology have led to serendipitous methods for rendering. Most early graphics devices used vector graphics, which meant that rendering spheres and surfaces was impossible. Michael Connolly's program "MS" calculated points on the surface-accessible surface of a molecule, and the points were rendered as dots with good visibility using the new vector graphics technology, such as the Evans and Sutherland PS300 series. Thin sections ("slabs") through the structural display showed very clearly the complementarity of the surfaces for molecules binding to active sites, and the "Connolly surface" became a universal metaphor.

The relationship between the art and science of molecular graphics is shown in the exhibitions sponsored by the Molecular Graphics Society. Some exhibits are created with molecular graphics programs alone, while others are collages, or involve physical materials. An example from Mike Hann (1994), inspired by Magritte's painting *Ceci n'est pas une pipe*, uses an image of a salmeterol molecule.

"*Ceci n'est pas une molecule*," writes Mike Hann, "serves to remind us that all of the graphics images presented here are not molecules, not even pictures of molecules, but pictures of icons which we believe represent some aspects of the molecule's properties."

Space-filling models

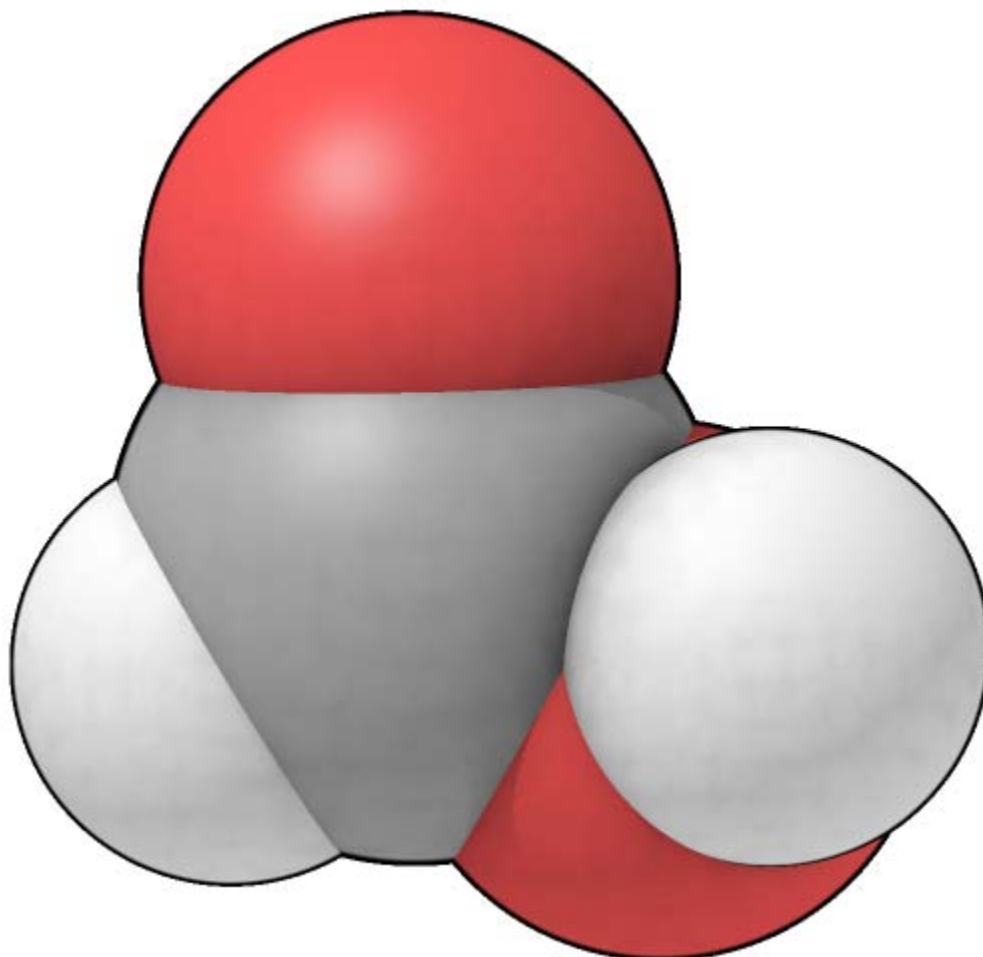


Fig. 4. Space-filling model of formic acid. Key: Hydrogen = white, carbon = black, oxygen = red.

Fig. 4 is a "space-filling" representation of formic acid, where atoms are drawn to suggest the amount of space they occupy. This is necessarily an icon: in the quantum mechanical representation of molecules, there are only (positively charged) nuclei and a "cloud" of negative electrons. The electron cloud defines an approximate size for the molecule, though there can be no single precise definition of size. For many years the size of atoms has been approximated by mechanical models (CPK), where the atoms have been represented by plastic spheres whose radius (van der Waals radius) describes a sphere within which "most" of the electron density can be found. These spheres could be clicked together to show the steric aspects of the molecule rather than the positions of the nuclei. Fig. 4 shows the intricacy required to make sure that all spheres intersect correctly, and also demonstrates a reflective model.

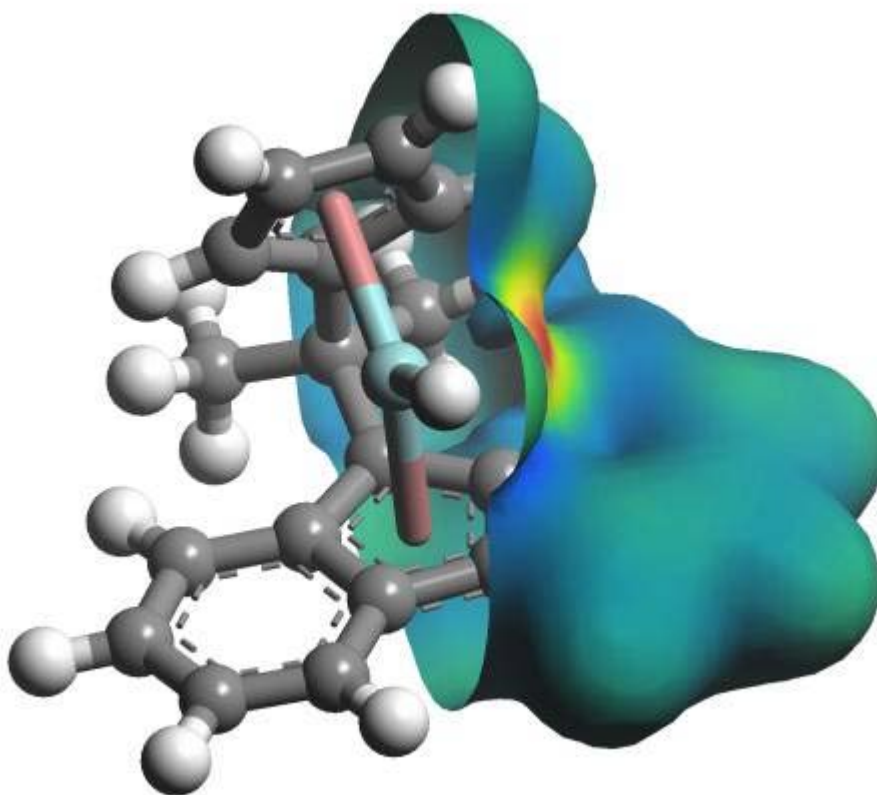


Fig. 5. A molecule (zirconocene) where part (left) is rendered as ball-and-stick and part (right) as an isosurface.

Since the atomic radii (e.g. in Fig. 4) are only slightly less than the distance between bonded atoms, the iconic spheres intersect, and in the CPK models, this was achieved by planar truncations along the bonding directions, the section being circular. When raster graphics became affordable, one of the common approaches was to replicate CPK models *in silico*. It is relatively straightforward to calculate the circles of intersection, but more complex to represent a model with hidden surface removal. A useful side product is that a conventional value for the molecular volume can be calculated.

The use of spheres is often for convenience, being limited both by graphics libraries and the additional effort required to compute complete electronic density or other space-filling quantities. It is now relatively common to see images of surfaces that have been colored to show quantities such as electrostatic potential. Common surfaces in molecular visualization include solvent-accessible ("Lee-Richards") surfaces, solvent-excluded ("Connolly") surfaces, and isosurfaces. The isosurface in Fig. 5 appears to show the electrostatic potential, with blue colors being negative and red/yellow (near the metal) positive (there is no absolute convention of coloring, and red/positive, blue/negative are often reversed). Opaque isosurfaces do not allow the atoms to be seen and identified and it is not easy to deduce them. Because of this, isosurfaces are often drawn with a degree of transparency.

Technology

Early interactive molecular computer graphics systems were vector graphics machines, which used stroke-writing vector monitors, sometimes even oscilloscopes. The electron beam does not sweep left-and-right as in a raster display. The display hardware followed a sequential list of digital drawing instructions (the display list), directly drawing at an angle one stroke for each molecular bond. When the list was complete, drawing would begin again from the top of the list, so if the list was long (a large number of molecular bonds), the display would flicker heavily. Later vector displays could rotate complex structures with smooth motion, since the orientation of all of the coordinates in the display list could be changed by loading just a few numbers into rotation registers in the display unit, and the display unit would multiply all coordinates in the display list by the contents of these registers as the picture was drawn.

The early black-and white vector displays could somewhat distinguish for example a molecular structure from its surrounding electron density map for crystallographic structure solution work by drawing the molecule brighter than the map. Color display makes them easier to tell apart. During the 1970s two-color stroke-writing Penetron tubes were available, but not used in molecular computer graphics systems. In about 1980 Evans & Sutherland made the first practical full-color vector displays for molecular graphics, typically attached to an E&S PS-300 display. This early color tube was expensive, because it was originally engineered to withstand the shaking of a flight-simulator motion base.

Color raster graphics display of molecular models began around 1978 as seen in this paper by Porter on spherical shading of atomic models. Early raster molecular graphics systems displayed static images that could take around a minute to generate. Dynamically rotating color raster molecular display phased in during 1982-1985 with the introduction of the Ikonas programmable raster display.

Molecular graphics has always pushed the limits of display technology, and has seen a number of cycles of integration and separation of compute-host and display. Early systems like Project MAC were bespoke and unique, but in the 1970s the MMS-X and similar systems used (relatively) low-cost terminals, such as the Tektronix 4014 series, often over dial-up lines to multi-user hosts. The devices could only display static pictures but were able to evangelize MG. In the late 1970s, it was possible for departments (such as crystallography) to afford their own hosts (e.g., PDP-11) and to attach a display (such as Evans & Sutherland's MPS) directly to the bus. The display list was kept on the host, and interactivity was good since updates were rapidly reflected in the display—at the cost of reducing most machines to a single-user system.

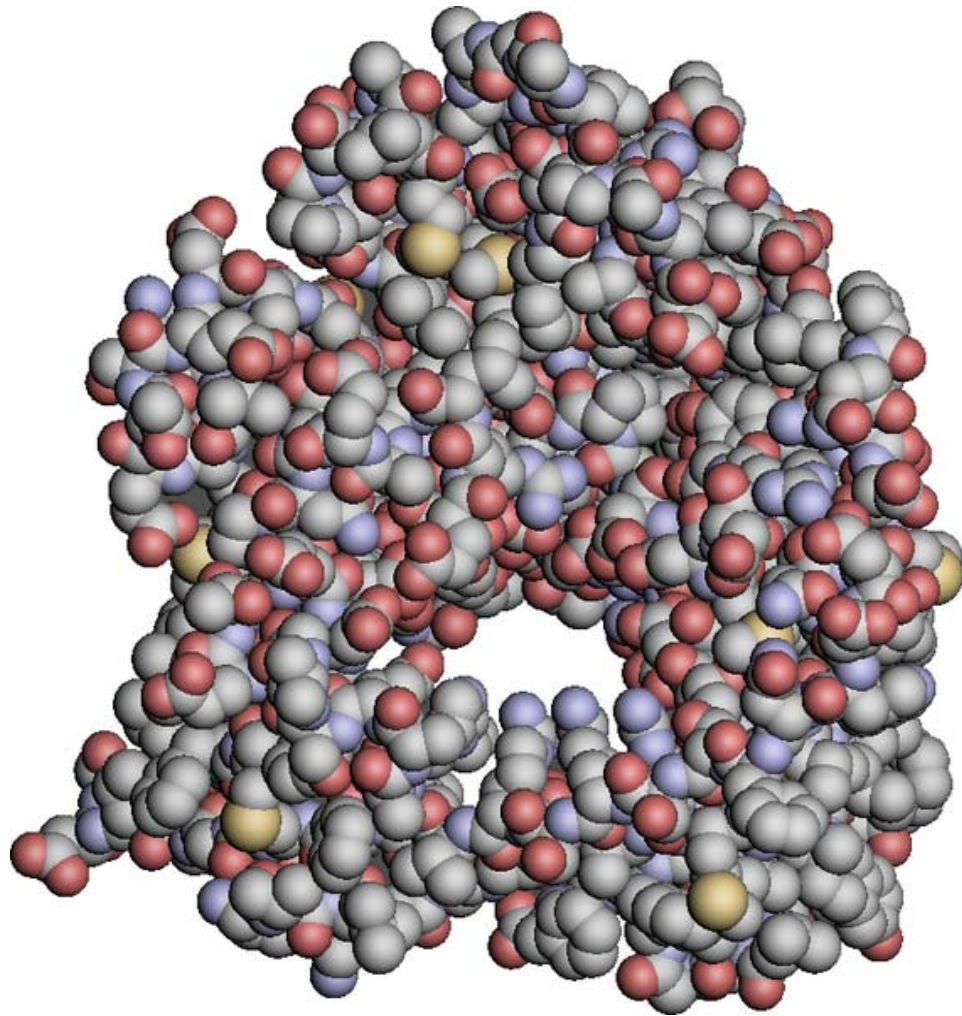
In the early 1980s, Evans & Sutherland (E&S) decoupled their PS300 display, which contained its own display information transformable through a dataflow architecture. Complex graphical objects could be downloaded over a serial line (e.g. 9600 baud) and then manipulated without impact on the host. The architecture was excellent for high performance display but very inconvenient for domain-specific calculations, such as

electron-density fitting and energy calculations. Many crystallographers and modellers spent arduous months trying to fit such activities into this architecture.

The benefits for MG were considerable, but by the later 1980s, UNIX workstations such as Sun-3 with raster graphics (initially at a resolution of 256 by 256) had started to appear. Computer-assisted drug design in particular required raster graphics for the display of computed properties such as atomic charge and electrostatic potential. Although E&S had a high-end range of raster graphics (primarily aimed at the aerospace industry) they failed to respond to the low-end market challenge where single users, rather than engineering departments, bought workstations. As a result the market for MG displays passed to Silicon Graphics, coupled with the development of minisupercomputers (e.g., CONVEX and Alliant) which were affordable for well-supported MG laboratories. Silicon Graphics provided a graphics language, IrisGL, which was easier to use and more productive than the PS300 architecture. Commercial companies (e.g., Biosym, Polygen/MSI) ported their code to Silicon Graphics, and by the early 1990s, this was the "industry standard". Dial boxes were often used as control devices.

Stereoscopic displays were developed based on liquid crystal polarized spectacles, and while this had been very expensive on the PS300, it now became a commodity item. A common alternative was to add a polarizable screen to the front of the display and to provide viewers with extremely cheap spectacles with orthogonal polarization for separate eyes. With projectors such as Barco, it was possible to project stereoscopic display onto special silvered screens and supply an audience of hundreds with spectacles. In this way molecular graphics became universally known within large sectors of chemical and biochemical science, especially in the pharmaceutical industry. Because the backgrounds of many displays were black by default, it was common for modelling sessions and lectures to be held with almost all lighting turned off.

In the last decade almost all of this technology has become commoditized. IrisGL evolved to OpenGL so that molecular graphics can be run on any machine. In 1992, Roger Sayle released his RasMol program into the public domain. RasMol contained a very high-performance molecular renderer that ran on Unix/X Window, and Sayle later ported this to the Windows and Macintosh platforms. The Richardsons developed kinemages and the Mage software, which was also multi-platform. By specifying the chemical MIME type, molecular models could be served over the Internet, so that for the first time MG could be distributed at zero cost regardless of platform. In 1995, Birkbeck College's crystallography department used this to run "Principles of Protein Structure", the first multimedia course on the Internet, which reached 100 to 200 scientists.



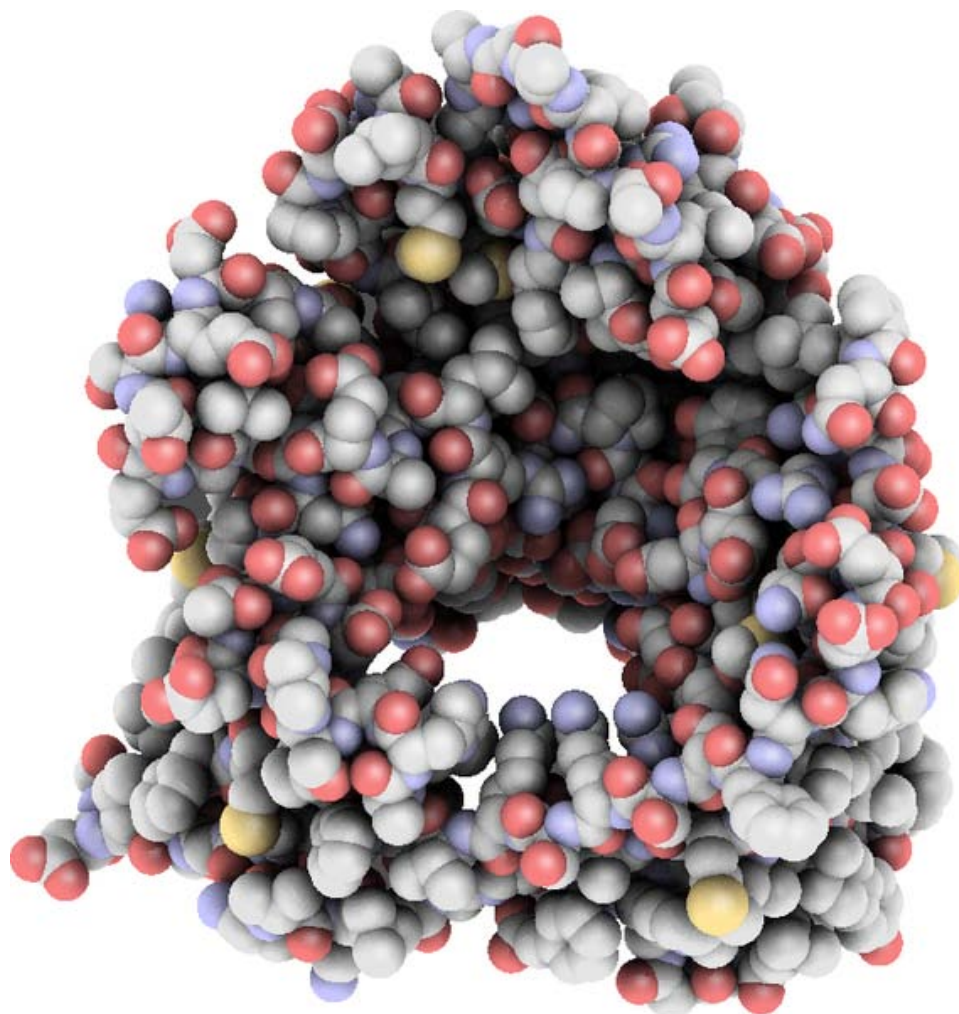


Fig. 6. A molecule of Porin (protein) shown without ambient occlusion (above) and with (below). Advanced rendering effects can improve the comprehension of the 3D shape of a molecule.

MG continues to see innovation that balances technology and art, and currently zero-cost or open source programs such as PyMOL and Jmol have very wide use and acceptance.

Recently the wide spread diffusion of advanced graphics hardware, has improved the rendering capabilities of the visualization tools. The capabilities of current shading languages allow the inclusion of advanced graphic effects (like ambient occlusion, cast shadows and non-photorealistic rendering techniques) in the interactive visualization of molecules. These graphic effects, beside being eye candy, can improve the comprehension of the three dimensional shapes of the molecules. An example of the effects that can be achieved exploiting recent graphics hardware can be seen in the simple open source visualization system QuteMol.

Algorithms

Reference frames

Drawing molecules requires a transformation between molecular coordinates (usually, but not always, in Angstrom units) and the screen. Because many molecules are chiral it is essential that the handedness of the system (almost always right-handed) is preserved. In molecular graphics the origin (0, 0) is usually at the lower left, while in many computer systems the origin is at top left. If the z-coordinate is out of the screen (towards the viewer) the molecule will be referred to right-handed axes, while the screen display will be left-handed.

Molecular transformations normally require:

- scaling of the display (but not the molecule).
- translations of the molecule and objects on the screen.
- rotations about points and lines.

Conformational changes (e.g. rotations about bonds) require rotation of one part of the molecule relative to another. The programmer must decide whether a transformation on the screen reflects a change of view or a change in the molecule or its reference frame.

Simple

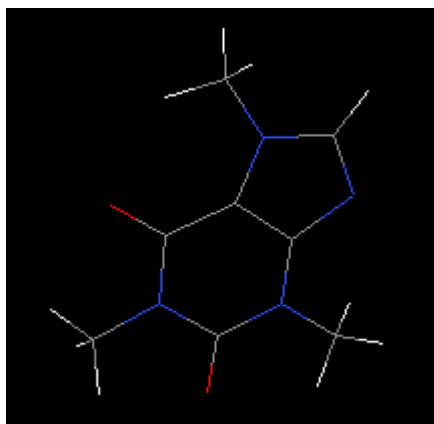


Fig. 7. Stick model of caffeine drawn in Jmol.

In early displays only vectors could be drawn e.g. (Fig. 7) which are easy to draw because no rendering or hidden surface removal is required.

On vector machines the lines would be smooth but on raster devices Bresenham's algorithm is used (note the "jaggies" on some of the bonds, which can be largely removed with antialiasing software.)

Atoms can be drawn as circles, but these should be sorted so that those with the largest z-coordinates (nearest the screen) are drawn last. Although imperfect, this often gives a reasonably attractive display. Other simple tricks which do not include hidden surface algorithms are:

- colouring each end of a bond with the same colour as the atom to which it is attached (Fig. 7).
- drawing less than the whole length of the bond (e.g. 10%-90%) to simulate the bond sticking out of a circle.
- adding a small offset white circle within the circle for an atom to simulate reflection.

Typical pseudocode for creating Fig. 7 (to fit the molecule exactly to the screen):

```
// assume:
// atoms with x, y, z coordinates (Angstrom) and elementSymbol
// bonds with pointers/references to atoms at ends
// table of colours for elementTypes
// find limits of molecule in molecule coordinates as xMin, yMin, xMax,
yMax
scale = min(xScreenMax/(xMax-xMin), yScreenMax/(yMax-yMin))
xOffset = -xMin * scale; yOffset = -yMin * scale
for (bond in $bonds) {
  atom0 = bond.getAtom(0)
  atom1 = bond.getAtom(1)
  x0 = xOffset+atom0.getX()*scale; y0 = yOffset+atom0.getY()*scale //
(1)
  x1 = xOffset+atom1.getX()*scale; y1 = yOffset+atom1.getY()*scale //
(2)
  x1 = atom1.getX(); y1 = atom1.getY()
  xMid = (x0 + x1) /2; yMid = (y0 + y1) /2;
  colour0 = ColourTable.getColour(atom0.getSymbol())
  drawLine (colour0, x0, y0, xMid, yMid)
  colour1 = ColourTable.getColour(atom1.getSymbol())
  drawLine (colour1, x1, y1, xMid, yMid)
}
```

Note that this assumes the origin is in the bottom left corner of the screen, with Y up the screen. Many graphics systems have the origin at the top left, with Y down the screen. In this case the lines (1) and (2) should have the y coordinate generation as:

```
y0 = yScreenMax - (yOffset+atom0.getY()*scale) // (1)
y1 = yScreenMax - (yOffset+atom1.getY()*scale) // (2)
```

Changes of this sort change the handedness of the axes so it is easy to reverse the chirality of the displayed molecule unless care is taken.

Advanced

For greater realism and better comprehension of the 3D structure of a molecule many computer graphics algorithms can be used. For many years molecular graphics has

stressed the capabilities of graphics hardware and has required hardware-specific approaches. With the increasing power of machines on the desktop, portability is more important and programs such as Jmol have advanced algorithms that do not rely on hardware. On the other hand recent graphics hardware is able to interactively render very complex molecule shapes with a quality that would not be possible with standard software techniques.

Chronology

Formative Era: 1960-1976

The formative years of molecular graphics can be summarized by

- Hard-copy molecular graphic in heavy, routine use by chemists (ORTEP, PLUTO, etc.);
- Interactive molecular graphics prototype systems producing little or no publishable chemistry.

The problems for interactive molecular graphics were

- The small address space of laboratory computers (12-bit and 16-bit) limited the size of programs and data that could be easily handled to below the size of most real problems;
- Interactive computer graphics displays were only just fast enough to display simple molecules, below the size of most real problems.

Developer(s)	Approximate date	Technology	Comments
Crystallographers	< 1960	Hand-drawn	Crystal structures, with hidden atom and bond removal. Often clinographic projections.
Johnson, Motherwell	ca 1970	Pen plotter	ORTEP, PLUTO. Very widely deployed for publishing crystal structures.
Cyrus Levinthal, Bob Langridge, Ward, Stots	1966	Project MAC display system, two-degree of freedom, spring-return velocity joystick for rotating the image.	First protein display on screen. System for interactively building protein structures.
Barry	1969	Link 300 computer with a dual trace oscilloscope display.	Interactive molecular structure viewing system. Early examples of dynamic rotation, intensity depth-cueing, and side-by-

			side stereo. Early use of the small angle approximations ($a = \sin a$, $1 = \cos a$) to speed up graphical rotation calculations.
Ortony	1971	Designed a stereo viewer (British patent appl. 13844/70) for molecular computer graphics.	Horizontal two-way (half-silvered) mirror combines images drawn on the upper and lower halves of a CRT. Crossed polarizers isolate the images to each eye.
Ortony	1971	Light pen, knob.	Interactive molecular structure viewing system. Select bond by turning another knob until desired bond lights up in sequence, a technique later used on the MMS-4 system below, or by picking with the light pen. Points in space are specified with a 3-D "bug" under dynamic control.
Barry, Graesser, Marshall	1971	CHEMAST: LINK 300 computer driving an oscilloscope. Two-axis joystick, similar to one used later by GRIP-75 (below).	Interactive molecular structure viewing system. Structures dynamically rotated using the joystick.
Tountas and Katz	1971	Adage AGT/50 display	Interactive molecular structure viewing system. Mathematics of nested rotation and for laboratory-space rotation.
Perkins, Piper, Tattam, White	1971	Honeywell DDP 516 computer, EAL TR48 analog computer, Lanelec oscilloscope, 7 linear potentiometers. Stereo.	Interactive molecular structure viewing system.
Wright	1972	GRIP-71 at UNC-CH: IBM System/360 Model 40 time-shared computer, IBM 2250 display, buttons, light pen, keyboard.	Discrete manipulation and energy relaxation of protein structures. Program code became the foundation of the GRIP-75 system below.

Barry and North	1972	Oxford Univ.: Ferranti Argus 500 computer, Ferranti model 30 display, keyboard, track ball, one knob. Stereo.	Prototype large-molecule crystallographic structure solution system. Track ball rotates a bond, knob brightens the molecule vs. electron density map.
North, Ford, Watson	Early 1970s	Leeds Univ.: DEC PDP-11/40 computer, Hewlett-Packard display. 16 knobs, keyboard, spring-return joystick. Stereo.	Prototype large-molecule crystallographic structure solution system. Six knobs rotate and translate a small molecule.
Barry, Bosshard, Ellis, Marshall, Fritch, Jacobi	1974	MMS-4: Washington Univ. at St. Louis, Link 300 computer and an LDS-1 / Link 300 display, custom display modules. Rotation joystick, knobs. Stereo.	Prototype large-molecule crystallographic structure solution system. Select bond to rotate by turning another knob until desired bond lights up in sequence.
Cohen and Feldmann	1974	DEC PDP-10 computer, Adage display, push buttons, keyboard, knobs	Prototype large-molecule crystallographic structure solution system.
Stellman	1975	Princeton: PDP-10 computer, LDS-1 display, knobs	Prototype large-molecule crystallographic structure solution system. Electron density map not shown; instead an "H Factor" figure of merit is updated as the molecular structure is manipulated.
Collins, Cotton, Hazen, Meyer, Morimoto	1975	CRYSNET , Texas A&M Univ. DEC PDP-11/40 computer, Vector General Series 3 display, knobs, keyboard. Stereo.	Prototype large-molecule crystallographic structure solution system. Variety of viewing modes: rocking, spinning, and several stereo display modes.
Cornelius and Kraut	1976 (approx.)	Univ. of Calif. at San Diego: DEC PDP-11/40 emulator (CalData 135), Evans and Sutherland Picture System display, keyboard, 6 knobs. Stereo.	Prototype large-molecule crystallographic structure solution system.

(Yale Univ.)	1976 (approx.)	PIGS: DEC PDP-11/70 computer, Evans and Sutherland Picture System 2 display, data tablet, knobs.	Prototype large-molecule crystallographic structure solution system. The tablet was used for most interactions.
Feldmann and Porter	1976	NIH: DEC PDP—11/70 computer. Evans and Sutherland Picture System 2 display, knobs. Stereo.	Interactive molecular structure viewing system. Intended to display interactively molecular data from the AMSOM – Atlas of Macromolecular Structure on Microfiche.
Rosenberger et al.	1976	MMS-X: Washington Univ. at St. Louis, TI 980B computer, Hewlett-Packard 1321A display, Beehive video terminal, custom display modules, pair of 3-D spring-return joysticks, knobs.	Prototype large-molecule crystallographic structure solution system. Successor to the MMS-4 system above. The 3-D spring-return joysticks either translate and rotate the molecular structure for viewing or a molecular substructure for fitting, mode controlled by a toggle switch.

A few paragraphs each on many of the display and crystallographic systems above are in a review chapter in Lipscomb.

MMS-4 and MMS-X

In a noteworthy attempt to overcome the low speed of graphics displays of the time took place at Washington University in St. Louis. Dave Barry's group attempted to leapfrog the state of the art in graphics displays by making custom display hardware to display images complex enough for large-molecule crystallographic structure solution, fitting molecules to their electron-density maps. The MMS-4 (table above) display modules were slow and expensive, so a second generation of modules was produced for the MMS-X (table above) system. The 16-bit computer initially driving the display was a limitation that may have delayed productive crystallography until the years of the "Mature Era" below.

Mature Era: 1977-present

Developer(s)	Approximate date	Technology	Comments
Britton, Lipscomb, Pique, Wright, Brooks	1977	GRIP-75 at UNC-CH: Time-shared IBM System/360 Model 75 computer, DEC PDP 11/45 computer, Vector	First large-molecule crystallographic structure solution.

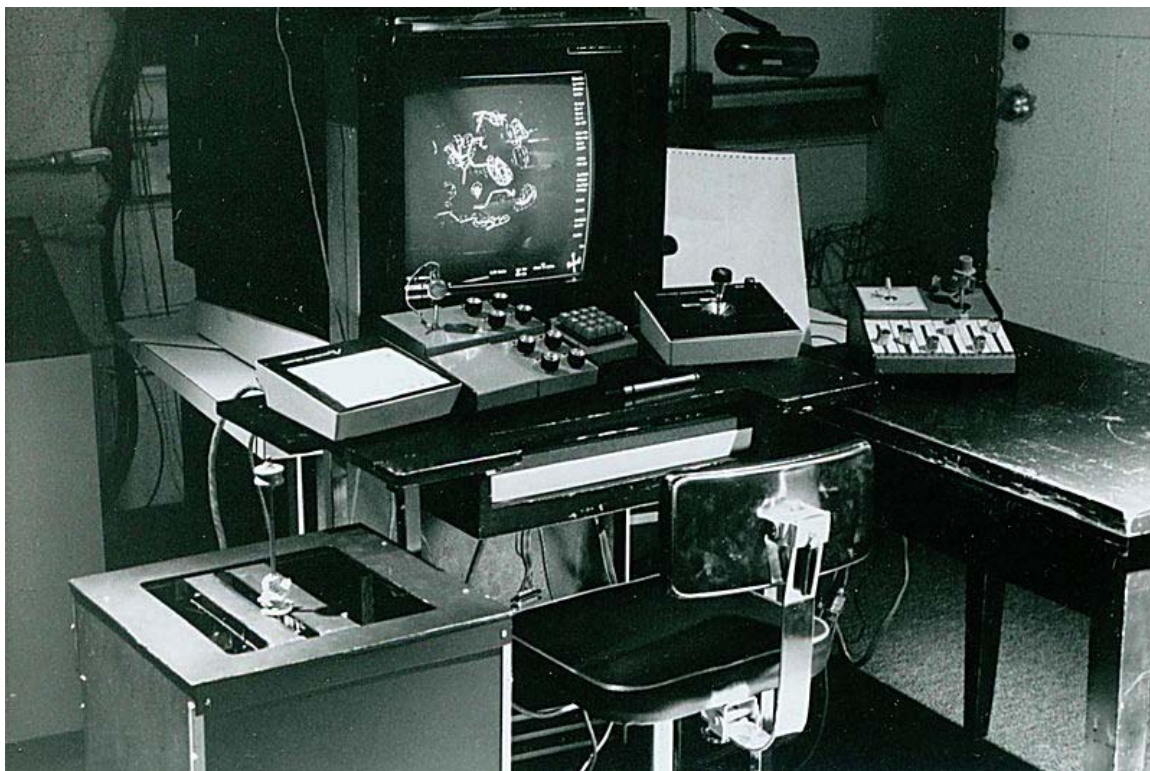
			General Series 3 display, 3-D movement box from A.M. Noll and 3-D spring return joystick for substructure manipulation, Measurement Systems nested joystick, knobs, sliders, buttons, keyboard, light pen.	
Jones	1978		FRODO and RING Max Plank Inst., Germany, RING: DEC PDP-11/40 and Siemens 4004 computers, Vector General 3404 display, 6 knobs.	Large-molecule crystallographic structure solution. FRODO may have run on a DEC VAX-780 as a follow-on to RING.
Diamond	1978		Bilder Cambridge, England, DEC PDP-11/50 computer, Evans and Sutherland Picture System display, tablet.	Large-molecule crystallographic structure solution. All input is by data tablet. Molecular structures built on-line with ideal geometry. Later passes stretch bonds with idealization.
Langridge, White, Marshall	Late 1970s		Departmental systems (PDP-11, Tektronix displays or DEC-VT11, e.g. MMS-X)	Mixture of commodity computing with early displays.
Davies, Hubbard	Mid-1980s		CHEM-X, HYDRA	Laboratory systems with multicolor, raster and vector devices (Sigmex, PS300).
Biosym, Tripos, Polygen	Mid-1980s		PS300 and lower cost dumb terminals (VT200, SIGMEX)	Commercial integrated modelling and display packages.
Silicon Graphics, Sun	Late 1980s		IRIS GL (UNIX) workstations	Commodity-priced single-user workstations with stereoscopic display.
EMBL - WHAT IF	1989, 2000		Machine independent	Nearly free, multifunctional, still fully supported, many free servers based on it
Sayle, Richardson	1992, 1993		RasMol, Kinemage	Platform-independent MG.
MDL (van Vliet,	1995–1998		Chime	proprietary C++ ; free

Maffett, Adler, Holt)			browser plugin for Mac (OS9) and PCs
ChemAxon	1998-	MarvinSketch & MarvinView. MarvinSpace (2005)	proprietary Java applet or stand-alone application.
Community efforts	2000-	Jmol, PyMol, Avogadro, PDB	Open-source Java applet or stand-alone application.
NOCH	2002-	NOC	Powerful and open source code molecular structure explorer
LION Bioscience / EMBL	2004-	SRS 3D	Free, open-source system based on Java3D. Integrates 3D structures with sequence and feature data (domains, SNPs, etc.).
San Diego Supercomputer Center	2006-	Sirius	Free for academic/non-profit institutions
Weizmann Institute of Science - Community efforts	2008-	Proteopedia	Collaborative

Productive use of interactive molecular graphics became possible in the Mature Era largely by two hardware advances:

- 24-bit and 32-bit computers became affordable to university departments, increasing the size of programs and data that could be easily handled to the size of many real problems. Early examples include GRIP-75, which used in part a time-shared IBM System/360 Model 75, and FRODO, which used a DEC VAX-780.
- Interactive computer graphics displays became fast enough to easily display complex molecules and electron-density maps. The first such displays widely used were the Vector General Series 3 and the Evans and Sutherland Picture System 2, MultiPicture System, and PS-300

First large-molecule crystallographic computer graphics structure solutions



GRIP-75 molecular graphics system console, Univ. of N. Carolina, Chapel Hill

Before computer graphics could do the job, mechanical methods were used to fit large molecules to their electron density maps. Contour circles around high electron density were drawn on large plastic sheets. Bingo chips were placed on the plastic sheets where atoms were interpreted to be in the case of the solution of carboxypeptidase A. This was superseded by the Richards Box in which an adjustable brass Kendrew molecular model was placed front of a 2-way mirror, behind which were plastic sheets of the electron density map. This optically superimposed the molecular model and the electron density map. The model was moved to within the contour lines of the superimposed map. Then, atomic coordinates were recorded using a plumb bob and a meter stick. These mechanical methods are no longer used.

The first large molecule whose atomic structure was *partly* determined on a molecular computer graphics system was Transfer RNA by Sung-Hou Kim's team in 1976. They used the GRIP-75 computer graphics system after initial fitting on a mechanical Richards Box.

The first large molecule whose atomic structure was *entirely* determined on a molecular computer graphics system was neurotoxin A from venom of the Philippines sea snake, by Tsernoglou, Petsko, and Tu, using GRIP-75, with a statement of being first in 1977.

The Richardson group, using GRIP-75, published partial atomic structure results of the protein superoxide dismutase in 1977.

GRIP-75 was gradually superseded by systems that crystallographers could have in their own labs (FRODO, RING, Builder, MMS-X, etc.). There was only one GRIP-75 system built, and crystallographers had to travel to use it, and after solving about two dozen molecular structures it was retired.

All of these early systems were sometimes called "Electronic Richards Boxes" because they worked likewise. Crystallographers manually moved parts of the molecule (using joysticks and knobs) to within the contour lines of the electronically superimposed map.

Nowadays fitting of the molecular structure to the electron density map is largely automated by algorithms with computer graphics a guide to the process. An example is the XtalView XFit program.

Chapter 4

Molecular Machine

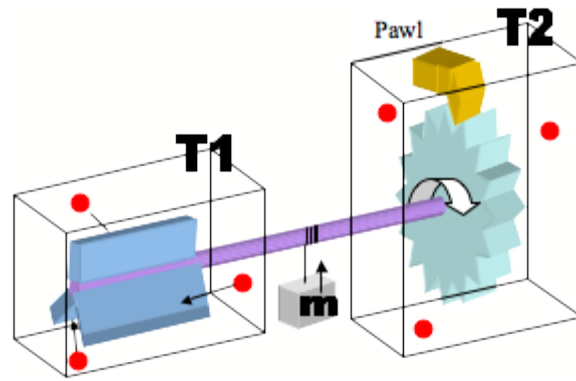
A **molecular machine**, or nanomachine, is defined as a discrete number of molecular components which perform mechanical-like movements (output) in response to specific stimuli (input). More generally, the expression is often applied to molecules that simply mimic functions which occur at the macroscopic level. The term is also common in nanotechnology, and a number of highly complex molecular machines have been proposed towards the goal of constructing a molecular assembler. Molecular machines can be divided into two broad categories: synthetic and biological.

Molecular systems capable of shifting a chemical or mechanical process away from equilibrium represent a potentially important branch of chemistry and nanotechnology. As the gradient generated from this process is able to perform useful work, by definition, these types of systems are examples of molecular machinery.

Historical Insight and Studies

There are two thought experiments that form the historical basis for molecular machines: Maxwell's demon and Feynman's Ratchet (or Brownian ratchet). Maxwell's Demon is well described elsewhere, and a slightly different interpretation of Richard Feynman's ratchet is given here.

Imagine a very small system (seen below) of two paddles or gears connected by a rigid axle and that it is possible to keep these two paddles at two different temperatures. One of the gears (at T_2) has a pawl that is rectifying the system motion, and therefore, the axle can only move in a clockwise rotation, and in doing so, it could lift a weight (m) upward upon ratcheting. Now imagine if the paddle in box T_1 was in a much hotter environment than the gear in box T_2 ; it would be expected that the kinetic energy of the gas molecules (red circles) hitting the paddle in T_1 would be much higher than the gas molecules hitting the gear at T_2 . Therefore, with lower kinetic energy of the gases in T_2 , there would be very little resistance from the molecules on colliding with the gear in the statistically opposite direction. Further, the ratcheting would allow for directionality, and slowly over time, the axle would rotate and ratchet, lifting the weight (m).



Schematic figure of Feynman's Ratchet

As described, this system may seem like a perpetual motion machine; however, the key ingredient is the heat gradient within the system. This ratchet does not threaten the second law of thermodynamics, because this temperature gradient must be maintained by some external means. Brownian motion of the gas particles provides the power to the machine, and the temperature gradient allows the machine to drive the system cyclically away from equilibrium. In Feynman's ratchet, random Brownian motion is not fought against, but instead, harnessed and rectified. Unfortunately, temperature gradients cannot be maintained over molecular scale distances because of molecular vibration redistributing the energy to other parts of the molecule. Furthermore, despite Feynman's machine doing useful work in lifting the mass, using Brownian motion to power a molecular level machine does not provide any insight on how that power (or potential energy of the lifted weight, m) can be used to perform nanoscale tasks.

Modern Insights and Studies

Unlike macroscopic motion, molecular systems are constantly undergoing significant dynamic motions subject to the laws of Brownian mechanics (or Brownian motion), and as such, harnessing molecular motion is a far more difficult process. At the macroscopic level, many machines operate in the gas phase, and often, air resistance is neglected, as it is insignificant, but analogously for a molecular system in a Brownian environment, molecular motion is similar "to walking in a hurricane, or swimming in molasses." The phenomenon of Brownian motion (observed by Robert Brown (botanist), 1827) was later explained by Albert Einstein in 1905. Einstein found that Brownian motion is a consequence of scale and not the nature of the surroundings. As long as thermal energy is applied to a molecule, it will undergo Brownian motion with the kinetic energy appropriate to that temperature. Therefore, like Feynman's strategy, when designing a molecular machine, it seems sensible to utilize Brownian motion rather than attempt to fight against it.

Like macroscopic machines, molecular machines typically have movable parts. However, while everyday macroscopic machines may provide inspiration for molecular machines, it is misleading to draw analogies between their design strategy; the dynamics of large and small length scales are simply too different. Harnessing Brownian motion and

making molecular level machines is regulated by the second law of thermodynamics, with its often counter-intuitive consequences, and as such, we need another inspiration.

Although it is a challenging process to harness Brownian motion, nature has provided us with several blueprints for molecular motion performing useful work. Nature has created many useful structures for compartmentalizing molecular systems, hence creating distinct non-equilibrium distributions; the cell membrane is an excellent example. Lipophilic barriers make use of a number of different mechanisms to power motion from one compartment to another.

Examples of molecular machines

From a synthetic perspective, there are two important types of molecular machines: molecular switches (or shuttles) and molecular motors. The major difference between the two systems is that a switch influences a system as a function of state, whereas a motor influences a system as function of trajectory. A switch (or shuttle) may appear to undergo translational motion, but returning a switch to its original position undoes any mechanical effect and liberates energy to the system. Furthermore, switches cannot use chemical energy to repetitively and progressively drive a system away from equilibrium where a motor can.

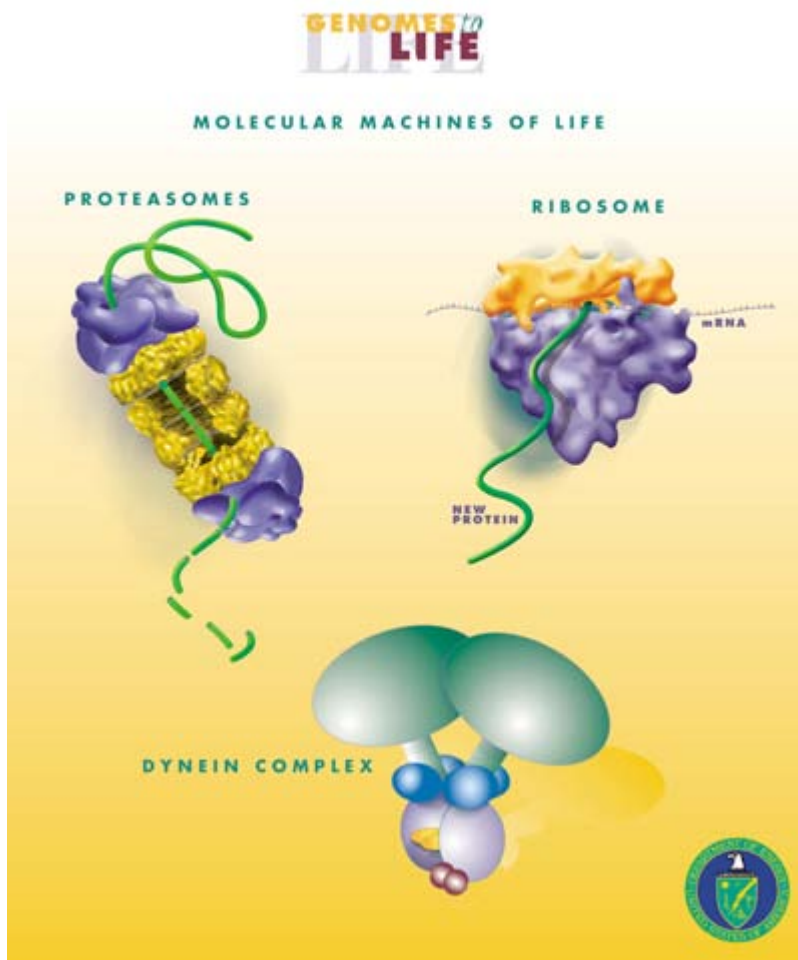
Synthetic

A wide variety of rather simple molecular machines have been synthesized by chemists. They can consist of a single molecule; however, they are often constructed for mechanically-interlocked molecular architectures, such as rotaxanes and catenanes.

- Molecular motors are molecules that are capable of unidirectional rotation motion powered by external energy input. A number of molecular machines have been synthesized powered by light or reaction with other molecules.
- A molecular propeller is a molecule that can propel fluids when rotated, due to its special shape that is designed in analogy to macroscopic propellers. It has several molecular-scale blades attached at a certain pitch angle around the circumference of a nanoscale shaft.
- A molecular switch is a molecule that can be reversibly shifted between two or more stable states. The molecules may be shifted between the states in response to changes in e.g. pH, light, temperature, an electrical current, microenvironment, or the presence of a ligand.
- A molecular shuttle is a molecule capable of shuttling molecules or ions from one location to another. A common molecular shuttle consists of a rotaxane where the macrocycle can move between two sites or stations along the dumbbell backbone.

- Molecular tweezers are host molecules capable of holding items between its two arms. The open cavity of the molecular tweezers binds items using non-covalent bonding including hydrogen bonding, metal coordination, hydrophobic forces, van der Waals forces, π - π interactions, and/or electrostatic effects. Examples of molecular tweezers have been reported that are constructed from DNA and are considered DNA machines.
- A molecular sensor is a molecule that interacts with an analyte to produce a detectable change. Molecular sensors combine molecular recognition with some form of reporter, so the presence of the item can be observed.
- A molecular logic gate is a molecule that performs a logical operation on one or more logic inputs and produces a single logic output. Unlike a molecular sensor, the molecular logic gate will only output when a particular combination of inputs are present.

Biological



Some biological molecular machines

The most complex molecular machines are found within cells. These include motor proteins, such as myosin, which is responsible for muscle contraction, kinesin, which moves cargo inside cells away from the nucleus along microtubules, and dynein, which produces the axonemal beating of motile cilia and flagella. These proteins and their nanoscale dynamics are far more complex than any molecular machines that have yet been artificially constructed.

The detailed mechanism of ciliary motility has been described by Satir in a 2008 review article. A high-level-abstraction summary is that, "[i]n effect, the [motile cilium] is a nanomachine composed of perhaps over 600 proteins in molecular complexes, many of which also function independently as nanomachines."

Theoretical

The construction of more complex molecular machines is an active area of theoretical research. A number of molecules, such as molecular propellers, have been designed, although experimental studies of these molecules are inhibited by the lack of methods to construct these molecules. These complex molecular machines form the basis of areas of nanotechnology, including molecular assembler.

Chapter 5

Molecular Assembler

A **molecular assembler** is a "proposed device able to guide chemical reactions by positioning reactive molecules with atomic precision", as defined by K. Eric Drexler. Some biological molecules such as ribosomes fit this definition, since they receive instructions from messenger RNA and then assemble specific sequences of amino acids to construct protein molecules. However, the term "molecular assembler" usually refers to theoretical human-made devices.

Beginning in 2007, the British Engineering and Physical Sciences Research Council funds development of ribosome-like molecular assemblers. Clearly, molecular assemblers are possible in this limited sense. A technology roadmap project, led by the Battelle Memorial Institute and hosted by several U.S. National Laboratories has explored a range of atomically precise fabrication technologies, including both early-generation and longer-term prospects for programmable molecular assembly; the report was released in December, 2007.

Likewise, the term "molecular assembler" has been used in science fiction and popular culture to refer to a wide range of fantastic atom-manipulating nanomachines, many of which may be physically impossible in reality. Much of the controversy regarding "molecular assemblers" results from the confusion in the use of the name for both technical concepts and popular fantasies. In 1992, Drexler introduced the related but better-understood term "molecular manufacturing," which he defined as the programmed "chemical synthesis of complex structures by mechanically positioning reactive molecules, not by manipulating individual atoms."

Here we, mostly discusses "molecular assemblers" in the popular sense. These include hypothetical machines that manipulate individual atoms and machines with organism-like self-replicating abilities, mobility, ability to consume food, and so forth. These are quite different from devices that merely (as defined above) "guide chemical reactions by positioning reactive molecules with atomic precision".

Because synthetic molecular assemblers have never been constructed and because of the confusion regarding the meaning of the term, there has been much controversy as to whether "molecular assemblers" are possible or simply science fiction. Confusion and

controversy also stem from their classification as nanotechnology, which is an active area of laboratory research which has already been applied to the production of real products; however, there had been, until recently, no research efforts into the actual construction of "molecular assemblers". A primary criticism of the computational research into products of advanced "molecular assemblers" is that the structures investigated are impossible to make today.

Nanofactories

A **nanofactory** is a proposed system in which nanomachines (resembling molecular assemblers, or industrial robot arms) would combine reactive molecules via mechanosynthesis to build larger atomically precise parts. These, in turn, would be assembled by positioning mechanisms of assorted sizes to build macroscopic (visible) but still atomically-precise products.

A typical nanofactory would fit in a desktop box, in the vision of K. Eric Drexler published in *Nanosystems: Molecular Machinery, Manufacturing and Computation* (1992), a notable work of "exploratory engineering". During the last decade, others have extended the nanofactory concept, including an analysis of nanofactory convergent assembly by Ralph Merkle, a systems design of a replicating nanofactory architecture by J. Storrs Hall, Forrest Bishop's "Universal Assembler", the patented exponential assembly process by Zyvex, and a top-level systems design for a 'primitive nanofactory' by Chris Phoenix (Director of Research at the Center for Responsible Nanotechnology). All of these nanofactory designs (and more) are summarized in Chapter 4 of *Kinematic Self-Replicating Machines* (2004) by Robert Freitas and Ralph Merkle. The Nanofactory Collaboration, founded by Robert Freitas and Ralph Merkle in 2000, is a focused ongoing effort involving 23 researchers from 10 organizations and 4 countries that is developing a practical research agenda specifically aimed at positionally-controlled diamond mechanosynthesis and diamondoid nanofactory development.

In 2005, a computer-animated short film of the nanofactory concept was produced by John Burch, in collaboration with Drexler. Such visions have been the subject of much debate, on several intellectual levels. No one has discovered an insurmountable problem with the underlying theories and no one has proved that the theories can be translated into practice. However, the debate continues, with some of it being summarized in the Molecular nanotechnology article.

If nanofactories could be built, severe disruption to the world economy would be one of many possible negative impacts, though it could be argued that this disruption would have little negative effect if everyone had such nanofactories. Great benefits also would be anticipated. Various works of science fiction have explored these and similar concepts. The potential for such devices was part of the mandate of a major UK study led by mechanical engineering professor Dame Ann Dowling. The report is now complete.

Self-replication

"Molecular assemblers" have been confused with self-replicating machines. To produce a practical quantity of a desired product, the nanoscale size of a typical science fiction universal molecular assembler requires an extremely large number of such devices. However, a single such theoretical molecular assembler might be programmed to self-replicate, constructing many copies of itself. This would allow an exponential rate of production. Then after sufficient quantities of the molecular assemblers were available, they would then be re-programmed for production of the desired product. However, if self-replication of molecular assemblers were not restrained then it might lead to competition with naturally occurring organisms. This has been called ecophagy or the grey goo problem.

One method to building molecular assemblers is to mimic evolutionary processes employed by biological systems. Biological evolution proceeds by random variation combined with culling of the less-successful variants and reproduction of the more-successful variants. Production of complex molecular assemblers might be evolved from simpler systems since "A complex system that works is invariably found to have evolved from a simple system that worked. . . . A complex system designed from scratch never works and can not be patched up to make it work. You have to start over, beginning with a system that works." However, most published safety guidelines include "recommendations against developing ... replicator designs which permit surviving mutation or undergoing evolution".

Most assembler designs keep the "source code" external to the physical assembler. At each step of a manufacturing process, that step is read from an ordinary computer file and "broadcast" to all the assemblers. If any assembler gets out of range of that computer, or when the link between that computer and the assemblers is broken, or when that computer is unplugged, the assemblers stop replicating. Such a "broadcast architecture" is one of the safety features recommended by the "Foresight Guidelines on Molecular Nanotechnology", and a map of the 137-dimensional replicator design space recently published by Freitas and Merkle provides numerous practical methods by which replicators can be safely controlled by good design.

Drexler and Smalley debate

One of the most outspoken critics of some concepts of "molecular assemblers" was Professor Richard Smalley (1943-2005) who won the Nobel prize for his contributions to the field of nanotechnology. Smalley believed that such assemblers were not physically possible and introduced scientific objections to them. His two principal technical objections were termed the "fat fingers problem" and the "sticky fingers problem". He believed these would exclude the possibility of "molecular assemblers" that worked by precision picking and placing of individual atoms. Drexler and coworkers responded to these two issues in a 2001 publication.

Smalley also believed that Drexler's speculations about apocalyptic dangers of self-replicating machines that have been equated with "molecular assemblers" would threaten the public support for development of nanotechnology. To address the debate between Drexler and Smalley regarding molecular assemblers *Chemical & Engineering News* published a point-counterpoint consisting of an exchange of letters that addressed the issues.

Regulation

Speculation on the power of systems that have been called "molecular assemblers" has sparked a wider political discussion on the implication of nanotechnology. This is in part due to the fact that nanotechnology is a very broad term and could include "molecular assemblers." Discussion of the possible implications of fantastic molecular assemblers has prompted calls for regulation of current and future nanotechnology. There are very real concerns with the potential health and ecological impact of nanotechnology that is being integrated in manufactured products. Greenpeace for instance commissioned a report concerning nanotechnology in which they express concern into the toxicity of nanomaterials that have been introduced in the environment. However, it makes only passing references to "assembler" technology. The UK Royal Society and UK Royal Academy of Engineering also commissioned a report entitled "Nanoscience and nanotechnologies: opportunities and uncertainties" regarding the larger social and ecological implications on nanotechnology. This report does not discuss the threat posed by potential so-called "molecular assemblers."

Formal scientific review

In 2006, U.S. National Academy of Sciences released the report of a study of molecular manufacturing as part of a longer report, *A Matter of Size: Triennial Review of the National Nanotechnology Initiative*. The study committee reviewed the technical content of *Nanosystems*, and in its conclusion states that no current theoretical analysis can be considered definitive regarding several questions of potential system performance, and that optimal paths for implementing high-performance systems cannot be predicted with confidence. It recommends experimental research to advance knowledge in this area:

"Although theoretical calculations can be made today, the eventually attainable range of chemical reaction cycles, error rates, speed of operation, and thermodynamic efficiencies of such bottom-up manufacturing systems cannot be reliably predicted at this time. Thus, the eventually attainable perfection and complexity of manufactured products, while they can be calculated in theory, cannot be predicted with confidence. Finally, the optimum research paths that might lead to systems which greatly exceed the thermodynamic efficiencies and other capabilities of biological systems cannot be reliably predicted at this time. Research funding that is based on the ability of investigators to produce experimental demonstrations that link to abstract models and guide long-term vision is most appropriate to achieve this goal."

Grey goo

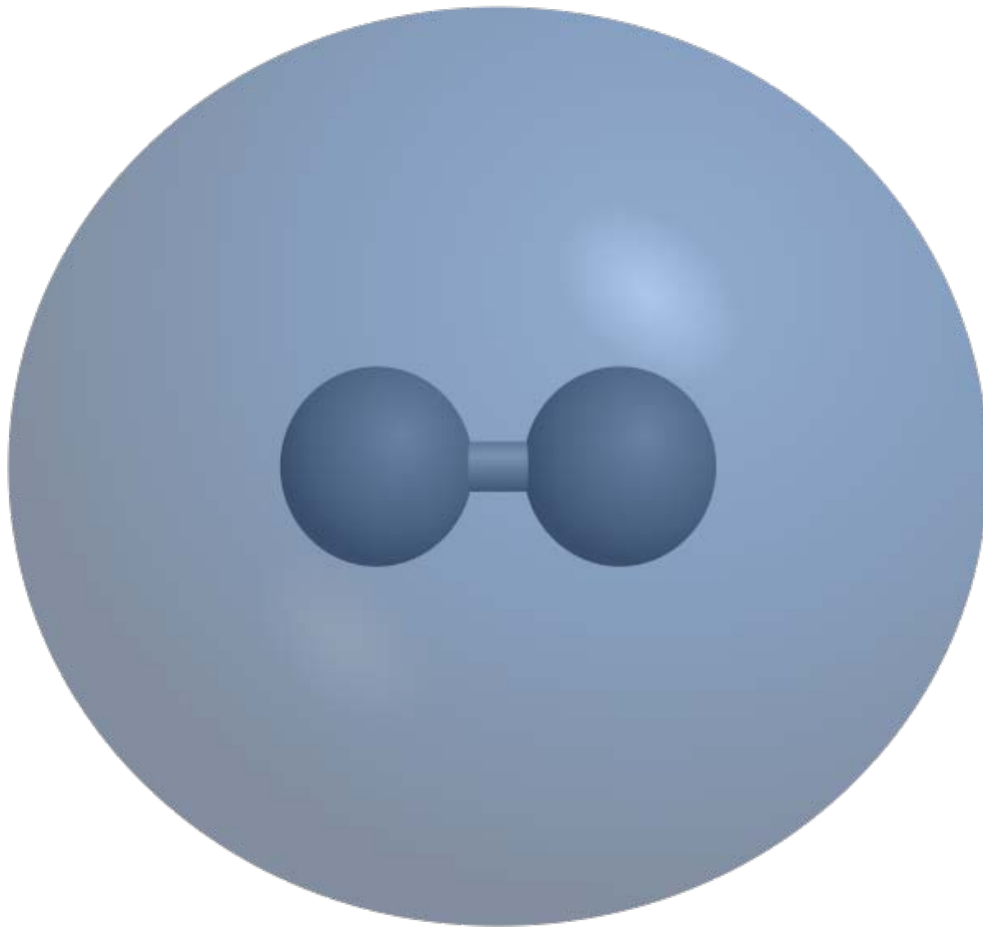
One potential scenario that has been envisioned is out-of-control self-replicating molecular assemblers in the form of grey goo which consumes carbon to continue its replication. If unchecked such mechanical replication could potentially consume whole ecoregions or the whole Earth (ecophagy), or it could simply outcompete natural lifeforms for necessary resources such as carbon, ATP, or UV light (which some nanomotor examples run on). It is worth noting that the ecophagy and 'grey goo' scenarios, like synthetic molecular assemblers, are based upon still-theoretical technologies that have not yet been demonstrated experimentally.

Chapter 6

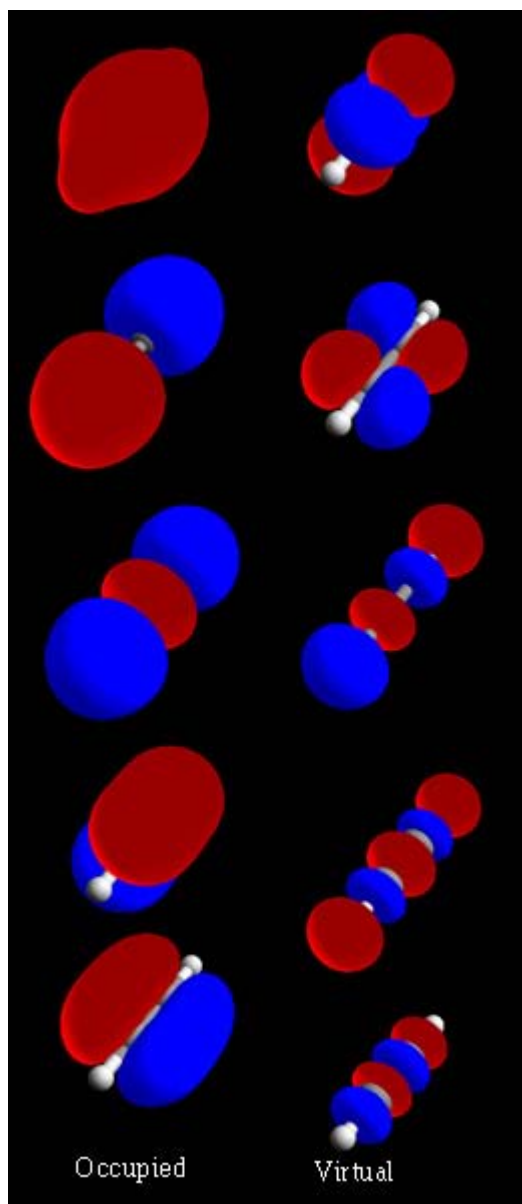
Molecular Orbital and Molecular Orbital Theory

Molecular orbital

In chemistry, a **molecular orbital** (or **MO**) is a mathematical function describing the wave-like behavior of an electron in a molecule. This function can be used to calculate chemical and physical properties such as the probability of finding an electron in any specific region. The term "orbital" was first used in English by Robert S. Mulliken as the English translation of Schrödinger's 'Eigenfunktion'. It has since been equated with the "region" generated with the function. Molecular orbitals are usually constructed by combining atomic orbitals or hybrid orbitals from each atom of the molecule, or other molecular orbitals from groups of atoms. They can be quantitatively calculated using the Hartree-Fock or Self-Consistent Field (SCF) methods.



H_2 1σ bonding molecular orbital



Complete acetylene ($\text{H-C}\equiv\text{C-H}$) molecular orbital set

Overview

A molecular orbital (MO) can specify the electron configuration of a molecule: the spatial distribution and energy of one (or one pair of) electron(s). Most commonly an MO is represented as a linear combination of atomic orbitals (the LCAO-MO method), especially in qualitative or very approximate usage. They are invaluable in providing a simple model of bonding in molecules, understood through molecular orbital theory.

Most present-day methods in computational chemistry begin by calculating the MOs of the system. A molecular orbital describes the behavior of one electron in the electric field generated by the nuclei and some average distribution of the other electrons. In the case

of two electrons occupying the same orbital, the Pauli principle demands that they have opposite spin. Necessarily this is an approximation, and highly accurate descriptions of the molecular electronic wave function do not have orbitals.

Qualitative discussion

For an imprecise, but qualitatively useful, discussion of the molecular structure, the molecular orbitals can be obtained from the "Linear combination of atomic orbitals molecular orbital method" ansatz. Here, the molecular orbitals are expressed as linear combinations of atomic orbitals.

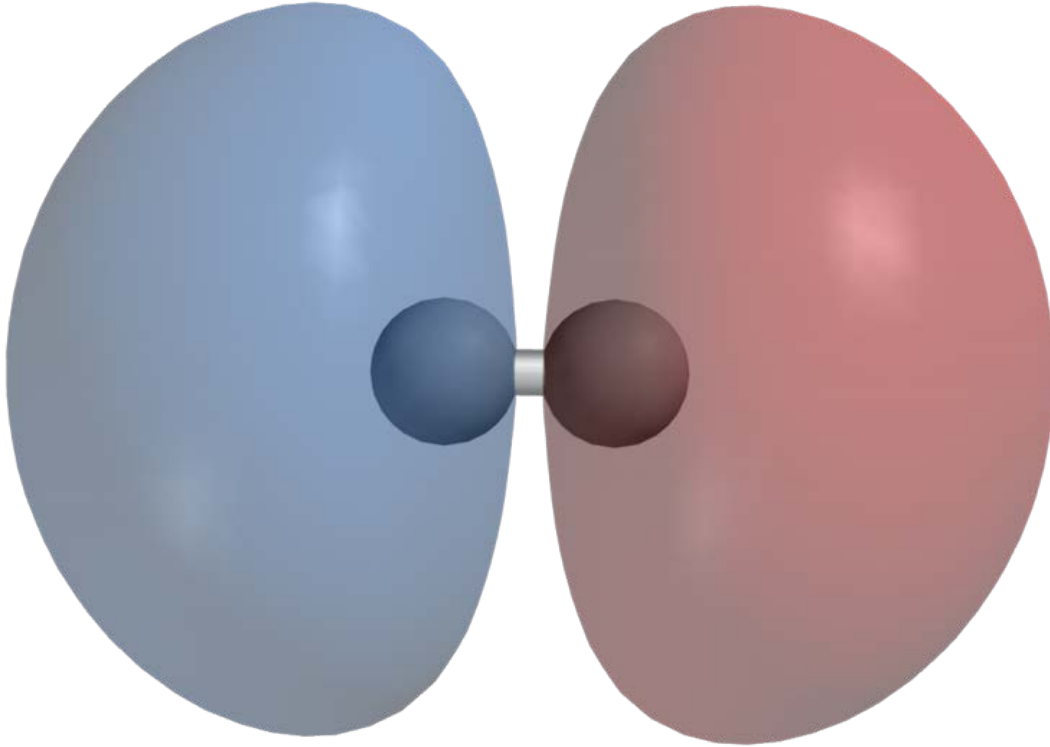
Molecular orbitals were first introduced by Friedrich Hund and Robert S. Mulliken in 1927 and 1928. The linear combination of atomic orbitals or "LCAO" approximation for molecular orbitals was introduced in 1929 by Sir John Lennard-Jones.. His ground-breaking paper showed how to derive the electronic structure of the fluorine and oxygen molecules from quantum principles. This qualitative approach to molecular orbital theory is part of the start of modern quantum chemistry.

Some properties:

- The number of molecular orbitals is equal to the number of atomic orbitals included in the linear expansion,
- If the molecule has some symmetry, the degenerate atomic orbitals (with the same atomic energy) are grouped in linear combinations (called **symmetry adapted atomic orbitals (SO)**) which belong to the representation of the symmetry group, so the wave functions that describe the group are known as **symmetry-adapted linear combinations (SALC)**.
- The number of molecular orbitals belonging to one group representation is equal to the number of symmetry-adapted atomic orbitals belonging to this representation,
- Within a particular representation, the symmetry-adapted atomic orbitals mix more if their atomic energy levels are closer.

Examples

H₂



H₂ 1σ* antibonding molecular orbital

As a simple example consider the hydrogen molecule, H₂, with the two atoms labelled H' and H''. The lowest-energy atomic orbitals, 1s' and 1s'', do not transform according to the symmetries of the molecule. However, the following symmetry adapted atomic orbitals do:

1s' - 1s'' Antisymmetric combination: negated by reflection, unchanged by other operations

1s' + 1s'' Symmetric combination: unchanged by all symmetry operations

The symmetric combination (called a bonding orbital) is lower in energy than the basis orbitals, and the antisymmetric combination (called an antibonding orbital) is higher. Because the H₂ molecule has two electrons, they can both go in the bonding orbital, making the system lower in energy (and hence more stable) than two free hydrogen atoms. This is called a covalent bond. The *bond order* is equal to the number of bonding electrons minus the number of antibonding electrons, divided by 2. In this example there

are 2 electrons in the bonding orbital and none in the antibonding orbital; the bond order is 1, and there is a single bond between the two hydrogen atoms.

He₂

On the other hand, consider the hypothetical molecule of He₂ with the atoms labelled He' and He''. Again, the lowest-energy atomic orbitals, 1s' and 1s'', do not transform according to the symmetries of the molecule, while the following symmetry adapted atomic orbitals do:

1s' - 1s'' Antisymmetric combination: negated by reflection, unchanged by other operations

1s' + 1s'' Symmetric combination: unchanged by all symmetry operations

Similar to the molecule H₂, the symmetric combination (called a bonding orbital) is lower in energy than the basis orbitals, and the antisymmetric combination (called an antibonding orbital) is higher. However, in its neutral ground state, each helium atom contains two electrons in its 1s orbital, combining for a total of four electrons. Two electrons fill the lower energy bonding orbital, while the remaining two fill the higher energy antibonding orbital. Thus, the resulting electron density around the molecule does not support the formation of a bond between the two atoms (called a sigma bond); therefore, the molecule does not exist. Another way of looking at it is that there are two bonding electrons and two antibonding electrons; therefore, the bond order is 0 and no bond exists.

Noble gases

Considering a hypothetical molecule of He₂, since the basis set of atomic orbitals is the same as in the case of H₂, we find that both the bonding and antibonding orbitals are filled, so there is no energy advantage to the pair. HeH would have a slight energy advantage, but not as much as H₂ + 2 He, so the molecule exists only a short while. In general, we find that atoms such as He that have completely full energy shells rarely bond with other atoms. Except for short-lived Van der Waals complexes, there are very few noble gas compounds known.

Ionic bonds

When the energy difference between the atomic orbitals of two atoms is quite large, one atom's orbitals contribute almost entirely to the bonding orbitals, and the other's almost entirely to the antibonding orbitals. Thus, the situation is effectively that some electrons have been transferred from one atom to the other. This is called an (mostly) ionic bond.

MO diagrams

For more complicated molecules, the wave mechanics approach loses utility in a qualitative understanding of bonding (although is still necessary for a quantitative

approach). The qualitative approach of MO uses a molecular orbital diagram. In this type of diagram, the molecular orbitals are represented by horizontal lines; the higher a line, the higher the energy of the orbital, and degenerate orbitals are placed on the same level with a space between them. Then, the electrons to be placed in the molecular orbitals are slotted in one by one, keeping in mind the Pauli exclusion principle and Hund's rule of maximum multiplicity (only 2 electrons, having opposite spins, per orbital; have as many unpaired electrons on one energy level as possible before starting to pair them).

HOMO and LUMO

The highest occupied molecular orbital and lowest unoccupied molecular orbital are often referred to as the HOMO and LUMO, respectively. The difference of the energies of the HOMO and LUMO, termed the band gap, can sometimes serve as a measure of the excitability of the molecule: the smaller the energy, the more easily it will be excited.

More quantitative approach

To obtain quantitative values for the molecular energy levels, one needs to have molecular orbitals which are such that the configuration interaction (CI) expansion converges fast towards the full CI limit. The most common method to obtain such functions is the Hartree–Fock method which expresses the molecular orbitals as eigenfunctions of the Fock operator. One usually solves this problem by expanding the molecular orbitals as linear combinations of gaussian functions centered on the atomic nuclei. The equation for the coefficients of these linear combinations is a generalized eigenvalue equation known as the Roothaan equations which are in fact a particular representation of the Hartree-Fock equation.

Simple accounts often suggest that experimental molecular orbital energies can be obtained by the methods of ultra-violet photoelectron spectroscopy for valence orbitals and X-ray photoelectron spectroscopy for core orbitals. This however is incorrect as these experiments measure the ionization energy, the difference in energy between the molecule and one of the ions resulting from the removal of one electron. Ionization energies are linked approximately to orbital energies by Koopmans' theorem. While the agreement between these two values can be close for some molecules, it can be very poor in other cases.

Molecular orbital theory

In chemistry, **molecular orbital (MO) theory** is a method for determining molecular structure in which electrons are not assigned to individual bonds between atoms, but are treated as moving under the influence of the nuclei in the whole molecule. In this theory, each molecule has a set of molecular orbitals, in which it is assumed that the molecular orbital wave function ψ_f may be written as a simple weighted sum of the n constituent atomic orbitals χ_i , according to the following equation:

$$\psi_j = \sum_{i=1}^n c_{ij} \chi_i$$

The c_{ij} coefficients may be determined numerically by substitution of this equation into the Schrödinger equation and application of the variational principle. This method is called the linear combination of atomic orbitals (LCAO) approximation and is used in computational chemistry. An additional unitary transformation can be applied on the system to accelerate the convergence in some computational schemes. Molecular orbital theory was seen as a competitor to valence bond theory in the 1930s, before it was realized that the two methods are closely related and that when extended they become equivalent.

History

Molecular orbital theory was developed, in the years after valence bond theory had been established (1927), primarily through the efforts of Friedrich Hund, Robert Mulliken, John C. Slater, and John Lennard-Jones. MO theory was originally called the Hund-Mulliken theory. The word *orbital* was introduced by Mulliken in 1932. By 1933, the molecular orbital theory had become accepted as a valid and useful theory. According to German physicist and physical chemist Erich Hückel, the first quantitative use of molecular orbital theory was the 1929 paper of Lennard-Jones. The first accurate calculation of a molecular orbital wavefunction was that made by Charles Coulson in 1938 on the hydrogen molecule. By 1950, molecular orbitals were completely defined as eigenfunctions (wave functions) of the self-consistent field Hamiltonian and it was at this point that molecular orbital theory became fully rigorous and consistent. This rigorous approach is known as the Hartree–Fock method for molecules although it had its origins in calculations on atoms. In calculations on molecules, the molecular orbitals are expanded in terms of an atomic orbital basis set, leading to the Roothaan equations. This led to the development of many ab initio quantum chemistry methods. In parallel, molecular orbital theory was applied in a more approximate manner using some empirically derived parameters in methods now known as semi-empirical quantum chemistry methods.

Overview

Molecular orbital (MO) theory uses a linear combination of atomic orbitals (LCAO) to represent molecular orbitals involving the whole molecule. These are often divided into bonding orbitals, anti-bonding orbitals, and non-bonding orbitals. A molecular orbital is merely a Schrödinger orbital which includes several, but often only two nuclei. If this orbital is of type in which the electron(s) in the orbital have a higher probability of being *between* nuclei than elsewhere, the orbital will be a bonding orbital, and will tend to hold the nuclei together. If the electrons tend to be present in a molecular orbital in which they spend more time elsewhere than between the nuclei, the orbital will function as an anti-bonding orbital and will actually weaken the bond. Electrons in non-bonding orbitals tend to be in deep orbitals (nearly atomic orbitals) associated almost entirely with one nucleus

or the other, and thus they spend equal time between nuclei or not. These electrons neither contribute nor detract from bond strength.

Molecular orbitals are further divided according to the types of atomic orbitals combining to form a bond. These orbitals are results of electron-nucleus interactions that are caused by the fundamental force of electromagnetism. Chemical substances will form a bond if their orbitals become lower in energy when they interact with each other. Different chemical bonds are distinguished that differ by electron cloud shape and by energy levels.

MO theory provides a global, delocalized perspective on chemical bonding. For example, in the MO theory for hypervalent molecules it is unnecessary to invoke a major role for d-orbitals, whereas valence bond theory normally uses hybridization with d-orbitals to explain hypervalency. In MO theory, *any* electron in a molecule may be found *anywhere* in the molecule, since quantum conditions allow electrons to travel under the influence of an arbitrarily large number of nuclei, so long as permitted by certain quantum rules. Although in MO theory *some* molecular orbitals may hold electrons which are more localized between specific pairs of molecular atoms, *other* orbitals may hold electrons which are spread more uniformly over the molecule. Thus, overall, bonding (and electrons) are far more delocalized (spread out) in MO theory, than is implied in valence bond (VB) theory. This makes MO theory more useful for the description of extended systems.

An example is that in the MO picture of benzene, composed of a hexagonal ring of 6 carbon atoms. In this molecule, 24 of the 30 total valence bonding electrons are located in 12 σ (sigma) bonding orbitals which are mostly located between pairs of atoms (C-C or C-H), similar to the valence bond picture. However, in benzene the remaining 6 bonding electrons are located in 3 π (pi) molecular bonding orbitals that are delocalized around the ring. Two are in an MO which has equal contributions from all 6 atoms. The other two orbitals have vertical nodes at right angles to each other. As in the VB theory, all of these 6 delocalized pi electrons reside in a larger space which exists above and below the ring plane. All carbon-carbon bonds in benzene are chemically equivalent. In MO theory this is a direct consequence of the fact that the 3 molecular pi orbitals form a combination which evenly spreads the extra 6 electrons over 6 carbon atoms.

In molecules such as methane, the 8 valence electrons are found in 4 MOs that are spread out over all 5 atoms. However, it is possible to approximate the MOs with 4 localized orbitals similar in shape to sp^3 hybrid orbitals predicted by VB theory. This is often adequate for σ (sigma) bonds, but it is not possible for the π (pi) orbitals. However, the delocalized MO picture is more appropriate for ionization and spectroscopic predictions. Upon ionization of methane, a single electron is taken from the MO which surrounds the whole molecule, weakening all 4 bonds equally. VB theory would predict that one electron is removed for an sp^3 orbital, resulting in the need for resonance between four valence bond structures, each of which has a one-electron bond.

As in benzene, in substances such as beta carotene, chlorophyll or heme, some electrons the π (pi) orbitals are spread out in molecular orbitals over long distances in a molecule,

giving rise to light absorption in lower energies (visible colors), a fact which is observed. This and other spectroscopic data for molecules are better explained in MO theory, with an emphasis on electronic states associated with multicenter orbitals, including mixing of orbitals premised on principles of orbital symmetry matching. The same MO principles also more naturally explain some electrical phenomena, such as high electrical conductivity in the planar direction of the hexagonal atomic sheets that exist in graphite. In MO theory, "resonance" (a mixing and blending of VB bond states) is a natural consequence of symmetry. For example, in graphite, as in benzene, it is not necessary to invoke the sp^2 hybridization and resonance of VB theory, in order to explain electrical conduction. Instead, MO theory simply recognizes that some electrons in the graphite atomic sheets are completely delocalized over arbitrary distances, and reside in very large *molecular orbitals* that cover an entire graphite sheet, and some electrons are thus as free to move and conduct electricity *in the sheet plane*, as if they resided in a metal.

Chapter 7

Molecular Model

A **molecular model** is a physical model that represents molecules and their processes. The creation of mathematical models of molecular properties and behaviour is **molecular modelling**, and their graphical depiction is **molecular graphics**, but these topics are closely linked and each uses techniques from the others. In this, "molecular model" will primarily refer to systems containing more than one atom and where nuclear structure is neglected. The electronic structure is often also omitted or represented in a highly simplified way.

Overview

Physical models of atomistic systems have played an important role in understanding chemistry and generating and testing hypotheses. Most commonly there is an explicit representation of atoms, though other approaches such as soap films and other continuous media have been useful. There are several motivations for creating physical models:

- as pedagogic tools for students or those unfamiliar with atomistic structures;
- as objects to generate or test theories (e.g., the structure of DNA);
- as analogue computers (e.g., for measuring distances and angles in flexible systems);
- as aesthetically pleasing objects on the boundary of art and science.

The construction of physical models is often a creative act, and many bespoke examples have been carefully created in the workshops of science departments. There is a very wide range of approaches to physical modelling, and this lists only the most common or historically important. The main strategies are:

- bespoke construction of a single model;
- use of common materials (plasticine, matchsticks) or children's toys (Tinkertoy(TM), Meccano, Lego, etc.);
- re-use of generic components in kits (ca. 1930s to present).

Models encompass a wide range of degrees of precision and engineering: some models such as J.D. Bernal's water are conceptual, while the macromodels of Pauling and Crick and Watson were created with much greater precision.

Molecular models have inspired molecular graphics, initially in textbooks and research articles and more recently on computers. Molecular graphics has replaced some functions of physical molecular models, but physical kits continue to be very popular and are sold in large numbers. Their unique strengths include:

- cheapness and portability;
- immediate tactile and visual messages;
- easy interactivity for many processes (e.g., conformational analysis and pseudorotation).

History

In the 1600s, Johannes Kepler speculated on the symmetry of snowflakes and also on the close packing of spherical objects such as fruit (this problem remained unsolved until very recently). The symmetrical arrangement of closely packed spheres informed theories of molecular structure in the late 1800s, and many theories of crystallography and solid state inorganic structure used collections of equal and unequal spheres to simulate packing and predict structure.

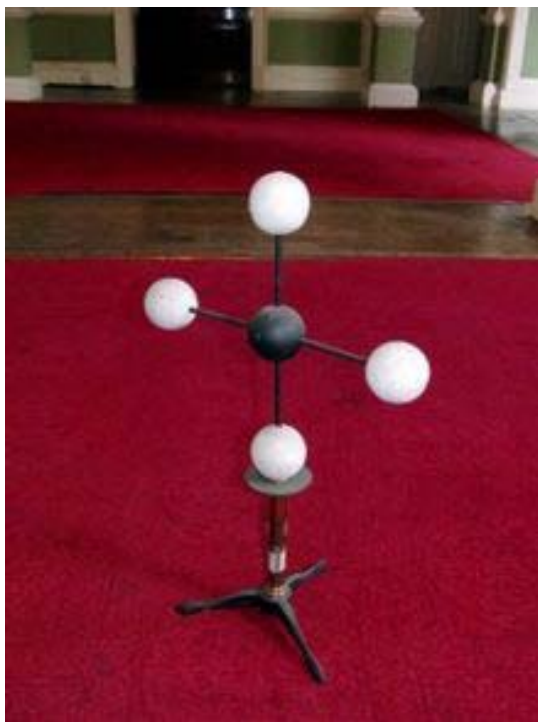


Fig. 1. Hofmann's model for methane.

John Dalton represented compounds as aggregations of circular atoms, and although Loschmidt did not create physical models, his diagrams based on circles are two-dimensional analogues of later models. Hofmann is credited with the first physical molecular model around 1860 (Fig. 1). Note how the size of the carbon appears smaller than the hydrogen. The importance of stereochemistry was not then recognised and the model is essentially topological (it should be a 3-dimensional tetrahedron).

J.H. van 't Hoff and J. le Bel introduced the concept of chemistry in space— stereochemistry in three dimensions. van 't Hoff built tetrahedral molecules representing the three-dimensional properties of carbon.

Models based on spheres

Robert Hooke proposed a relationship between crystals and the packing of spheres . R. Häüy argued that the structures of crystals involved regular lattices of repeating units with shapes similar to the macroscopic crystal. Barlow, who jointly developed the theories of space groups, proposed models of crystals based on sphere packings (ca. 1890).



Fig. 2. Sodium chloride (NaCl) lattice, showing close-packed spheres representing a face-centered cubic AB lattice similar to that of NaCl and most other alkali halides. In this

model the spheres are equal sizes whereas more "realistic" models would have different radii for cations and anions.

The binary compounds sodium chloride (NaCl) and caesium chloride (CsCl) have cubic structures but have different space groups. This can be rationalised in terms of close packing of spheres of different sizes. For example, NaCl can be described as close-packed chloride ions (in a face-centered cubic lattice) with sodium ions in the octahedral holes. After the development of X-ray crystallography as a tool for determining crystal structures, many laboratories built models based on spheres. With the development of plastic or polystyrene balls it is now easy to create such models.

Models based on ball-and-stick

The concept of the chemical bond as a direct link between atoms can be modelled by linking balls (atoms) with sticks/rods (bonds). This has been extremely popular and is still widely used today. Initially atoms were made of spherical wooden balls with specially drilled holes for rods. Thus carbon can be represented as a sphere with four holes at the tetrahedral angles $\cos^{-1}(-1/3) \approx 109.47^\circ$.

A problem with rigid bonds and holes is that systems with arbitrary angles could not be built. This can be overcome with flexible bonds, originally helical springs but now usually plastic. This also allows double and triple bonds to be approximated by multiple single bonds (Fig. 3).

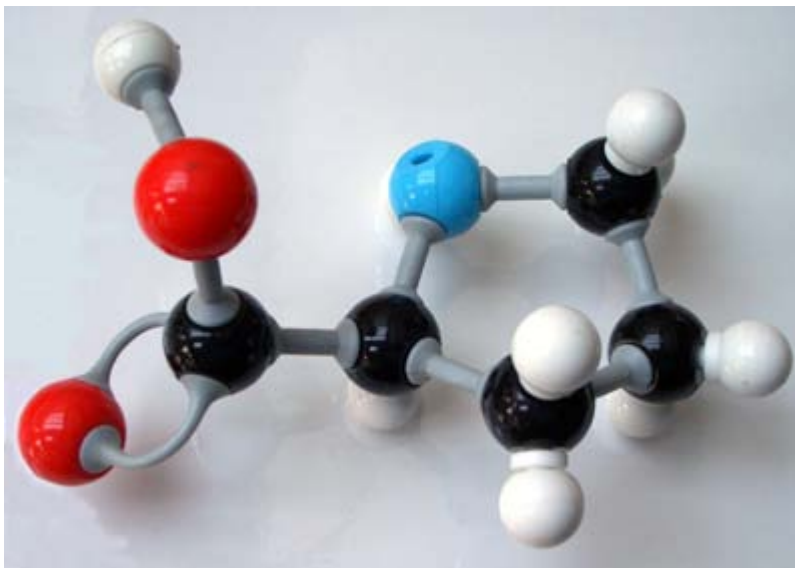


Fig 3. A modern plastic ball and stick model. The molecule shown is proline.

Figure 3 represents a ball-and-stick model of proline. The balls have colours: **black** represents carbon (C); **red**, oxygen (O); **blue**, nitrogen (N); and white, hydrogen (H). Each ball is drilled with as many holes as its conventional valence (C: 4; N: 3; O: 2; H: 1) directed towards the vertices of a tetrahedron. Single bonds are represented by (fairly)

rigid grey rods. Double and triple bonds use two longer flexible bonds which restrict rotation and support conventional cis/trans stereochemistry.

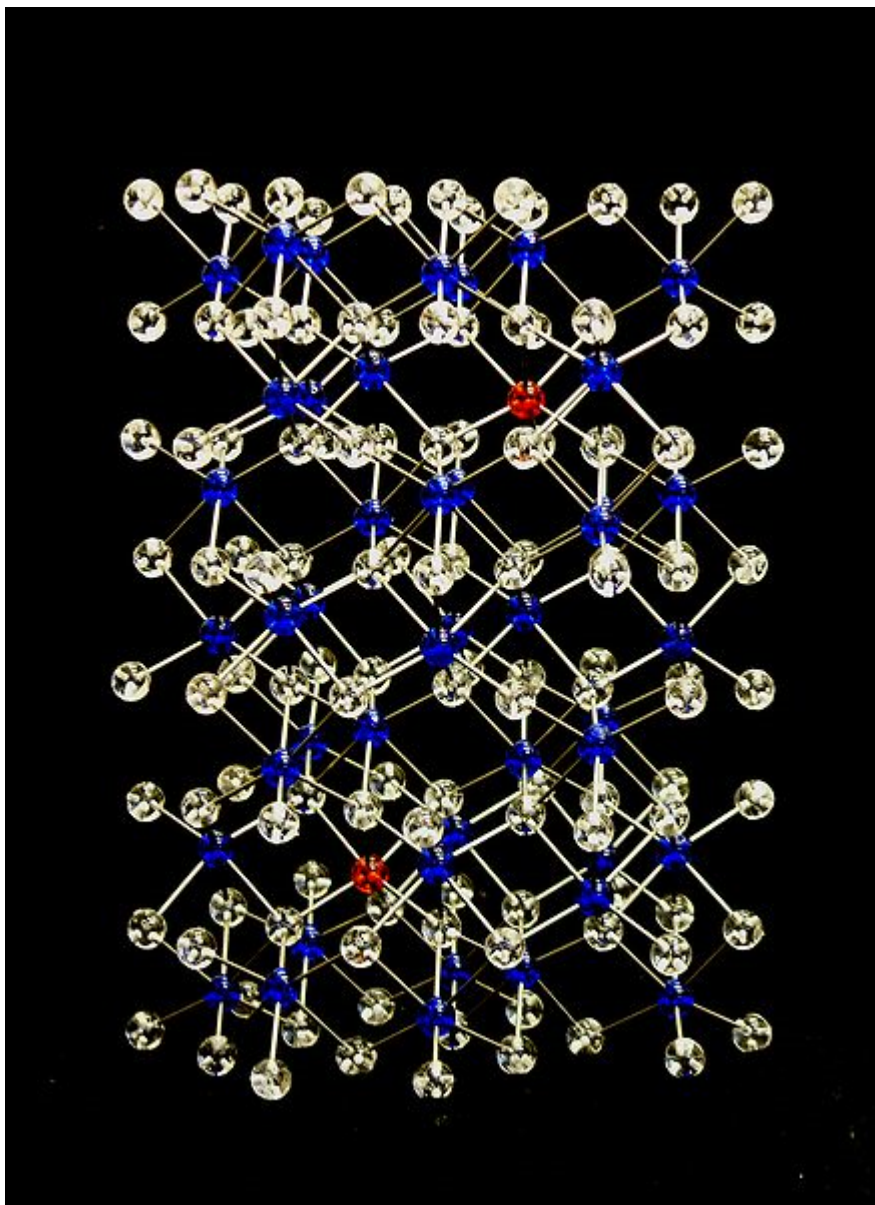


Fig. 4. Beevers ball and stick model of ruby (Cr-doped corundum) made with acrylic balls and stainless steel rods.

However, most molecules require holes at other angles and specialist companies manufacture kits and bespoke models. One of the earlier companies was Woosters at Bottisham, Cambridgeshire, UK. Besides tetrahedral, trigonal and octahedral holes, there were all-purpose balls with 24 holes. These models allowed rotation about the single rod bonds, which could be both an advantage (showing molecular flexibility) and a disadvantage (models are floppy). The approximate scale was 5 cm per ångström (0.5 m/nm or 500,000,000:1), but was not consistent over all elements.

Arnold Beevers in Edinburgh (now operating as Miramodus) created small models using PMMA balls and stainless steel rods. By using individually drilled balls with precise bond angles and bond lengths in these models, large crystal structures to be accurately created, but with light and rigid form. Figure 4 shows a unit cell of ruby in this style.

Skeletal models

Crick and Watson's DNA model and the protein-building kits of Kendrew were among the first skeletal models. These were based on atomic components where the valences were represented by rods; the atoms were points at the intersections. Bonds were created by linking components with tubular connectors with locking screws.

Andre Dreiding introduced a molecular modelling kit (ca. 1975) which dispensed with the connectors. A given atom would have solid and hollow valence spikes. The solid rods clicked into the tubes forming a bond, usually with free rotation. These were and are very widely used in organic chemistry departments and were made so accurately that interatomic measurements could be made by ruler.

More recently, inexpensive plastic models (such as Orbit) use a similar principle. A small plastic sphere has protuberances onto which plastic tubes can be fitted. The flexibility of the plastic means that distorted geometries can be made.

Polyhedral models

Many inorganic solids consist of atoms surrounded by a coordination sphere of electronegative atoms (e.g. PO_4 tetrahedra, TiO_6 octahedra). Structures can be modelled by gluing together polyhedra made of paper or plastic.

Composite models

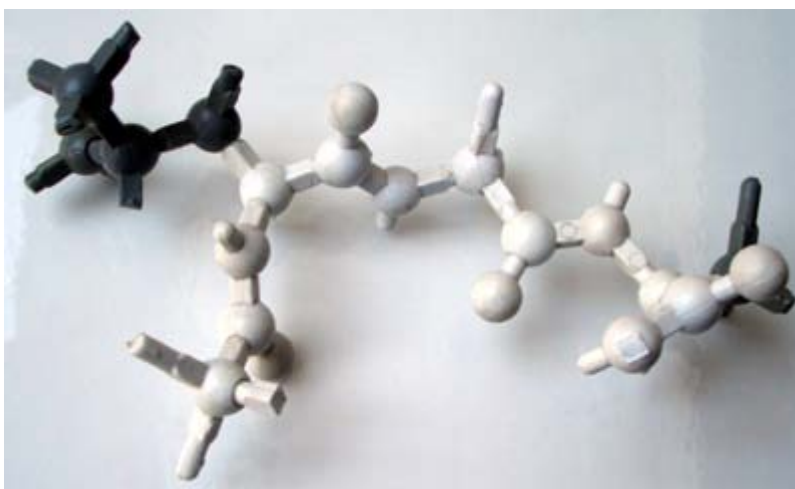


Fig. 5. A Nicholson model, showing a short part of protein backbone (white) with side chains (grey). Note the snapped stubs representing hydrogen atoms.

A good example of composite models is the Nicholson approach, widely used from the late 1970s for building models of biological macromolecules. The components are primarily amino acids and nucleic acids with preformed residues representing groups of atoms. Many of these atoms are directly moulded into the template, and fit together by pushing plastic stubs into small holes. The plastic grips well and makes bonds difficult to rotate, so that arbitrary torsion angles can be set and retain their value. The conformations of the backbone and side chains are determined by pre-computing the torsion angles and then adjusting the model with a protractor.

The plastic is white and can be painted to distinguish between O and N atoms. Hydrogen atoms are normally implicit and modelled by snipping off the spokes. A model of a typical protein with approximately 300 residues could take a month to build. It was common for laboratories to build a model for each protein solved. By 2005, so many protein structures were being determined that relatively few models were made.

Computer-based models

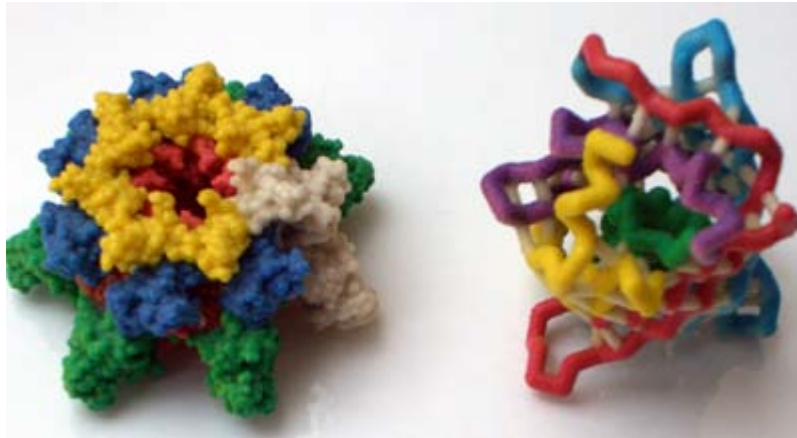


Fig. 6. Integrated protein models.

With the development of computer-based physical modelling, it is now possible to create complete single-piece models by feeding the coordinates of a surface into the computer. Figure 6 shows models of anthrax toxin, left (at a scale of approximately 20 Å/cm or 1:5,000,000) and green fluorescent protein, right (5 cm high, at a scale of about 4 Å/cm or 1:25,000,000) from 3D Molecular Design. Models are made of plaster or starch, using a rapid prototyping process.

It has also recently become possible to create accurate molecular models inside glass blocks using a technique known as subsurface laser engraving. The image at right (Fig. 7) shows the 3D structure of an *E. coli* protein (DNA polymerase beta-subunit, PDB code 1MMI) etched inside a block of glass by British company Luminorum Ltd.

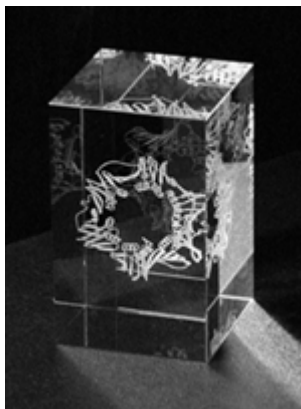


Fig. 7. Protein model in glass.

Common colors

Some of the most common colors used in molecular models are as follows:

Hydrogen	White
Alkali Metals	Violet
Alkaline-Earth Metals	Dark Green
Boron, Most Transition Metals	Peach/Salmon
Carbon	Black
Nitrogen	Dark Blue
Oxygen	Red
Fluorine, Chlorine	Green
Bromine	Dark Red
Iodine	Dark Violet
Noble Gases	Cyan
Phosphorus	Orange
Sulfur	Yellow
Titanium	Gray

Chronology

This table is an incomplete chronology of events where physical molecular models provided major scientific insights.

developer(s)	approximate date	technology	comments
Kepler			sphere packing, symmetry of snowflakes.
Loschmidt		2-D graphics	representation of atoms and bonds by touching circles
Hofmann		ball-and-stick	first recognisable physical molecular model
van't Hoff		paper?	representation of atoms as tetrahedra supported the development of stereochemistry
Bernal		Plasticine and spokes	model of liquid water
Corey, Pauling, Koltun (CPK coloring)		Space filling models of alpha-helix, etc.	Pauling's "Nature of the Chemical Bond" covered all aspects of molecular structure and influenced many aspects of models
Crick and Watson		spikes, flat templates and connectors with screws	model of DNA
Molecular graphics	ca 1960	display on computer screens	complements rather than replaces physical models

Chapter 8

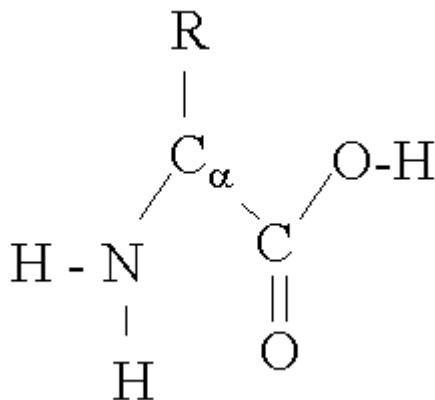
Protein Structure

Proteins are an important class of biological macromolecules present in all organisms. All proteins are polymers of amino acids. Classified by their physical size, proteins are nanoparticles (definition: 1–100 nm). Each protein polymer – also known as a polypeptide – consists of a sequence of 20 different L- α -amino acids, also referred to as residues. For chains under 40 residues the term peptide is frequently used instead of protein. To be able to perform their biological function, proteins fold into one or more specific spatial conformations, driven by a number of non-covalent interactions such as hydrogen bonding, ionic interactions, Van Der Waals forces, and hydrophobic packing. To understand the functions of proteins at a molecular level, it is often necessary to determine their three-dimensional structure. This is the topic of the scientific field of structural biology, which employs techniques such as X-ray crystallography, NMR spectroscopy, and dual polarisation interferometry to determine the structure of proteins.

Protein structures range in size from tens to several thousand residues. Very large aggregates can be formed from protein subunits: for example, many thousand actin molecules assemble into a microfilament.

A protein may undergo reversible structural changes in performing its biological function. The alternative structures of the same protein are referred to as different conformations, and transitions between them are called conformational changes.

Protein covalent structure and stereochemistry

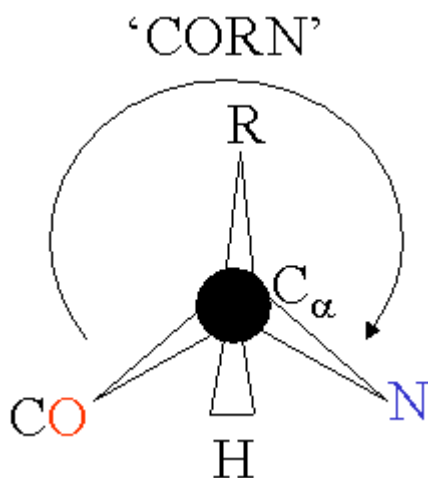


An α -amino acid. The C_αH atom is omitted in the diagram.

Protein amino acids are combined into a single polypeptide chain in a condensation reaction. This reaction is catalysed by the ribosome in a process known as translation.

Amino acid residues

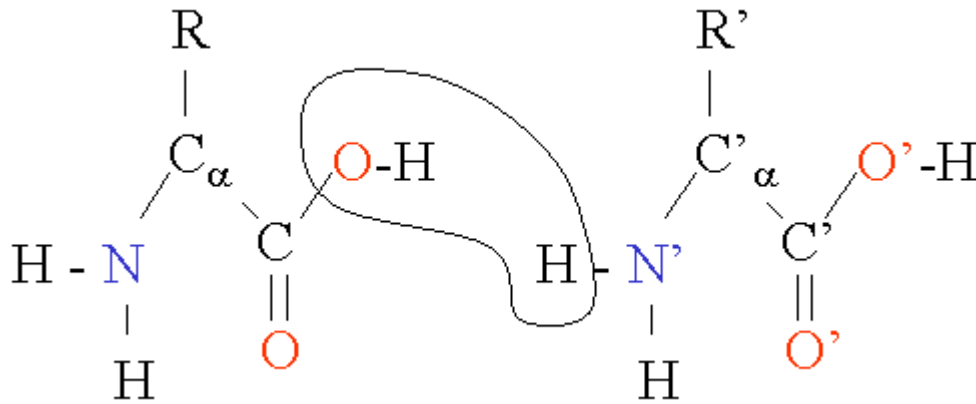
Each α -amino acid consists of a backbone part that is present in all the amino acid types, and a side chain that is unique to each type of residue. An exception from this rule is proline, where the hydrogen atom is replaced by a bond to the side chain. Because the carbon atom is bound to four different groups it is chiral, however only one of the isomers occur in biological proteins. Glycine however, is not chiral since its side chain is a hydrogen atom. A simple mnemonic for correct L-form is "CORN": when the C_α atom is viewed with the H in front, the residues read "CO-R-N" in a clockwise direction.



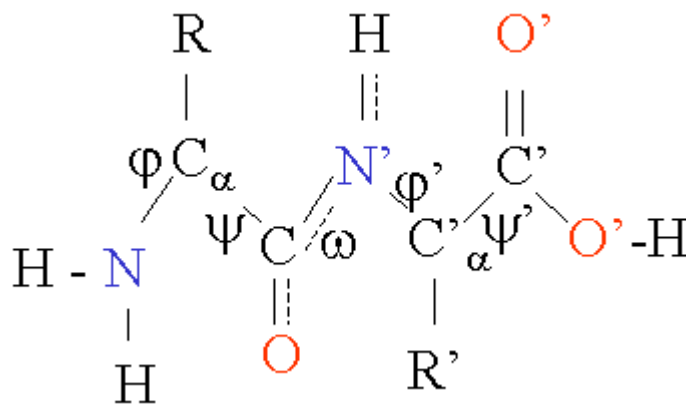
CO-R-N rule

The 20 naturally occurring amino acids have different physical and chemical properties, including their electrostatic charge, pKa, hydrophobicity, size and specific functional groups. These properties play a major role in molding protein structure.

The peptide bond



Two amino acids



Bond angles for ψ and ω

The peptide bond tends to be planar due to the delocalization of the electrons from the double bond. The rigid peptide dihedral angle, ω (the bond between C_1 and N) is always close to 180 degrees. The dihedral angles phi ϕ (the bond between N and C_α) and psi ψ (the bond between C_α and C_1) can have a certain range of possible values. These angles are the internal degrees of freedom of a protein, they control the protein's conformation. They are restrained by geometry to allowed ranges typical for particular secondary structure elements, and represented in a Ramachandran plot. A few important bond lengths are given in the table below.

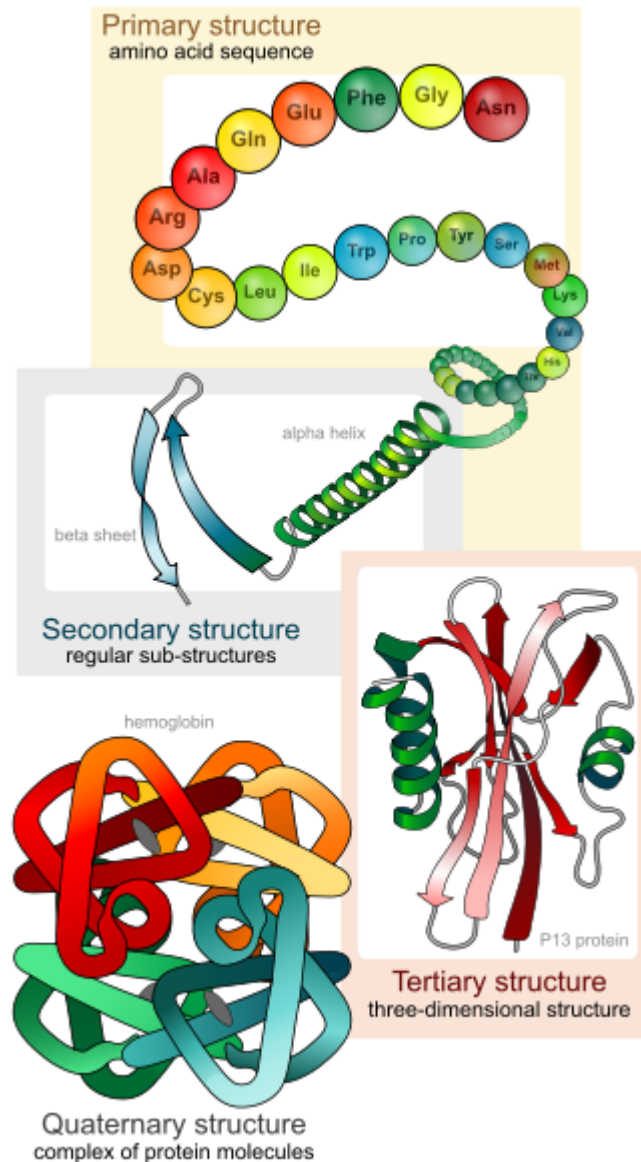
Peptide bond	Average length	Single bond	Average length	Hydrogen bond	Average (± 30)
C α - C	153 pm	C - C	154 pm	O-H --- O-H	280 pm
C - N	133 pm	C - N	148 pm	N-H --- O=C	290 pm
N - C α	146 pm	C - O	143 pm	O-H --- O=C	280 pm

Side-chain conformation

The atoms along the side chain are named with Greek letters in Greek alphabetical order: α , β , γ , δ , ϵ , and so on. C α refers to the carbon atom of the backbone closest to the carbonyl group of that amino acid, C β the second closest and so on. The dihedral angles around the bonds between these atoms are named χ_1 , χ_2 , χ_3 , etc. The dihedral angle of the first movable atom of the side chain, γ , defined as N-C α -C β -X γ , is named χ_1 . Side chains tend to adopt different staggered conformations called *gauche(-)*, *trans*, and *gauche(+)*, which corresponds to rotation angles of 60° , 180° , and -60° , respectively, around the sp³-sp³ bonds.

The diversity of side-chain conformations is often expressed in rotamer libraries. A rotamer library is a collection of rotamers for each residue type. Side-chain dihedral angles are not evenly distributed, but for most side chain types, the χ angles occur in tight clusters around certain values. Rotamer libraries therefore are usually derived from statistical analysis of side-chain conformations in known structures of proteins by clustering observed conformations or by dividing dihedral angle space into bins, and determining an average conformation in each bin.

Levels of protein structure



Protein structure, from primary to quaternary structure.

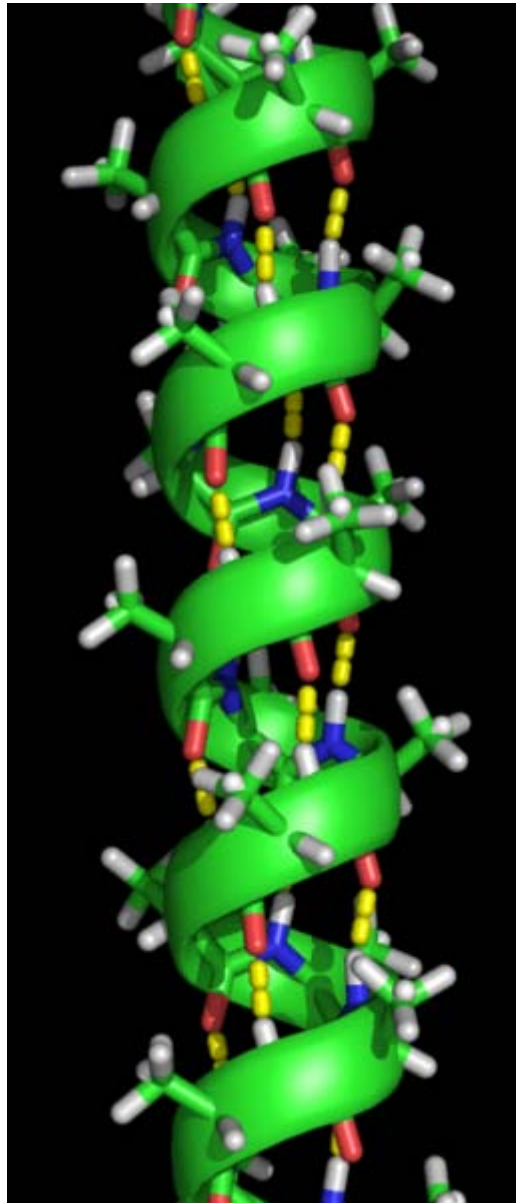
There are four distinct levels of protein structure.

Primary structure

The primary structure refers to the sequence of the different amino acids of the peptide or protein. The primary structure is held together by covalent or peptide bonds, which are made during the process of protein biosynthesis or translation. The two ends of the polypeptide chain are referred to as the carboxyl terminus (C-terminus) and the amino terminus (N-terminus) based on the nature of the free group on each extremity. Counting

of residues always starts at the N-terminal end (NH₂-group), which is the end where the amino group is not involved in a peptide bond. The primary structure of a protein is determined by the gene corresponding to the protein. A specific sequence of nucleotides in DNA is transcribed into mRNA, which is read by the ribosome in a process called translation. The sequence of a protein is unique to that protein, and defines the structure and function of the protein. The sequence of a protein can be determined by methods such as Edman degradation or tandem mass spectrometry. Often however, it is read directly from the sequence of the gene using the genetic code. Post-translational modifications such as disulfide formation, phosphorylations and glycosylations are usually also considered a part of the primary structure, and cannot be read from the gene.

Secondary structure



An alpha-helix with hydrogen bonds (yellow dots)

Secondary structure refers to highly regular local sub-structures. Two main types of secondary structure, the alpha helix and the beta strand, were suggested in 1951 by Linus Pauling and coworkers. These secondary structures are defined by patterns of hydrogen bonds between the main-chain peptide groups. They have a regular geometry, being constrained to specific values of the dihedral angles ψ and ϕ on the Ramachandran plot. Both the alpha helix and the beta-sheet represent a way of saturating all the hydrogen bond donors and acceptors in the peptide backbone. Some parts of the protein are ordered but do not form any regular structures. They should not be confused with random coil, an unfolded polypeptide chain lacking any fixed three-dimensional structure. Several sequential secondary structures may form a "supersecondary unit".

Tertiary structure

Tertiary structure refers to three-dimensional structure of a single protein molecule. The alpha-helices and beta-sheets are folded into a compact globule. The folding is driven by the *non-specific* hydrophobic interactions (the burial of hydrophobic residues from water), but the structure is stable only when the parts of a protein domain are locked into place by *specific* tertiary interactions, such as salt bridges, hydrogen bonds, and the tight packing of side chains and disulfide bonds. The disulfide bonds are extremely rare in cytosolic proteins, since the cytosol is generally a reducing environment.

Quaternary structure

Quaternary structure is a larger assembly of several protein molecules or polypeptide chains, usually called subunits in this context. The quaternary structure is stabilized by the same non-covalent interactions and disulfide bonds as the tertiary structure. Complexes of two or more polypeptides (i.e. multiple subunits) are called multimers. Specifically it would be called a dimer if it contains two subunits, a trimer if it contains three subunits, and a tetramer if it contains four subunits. The subunits are frequently related to one another by symmetry operations, such as a 2-fold axis in a dimer. Multimers made up of identical subunits are referred to with a prefix of "homo-" (e.g. a homotetramer) and those made up of different subunits are referred to with a prefix of "hetero-" (e.g. a heterotetramer, such as the two alpha and two beta chains of hemoglobin). Many proteins do not have the quaternary structure and function as monomers.

Domains, motifs, and folds in protein structure

Proteins are frequently described as consisting from several structural units.

- A **structural domain** is an element of the protein's overall structure that is self-stabilizing and often folds independently of the rest of the protein chain. Many domains are not unique to the protein products of one gene or one gene family but instead appear in a variety of proteins. Domains often are named and singled out because they figure prominently in the biological function of the protein they belong to; for example, the "calcium-binding domain of calmodulin". Because

they are independently stable, domains can be "swapped" by genetic engineering between one protein and another to make chimeras.

- The **structural and sequence motifs** refer to short segments of protein three-dimensional structure or amino acid sequence that were found in a large number of different proteins.
- The **supersecondary structure** refers to a specific combination of secondary structure elements, such as beta-alpha-beta units or helix-turn-helix motif. Some of them may be also referred to as structural motifs.
- **Protein fold** refers to the general protein architecture, like helix bundle, beta-barrel, Rossman fold or different "folds" provided in the Structural Classification of Proteins database.

Despite the fact that there are about 100,000 different proteins expressed in eukaryotic systems, there are many fewer different domains, structural motifs and folds. This is partly a consequence of evolution, since genes or parts of genes can be doubled or moved around within the genome. This means that, for example, a protein domain might be moved from one protein to another thus giving the protein a new function. Because of these mechanisms, pathways and mechanisms tend to be reused in several different proteins.

Protein folding

An unfolded polypeptide folds into its characteristic three-dimensional structure from random coil.

Protein structure determination

Around 90% of the protein structures available in the Protein Data Bank have been determined by X-ray crystallography. This method allows one to measure the 3D density distribution of electrons in the protein (in the crystallized state) and thereby infer the 3D coordinates of all the atoms to be determined to a certain resolution. Roughly 9% of the known protein structures have been obtained by Nuclear Magnetic Resonance techniques. The secondary structure composition can be determined via circular dichroism or dual polarisation interferometry. Cryo-electron microscopy has recently become a means of determining protein structures to high resolution (less than 5 angstroms or 0.5 nanometer) and is anticipated to increase in power as a tool for high resolution work in the next decade. This technique is still a valuable resource for researchers working with very large protein complexes such as virus coat proteins and amyloid fibers.

Structure classification

Protein structures can be classified based on their similarity or a common evolutionary origin. SCOP and CATH databases provide two different structural classifications of proteins.

Computational prediction of protein structure

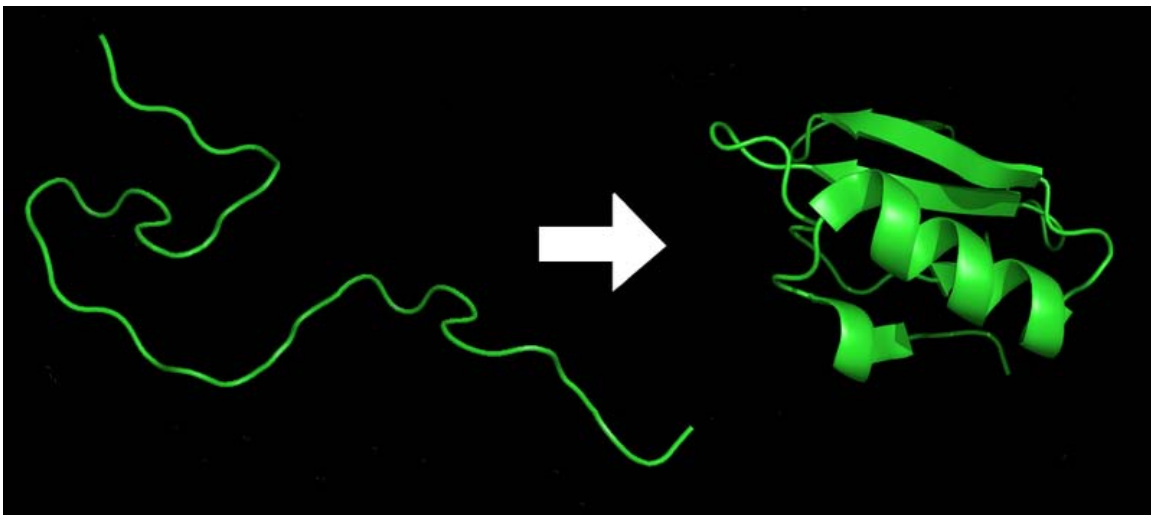
The generation of a protein sequence is much easier than the determination of a protein structure. However, the structure of a protein gives much more insight in the function of the protein than its sequence. Therefore, a number of methods for the computational prediction of protein structure from its sequence have been developed. *Ab initio* prediction methods use just the sequence of the protein. Threading and Homology Modeling methods can build a 3D model for a protein of unknown structure from experimental structures of evolutionary related proteins.

Protein structure related software

There are software to aid researchers working on, often overlapping, different aspects of protein structure. The most basic functionality is providing structure visualization. Analysis of protein structure can be facilitated by software that aligns structures. In the absence of existing structures for a given protein sequence, there are methods to predict or to model the structure of such sequences based on known protein structures. And given models of known or predicted structures, one can use software to verify them for errors, predict protein conformational changes, or predict substrate binding sites.

Chapter 9

Protein Folding



Protein before and after folding.

Protein folding is the physical process by which a polypeptide folds into its characteristic and functional three-dimensional structure from random coil. Each protein exists as an unfolded polypeptide or random coil when translated from a sequence of mRNA to a linear chain of amino acids. This polypeptide lacks any developed three-dimensional structure (the left hand side of the neighboring figure). Amino acids interact with each other to produce a well-defined three-dimensional structure, the folded protein (the right hand side of the figure), known as the native state. The resulting three-dimensional structure is determined by the amino acid sequence (Anfinsen's dogma).

The correct three-dimensional structure is essential to function, although some parts of functional proteins may remain unfolded. Failure to fold into the intended shape usually produces inactive proteins with different properties including toxic prions. Several neurodegenerative and other diseases are believed to result from the accumulation of *misfolded* (incorrectly folded) proteins. Many allergies are caused by the folding of the proteins, for the immune system does not produce antibodies for certain protein structures.

Known facts

Relationship between folding and amino acid sequence

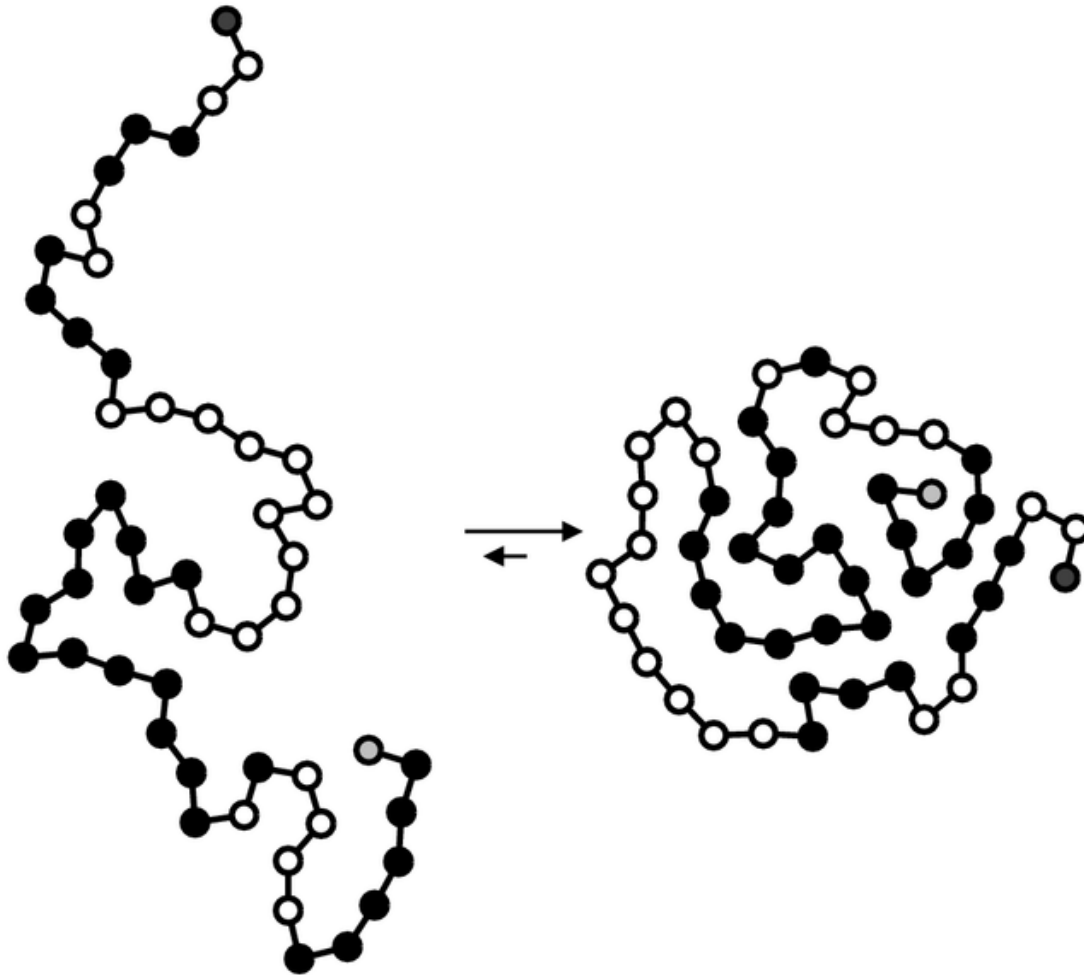


Illustration of the main driving force behind protein structure formation. In the compact fold (to the right), the hydrophobic amino acids (shown as black spheres) are in general shielded from the solvent.

The amino-acid sequence (or primary structure) of a protein determines its native conformation. A protein molecule folds spontaneously during or after biosynthesis. While these macromolecules may be regarded as "folding themselves", the process also depends on the solvent (water or lipid bilayer), the concentration of salts, the temperature, and the presence of molecular chaperones.

Folded proteins usually have a hydrophobic core in which side chain packing stabilizes the folded state, and charged or polar side chains occupy the solvent-exposed surface where they interact with surrounding water. Minimizing the number of hydrophobic side-chains exposed to water is an important driving force behind the folding process.

Formation of intramolecular hydrogen bonds provides another important contribution to protein stability. The strength of hydrogen bonds depends on their environment, thus H-bonds enveloped in a hydrophobic core contribute more than H-bonds exposed to the aqueous environment to the stability of the native state.

The process of folding *in vivo* often begins co-translationally, so that the N-terminus of the protein begins to fold while the C-terminal portion of the protein is still being synthesized by the ribosome. Specialized proteins called chaperones assist in the folding of other proteins. A well studied example is the bacterial GroEL system, which assists in the folding of globular proteins. In eukaryotic organisms chaperones are known as heat shock proteins. Although most globular proteins are able to assume their native state unassisted, chaperone-assisted folding is often necessary in the crowded intracellular environment to prevent aggregation; chaperones are also used to prevent misfolding and aggregation which may occur as a consequence of exposure to heat or other changes in the cellular environment.

There are two models of protein folding that are currently being confirmed. The first is the diffusion collision model in which a nucleus is formed, then the secondary structure, and finally these secondary structures are collided together and pack tightly together. The next model is the nucleation-condensation model, in which the secondary and tertiary structure of the protein is made at the same time. Finally, recent studies have shown that some proteins show characteristics of both of these folding models.

For the most part, scientists have been able to study many identical molecules folding together *en masse*. At the coarsest level, it appears that in transitioning to the native state, a given amino acid sequence takes on roughly the same route and proceeds through roughly the same intermediates and transition states. Often folding involves first the establishment of regular secondary and supersecondary structures, particularly alpha helices and beta sheets, and afterwards tertiary structure. Formation of quaternary structure usually involves the "assembly" or "coassembly" of subunits that have already folded. The regular alpha helix and beta sheet structures fold rapidly because they are stabilized by intramolecular hydrogen bonds, as was first characterized by Linus Pauling. Protein folding may involve covalent bonding in the form of disulfide bridges formed between two cysteine residues or the formation of metal clusters. Shortly before settling into their more energetically favourable native conformation, molecules may pass through an intermediate "molten globule" state.

The essential fact of folding, however, remains that the amino acid sequence of each protein contains the information that specifies both the native structure and the pathway to attain that state. This is not to say that nearly identical amino acid sequences always fold similarly. Conformations differ based on environmental factors as well; similar proteins fold differently based on where they are found. Folding is a spontaneous process independent of energy inputs from nucleoside triphosphates. The passage of the folded state is mainly guided by hydrophobic interactions, formation of intramolecular hydrogen bonds, and van der Waals forces, and it is opposed by conformational entropy.

Disruption of the native state

Under some conditions proteins will not fold into their biochemically functional forms. Temperatures above or below the range that cells tend to live in will cause thermally unstable proteins to unfold or "denature" (this is why boiling makes an egg white turn opaque). High concentrations of solutes, extremes of pH, mechanical forces, and the presence of chemical denaturants can do the same. Protein thermal stability is far from constant, however. For example, hyperthermophilic bacteria have been found that grow at temperatures as high as 122 °C, which of course requires that their full complement of vital proteins and protein assemblies be stable at that temperature or above.

A fully denatured protein lacks both tertiary and secondary structure, and exists as a so-called random coil. Under certain conditions some proteins can refold; however, in many cases denaturation is irreversible. Cells sometimes protect their proteins against the denaturing influence of heat with enzymes known as chaperones or heat shock proteins, which assist other proteins both in folding and in remaining folded. Some proteins never fold in cells at all except with the assistance of chaperone molecules, which either isolate individual proteins so that their folding is not interrupted by interactions with other proteins or help to unfold misfolded proteins, giving them a second chance to refold properly. This function is crucial to prevent the risk of precipitation into insoluble amorphous aggregates.

Incorrect protein folding and neurodegenerative disease

Aggregated proteins are associated with prion-related illnesses such as Creutzfeldt-Jakob disease, bovine spongiform encephalopathy (mad cow disease), amyloid-related illnesses such as Alzheimer's disease and familial amyloid cardiomyopathy or polyneuropathy, as well as intracytoplasmic aggregation diseases such as Huntington's and Parkinson's disease. These age onset degenerative diseases are associated with the multimerization of misfolded proteins into insoluble, extracellular aggregates and/or intracellular inclusions including cross-beta sheet amyloid fibrils; it is not clear whether the aggregates are the cause or merely a reflection of the loss of protein homeostasis, the balance between synthesis, folding, aggregation and protein turnover. Misfolding and excessive degradation instead of folding and function leads to a number of proteopathy diseases such as antitrypsin-associated emphysema, cystic fibrosis and the lysosomal storage diseases, where loss of function is the origin of the disorder. While protein replacement therapy has historically been used to correct the latter disorders, an emerging approach is to use pharmaceutical chaperones to fold mutated proteins to render them functional.

Effect of external factors on the folding of Proteins

Several external factors such as temperature, external fields (electric, magnetic), , molecular crowding , limitation of space could have a big influence on the folding of proteins . Modification of the local minima by external factors can also induce modifications of the folding trajectory.

The Levinthal paradox and kinetics

The Levinthal paradox observes that if a protein were to fold by sequentially sampling all possible conformations, it would take an astronomical amount of time to do so, even if the conformations were sampled at a rapid rate (on the nanosecond or picosecond scale). Based upon the observation that proteins fold much faster than this, Levinthal then proposed that a random conformational search does not occur, and the protein must, therefore, fold through a series of meta-stable intermediate states.

The duration of the folding process varies dramatically depending on the protein of interest. When studied outside the cell, the slowest folding proteins require many minutes or hours to fold primarily due to proline isomerization, and must pass through a number of intermediate states, like checkpoints, before the process is complete. On the other hand, very small single-domain proteins with lengths of up to a hundred amino acids typically fold in a single step. Time scales of milliseconds are the norm and the very fastest known protein folding reactions are complete within a few microseconds.

Energy landscape theory of protein folding

The protein folding phenomenon was largely an experimental endeavor until the formulation of an energy landscape theory of proteins by Joseph Bryngelson and Peter Wolynes in the late 1980s and early 1990s. This approach introduced the *principle of minimal frustration*. This principle says that nature has chosen amino acid sequences so that the folded state of the protein is very stable. Additionally, the undesired interactions between amino acids along the folding pathway are reduced making the acquisition of the folded state a very fast process. Even though nature has reduced the level of *frustration* in proteins, some degree of it remains up to now as can be observed in the presence of local minima in the energy landscape of proteins. A consequence of these evolutionarily selected sequences is that proteins are generally thought to have globally "funneled energy landscapes" (coined by José Onuchic) that are largely directed towards the native state. This "folding funnel" landscape allows the protein to fold to the native state through any of a large number of pathways and intermediates, rather than being restricted to a single mechanism. The theory is supported by both computational simulations of model proteins and experimental studies, and it has been used to improve methods for protein structure prediction and design. The description of protein folding by the leveling free-energy landscape is also consistent with the 2nd law of thermodynamics. Physically, thinking of landscapes in terms of visualizable potential or total energy surfaces simply with maxima, saddle points, minima and funnels, rather like geographic landscapes, is perhaps a little misleading. The relevant description is really a highly dimensional phase space in which manifolds might take a variety of more complicated topological forms: see for example.

Techniques for studying protein folding

Circular dichroism

Circular dichroism is one of the most general and basic tools to study protein folding. Circular dichroism spectroscopy measures the absorption of circularly polarized light. In proteins, structures such as alpha helices and beta sheets are chiral, and thus absorb such light. The absorption of this light acts as a marker of the degree of foldedness of the protein ensemble. This technique can be used to measure equilibrium unfolding of the protein by measuring the change in this absorption as a function of denaturant concentration or temperature. A denaturant melt measures the free energy of unfolding as well as the protein's m value, or denaturant dependence. A temperature melt measures the melting temperature (T_m) of the protein. This type of spectroscopy can also be combined with fast-mixing devices, such as stopped flow, to measure protein folding kinetics and to generate chevron plots.

Dual Polarisation Interferometry

Dual polarisation interferometry is a surface based technique for measuring the optical properties of molecular layers. When used to characterise protein folding, it measures the conformation by determining the overall size of a monolayer of the protein and its density in real time at sub-Angstrom resolution. Although real time, measurement of the kinetics of protein folding are limited to processes that occur slower than ~ 10 Hz. Similar to circular dichroism the stimulus for folding can be a denaturant or temperature.

Vibrational circular dichroism of proteins

The more recent developments of vibrational circular dichroism (VCD) techniques for proteins, currently involving Fourier transform (FFT) instruments, provide powerful means for determining protein conformations in solution even for very large protein molecules. Such VCD studies of proteins are often combined with X-ray diffraction of protein crystals, FT-IR data for protein solutions in heavy water (D_2O), or *ab initio* quantum computations to provide unambiguous structural assignments that are unobtainable from CD.

Modern studies of folding with high time resolution

The study of protein folding has been greatly advanced in recent years by the development of fast, time-resolved techniques. These are experimental methods for rapidly triggering the folding of a sample of unfolded protein, and then observing the resulting dynamics. Fast techniques in widespread use include neutron scattering, ultrafast mixing of solutions, photochemical methods, and laser temperature jump spectroscopy. Among the many scientists who have contributed to the development of these techniques are Jeremy Cook, Heinrich Roder, Harry Gray, Martin Gruebele, Brian Dyer, William Eaton, Sheena Radford, Chris Dobson, Alan Fersht, Bengt Nölting and Lars Konermann.

Computational prediction of protein tertiary structure

De novo or *ab initio* techniques for computational protein structure prediction are related to, but strictly distinct from, studies involving protein folding. Molecular Dynamics (MD) is an important tool for studying protein folding and dynamics in silico. Because of computational cost, *ab initio* MD folding simulations with explicit water are limited to peptides and very small proteins. MD simulations of larger proteins remain restricted to dynamics of the experimental structure or its high-temperature unfolding. In order to simulate long time folding processes (beyond about 1 microsecond), like folding of small-size proteins (about 50 residues) or larger, some approximations or simplifications in protein models need to be introduced. An approach using reduced protein representation (pseudo-atoms representing groups of atoms are defined) and statistical potential is not only useful in protein structure prediction, but is also capable of reproducing the folding pathways.

There are distributed computing projects which use idle CPU or GPU time of personal computers to solve problems such as protein folding or prediction of protein structure. One such prominent example being the Folding@Home project. People can run these programs on their computer or PlayStation 3 to support them.

Experimental techniques of protein structure determination

Folded structures of proteins are routinely determined by X-ray crystallography and NMR. Dynamic methods to characterise protein folding such as dual polarisation interferometry and CD provide a measurement of conformation and conformational change rather than structure.

Chapter 10

Protein Design & Fusion Protein

Protein Design

Protein design is the design of new protein molecules, either from scratch or by making calculated variations on a known structure. The use of rational protein design techniques is a major aspect of protein engineering.

The design of minimalist computer models of proteins (lattice proteins), and the secondary structural modification of real proteins, began in the mid-1990s. The *de novo* design of real proteins became possible shortly afterwards, and the 21st century has seen the creation of small proteins with real biological functions including chiroselective catalysis, ion detection, and antiviral behaviour. There is great hope that the design of these and larger proteins will have applications in medicine and bioengineering. Recent computational redesign was capable of experimentally switching the cofactor specificity of *Candida boidinii* xylose reductase from NADPH to NADH.

Overview

The number of possible amino acid sequences is enormous, but only a subset of them will fold reliably and quickly to a single native state. Protein design involves identifying novel sequences within this subset, in particular those with a physiologically active native state. Physically, the native state of a protein is the conformational free energy minimum for the chain. Therefore protein design is the search for sequences which have the chosen structure as a free energy minimum. In a sense it is the reverse of structure prediction: a tertiary structure is specified, and a sequence is identified which will fold to it. Hence it is also referred to as *inverse folding*. Prion diseases like Mad Cow Disease are helpful examples of how important it is that designer proteins possess only one possible stable conformation. In Mad Cow Disease, there exists a healthy protein with a fatal weakness: there is another conformation this protein can "comfortably" take; the abnormally folded shape has very little free energy and is therefore very stable. For reasons that are not yet fully understood, this mis-folded prion protein has the ability to catalyze other proteins of its type to also adopt the mis-folded prion shape, which results in a disease-generating cascade of previously functional proteins quickly becoming misfolded. They lose the

ability to perform their intended function in the new conformation, and have a tendency to form aggregates called plaques. The buildup of these aggregates in the brain leads to progressive neuronal death, and eventually death of the entire organism.

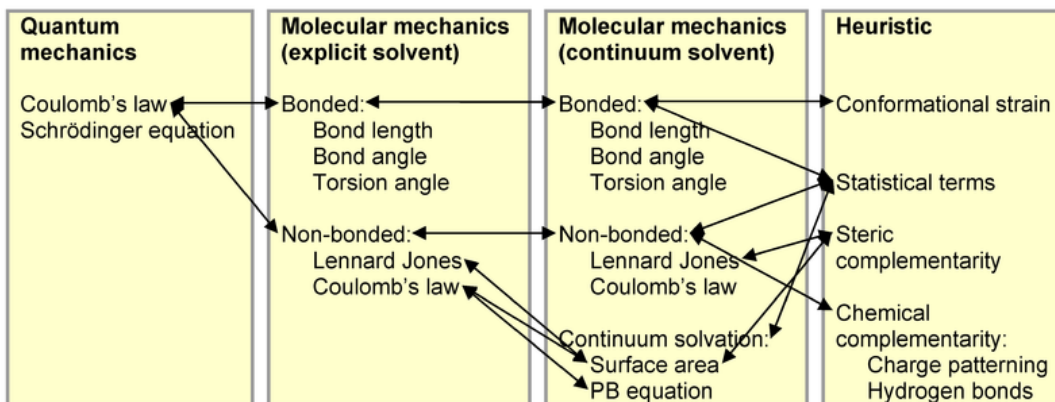
It is therefore easy to see both how important it is that a designer protein have only one possible stable tertiary structure, and that researchers exercise extreme diligence to ensure that this remains the case in all environments - especially *in vivo*.

Protein design requires an understanding of the molecular interactions that stabilize proteins in specific folded configurations; experience has shown, however, that it does not require an understanding of the dynamical process by which proteins fold,

Protein design can be accomplished using computer models, which, while simplifying the problem, are able to generate sequences that fold to the desired structure. Computational protein design algorithms search the sequence-conformation space for sequences that are low in energy when folded to the target structure. This search space is large; currently the most challenging requirement for computational protein design is a fast, yet accurate, energy function that can distinguish optimal sequences from similar suboptimal ones. Using computational methods, a protein with a novel fold has been designed, as well as sensors for unnatural molecules.

On the other hand, it is widely believed that not all possible protein structures are *designable*, which means that there are compact configurations of the chain which no sequences can fold to. In particular, conformations which are poor in secondary structures are unlikely to be designable. The designability of given structures is an issue that is still poorly understood.

Models of protein structure and function used in protein design



Comparison of various potential energy functions

Computational protein design algorithms use models of protein energetics to evaluate how mutations would affect a protein's structure and function. These energy functions typically include a combination of molecular mechanics, statistical (i.e. knowledge-

based), and other empirical terms. However, the trend has been towards using more physically based potential energy functions.

Ancestral sequence reconstruction

Ancestral reconstruction techniques have been used to design proteins with putative ancient functions.

Software

Iterative Protein Redesign and Optimization. IPRO redesigns proteins to increase or give specificity to native or novel substrates and cofactors. This is done by repeatedly randomly perturbing the backbones of the proteins around specified design positions, identifying the lowest energy combination of rotamers, and determining if the new design has a lower binding energy than previous ones. The iterative nature of this process allows IPRO to make additive mutations to the protein sequence that collectively improve the specificity towards the desired substrates and/or cofactors. Experimental testing of predictions by IPRO successfully switched the cofactor preference of *Candida boidinii* xylose reductase from NADPH to NADH. Details on how to download the software implemented in Python and experimental testing of predictions are outlined in the following paper.

EGAD: A Genetic Algorithm for protein Design. A free, open-source software package for protein design and prediction of mutation effects on protein folding stabilities and binding affinities. EGAD can also consider multiple structures simultaneously for designing specific binding proteins or locking proteins into specific conformational states. In addition to natural protein residues, EGAD can also consider free-moving ligands with or without rotatable bonds. EGAD can be used with single or multiple processors.

RosettaDesign. A software package, under active development and free for academic use, that has seen extensive successful use. RosettaDesign is accessible via a web server.

SHARPEN. A permissive open-source library for protein design and structure prediction. SHARPEN offers a variety of combinatorial optimization methods (e.g. Monte Carlo, Simulated Annealing, FASTER) and can score proteins using the successful Rosetta all-atom force field or molecular mechanics force fields (OPLSaa). In addition to the protein modeling library, SHARPEN includes tools for scalable distributed computing.

WHAT IF software for protein modelling, design, validation, and visualisation.

Abalone software for protein modelling and visualisation.

Fusion Protein

Fusion proteins or **chimeric proteins** are proteins created through the joining of two or more genes which originally coded for separate proteins. Translation of this *fusion gene* results in a single polypeptide with functional properties derived from each of the original proteins. *Recombinant fusion proteins* are created artificially by recombinant DNA technology for use in biological research or therapeutics. *Chimeric mutant proteins* occur naturally when a complex mutation, such as a chromosomal translocation, tandem duplication, or retrotransposition creates a novel coding sequence containing parts of the coding sequences from two different genes. Naturally occurring fusion proteins are commonly found in cancer cells, where they may function as oncoproteins. The bcr-abl fusion protein is a well-known example of an oncogenic fusion protein, and is considered to be the primary oncogenic driver of chronic myelogenous leukemia.

Functions

Some fusion proteins combine whole peptides and therefore contains all functional domains of the original proteins. However, other fusion proteins, especially those that are naturally occurring, combine only portions of coding sequences and therefore do not maintain the original functions of the parental genes that formed them.

Many whole gene fusions are fully functional, and can still act to replace the original peptides. Some, however, experience interactions between the two proteins that can modify their functions. Beyond these effects, some gene fusions may cause regulatory changes that alter when and where these genes act. For partial gene fusions, the shuffling of different active sites and binding domains can potentially result in new proteins with novel functions.

Recombinant technology

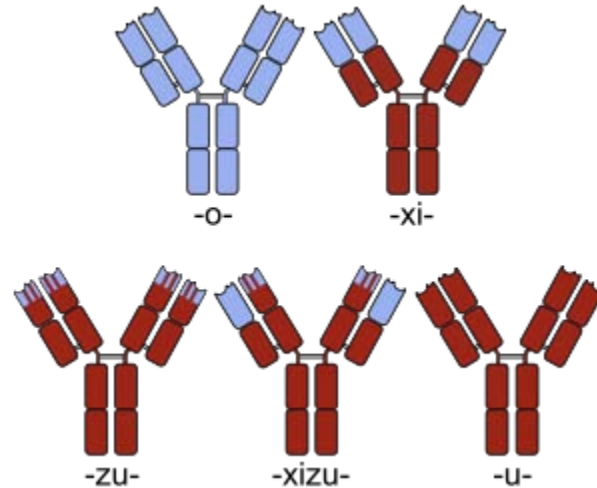
A **recombinant fusion protein** is a protein created through genetic engineering of a fusion gene. This typically involves removing the stop codon from a cDNA sequence coding for the first protein, then appending the cDNA sequence of the second protein in frame through ligation or overlap extension PCR. That DNA sequence will then be expressed by a cell as a single protein. The protein can be engineered to include the full sequence of both original proteins, or only a portion of either.

If the two entities are proteins, often linker (or "spacer") peptides are also added which make it more likely that the proteins fold independently and behave as expected.

Especially in the case where the linkers enable protein purification, linkers in protein or peptide fusions are sometimes engineered with cleavage sites for proteases or chemical agents which enable the liberation of the two separate proteins. This technique is often used for identification and purification of proteins, by fusing a GST protein, FLAG peptide, or a hexa-his peptide (6xHis-tag) which can be isolated using affinity

chromatography with nickel or cobalt resins. Fusion proteins can also be manufactured with toxins or antibodies attached to them in order to study disease development.

Chimeric protein drugs



Sketches of mouse (top left), chimeric (top right) and humanized (bottom left) monoclonal antibodies. Human parts are shown in brown, non-human parts in blue.

The purpose of creating fusion proteins in drug development is to impart properties from each of the "parent" proteins to the resulting chimeric protein. Several chimeric protein drugs are currently available for medical use.

Many chimeric protein drugs are monoclonal antibodies whose specificity for a target molecule was developed using mice and hence were initially "mouse" antibodies. As non-human proteins, mouse antibodies tend to evoke an immune reaction if administered to humans. The chimerization process involves engineering the replacement of segments of the antibody molecule that distinguish it from a human antibody. For example, human constant domains can be introduced, thereby eliminating most of the potentially immunogenic portions of the drug without altering its specificity for the intended therapeutic target. Antibody nomenclature indicates this type of modification by inserting *-xi-* into the non-proprietary name (e.g. abci-*xi*-mab). If parts of the variable domains are also replaced by human portions, *humanized* antibodies are obtained. Although not conceptually distinct from chimeras, this type is indicated using *-zu-* such as in dacli-*zu*-mab.

In addition to chimeric and humanized antibodies, there are other pharmaceutical purposes for the creation of chimeric constructs. Etanercept, for example, is a TNF α blocker created through the combination of a tumor necrosis factor receptor (TNFR) with the immunoglobulin G1 Fc segment. TNFR provides specificity for the drug target and the antibody Fc segment is believed to add stability and deliverability of the drug.

Natural occurrence

Naturally occurring fusion genes are most commonly created when a chromosomal translocation replaces the terminal exons of one gene with intact exons from a second gene. This creates a single gene which can be transcribed, spliced, and translated to produce a functional fusion protein. Many important cancer-promoting oncogenes are fusion genes produced in this way.

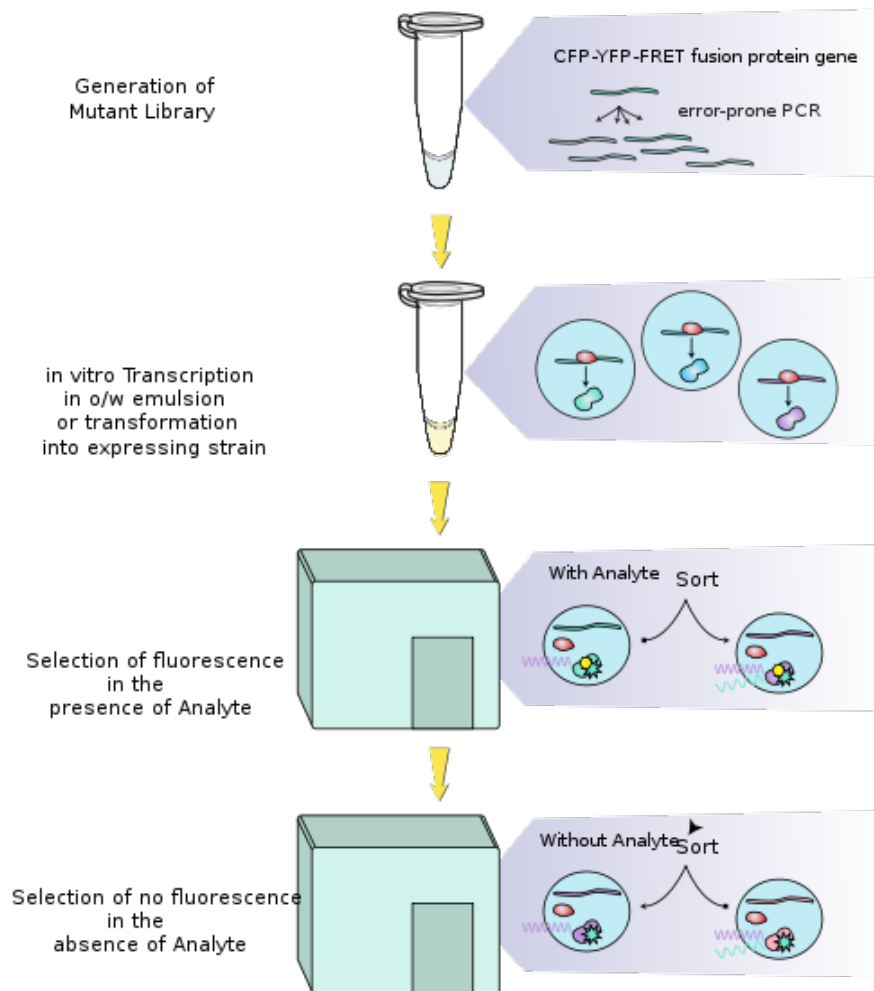
Examples include:

- Gag-onc fusion protein
- Bcr-abl fusion protein
- Tpr-met fusion protein

Antibodies are fusion proteins produced by VDJ recombination.

Chapter 11

Directed Evolution



An example of a possible round to evolve a protein based fluorescent sensor for a specific analyte using two consecutive FACS sortings

Directed evolution is a method used in protein engineering to harness the power of natural selection to evolve proteins or RNA with desirable properties not found in nature.

A typical experiment

A typical directed evolution experiment involves three steps:

1. *Diversification*: The gene encoding the protein of interest is mutated and/or recombined at random to create a large library of gene variants. Techniques commonly used in this step are error-prone PCR and DNA shuffling.
2. *Selection*: The library is tested for the presence of mutants (variants) possessing the desired property using a screen or selection. Screens enable the researcher to identify and isolate high-performing mutants by hand, while selections automatically eliminate all nonfunctional mutants.
3. *Amplification*: The variants identified in the selection or screen are replicated manyfold, enabling researchers to sequence their DNA in order to understand what mutations have occurred.

Together, these three steps are termed a "round" of directed evolution. Most experiments will perform more than one round. In these experiments, the "winners" of the previous round are diversified in the next round to create a new library. At the end of the experiment, all evolved protein or RNA mutants are characterized using biochemical methods.

Likelihood of success

The likelihood of success in a directed evolution experiment is directly related to the total library size, as evaluating more mutants increases the chances of finding one with the desired properties. Performing multiple rounds of evolution is useful not only because a new library of mutants is created in each round, but because each new library uses better mutants as templates. The experiment is analogous to climbing a hill on a landscape where elevation is a function of the desired property. The goal is to reach the summit, which represents the best mutant. Each round of selection samples mutants on all sides of the starting template and selects the mutant with the highest elevation, thereby climbing the hill. A new round samples mutants on all sides of this new template and picks the highest of these, and so on until the summit is reached.

In vivo and in vitro

Directed evolution can be performed in living cells (*in vivo* evolution) or may not involve cells at all (*in vitro* evolution). *In vivo* evolution has the advantage of selecting for properties in a cellular environment, which is useful when the evolved protein or RNA is to be used in living organisms, but *in vitro* evolution is often more versatile in the types of selections that can be performed. Furthermore, *in vitro* evolution experiments can generate larger libraries because the library DNA need not be inserted into cells, the currently limiting step.

Advantages

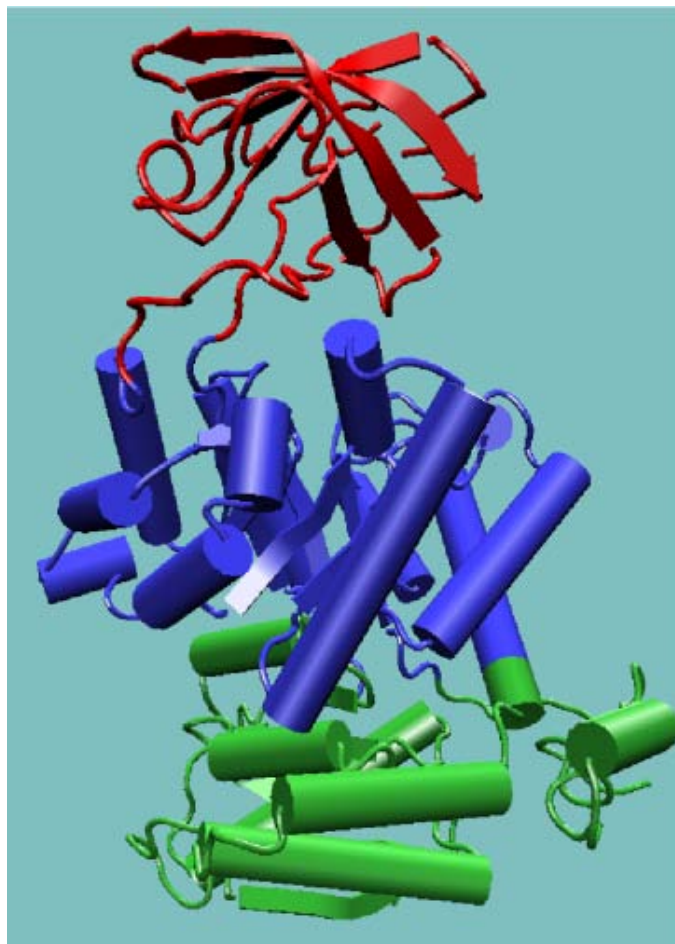
The advantage of the directed evolution approach is that the researcher need not understand the mechanism of the desired activity in order to improve it. An alternative method is rational design of site-directed mutagenesis based on X-ray crystallography data.

Uses

Most directed evolution projects seek to evolve properties that are useful to humans in an agricultural, medical or industrial context (biocatalysis). It is thus possible to use this method to optimize properties that were not selected for in the original organism. This may include catalytic activity, catalytic specificity, thermostability and many others.

Chapter 12

Protein Domain



Pyruvate kinase, a protein from three domains (PDB 1pkn)

A **protein domain** is a part of protein sequence and structure that can evolve, function, and exist independently of the rest of the protein chain. Each domain forms a compact three-dimensional structure and often can be independently stable and folded. Many proteins consist of several structural domains. One domain may appear in a variety of

different proteins. Molecular evolution uses domains as building blocks and these may be recombined in different arrangements to create proteins with different functions. Domains vary in length from between about 25 amino acids up to 500 amino acids in length. The shortest domains such as zinc fingers are stabilized by metal ions or disulfide bridges. Domains often form functional units, such as the calcium-binding EF hand domain of calmodulin. Because they are independently stable, domains can be "swapped" by genetic engineering between one protein and another to make chimeric proteins.

Background

The concept of the **domain** was first proposed in 1973 by Wetlaufer after X-ray crystallographic studies of hen lysozyme and papain and by limited proteolysis studies of immunoglobulins. Wetlaufer defined domains as stable units of protein structure that could fold autonomously. In the past domains have been described as units of:

- compact structure
- function and evolution
- folding.

Each definition is valid and will often overlap, i.e. a compact structural domain that is found amongst diverse proteins is likely to fold independently within its structural environment. Nature often brings several domains together to form multidomain and multifunctional proteins with a vast number of possibilities. In a multidomain protein, each domain may fulfil its own function independently, or in a concerted manner with its neighbours. Domains can either serve as modules for building up large assemblies such as virus particles or muscle fibres, or can provide specific catalytic or binding sites as found in enzymes or regulatory proteins.

An appropriate example is pyruvate kinase, a glycolytic enzyme that plays an important role in regulating the flux from fructose-1,6-biphosphate to pyruvate. It contains an all- β regulatory domain, an α/β -substrate binding domain and an α/β -nucleotide binding domain, connected by several polypeptide linkers. Each domain in this protein occurs in diverse sets of protein families.

The central α/β -barrel substrate binding domain is one of the most common enzyme folds. It is seen in many different enzyme families catalysing completely unrelated reactions. The α/β -barrel is commonly called the TIM barrel named after triose phosphate isomerase, which was the first such structure to be solved. It is currently classified into 26 homologous families in the CATH domain database. The TIM barrel is formed from a sequence of β - α - β motifs closed by the first and last strand hydrogen bonding together, forming an eight stranded barrel. There is debate about the evolutionary origin of this domain. One study has suggested that a single ancestral enzyme could have diverged into several families, while another suggests that a stable TIM-barrel structure has evolved through convergent evolution.

The TIM-barrel in pyruvate kinase is 'discontinuous', meaning that more than one segment of the polypeptide is required to form the domain. This is likely to be the result of the insertion of one domain into another during the protein's evolution. It has been shown from known structures that about a quarter of structural domains are discontinuous. The inserted β -barrel regulatory domain is 'continuous', made up of a single stretch of polypeptide.

Covalent association of two domains represents a functional and structural advantage since there is an increase in stability when compared with the same structures non-covalently associated. Other advantages are the protection of intermediates within inter-domain enzymatic clefts that may otherwise be unstable in aqueous environments, and a fixed stoichiometric ratio of the enzymatic activity necessary for a sequential set of reactions.

Domains are units of protein structure

The primary structure (string of amino acids) of a protein ultimately encodes its uniquely folded 3D conformation. The most important factor governing the folding of a protein into 3D structure is the distribution of polar and non-polar side chains. Folding is driven by the burial of hydrophobic side chains into the interior of the molecule so to avoid contact with the aqueous environment. Generally proteins have a core of hydrophobic residues surrounded by a shell of hydrophilic residues. Since the peptide bonds themselves are polar they are neutralised by hydrogen bonding with each other when in the hydrophobic environment. This gives rise to regions of the polypeptide that form regular 3D structural patterns called secondary structure. There are two main types of secondary structure: α -helices and β -sheets.

Some simple combinations of secondary structure elements have been found to frequently occur in protein structure and are referred to as supersecondary structure or motifs. For example, the β -hairpin motif consists of two adjacent antiparallel β -strands joined by a small loop. It is present in most antiparallel β structures both as an isolated ribbon and as part of more complex β -sheets. Another common super-secondary structure is the β - α - β motif, which is frequently used to connect two parallel β -strands. The central α -helix connects the C-termini of the first strand to the N-termini of the second strand, packing its side chains against the β -sheet and therefore shielding the hydrophobic residues of the β -strands from the surface.

Structural alignment is an important tool for determining domains.

Tertiary structure of domains

Several motifs pack together to form compact, local, semi-independent units called domains. The overall 3D structure of the polypeptide chain is referred to as the protein's tertiary structure. Domains are the fundamental units of tertiary structure, each domain containing an individual hydrophobic core built from secondary structural units connected by loop regions. The packing of the polypeptide is usually much tighter in the

interior than the exterior of the domain producing a solid-like core and a fluid-like surface. In fact, core residues are often conserved in a protein family, whereas the residues in loops are less conserved, unless they are involved in the protein's function. Protein tertiary structure can be divided into four main classes based on the secondary structural content of the domain.

- All- α domains have a domain core built exclusively from α -helices. This class is dominated by small folds, many of which form a simple bundle with helices running up and down.
- All- β domains have a core comprising of antiparallel β -sheets, usually two sheets packed against each other. Various patterns can be identified in the arrangement of the strands, often giving rise to the identification of recurring motifs, for example the Greek key motif.
- $\alpha+\beta$ domains are a mixture of all- α and all- β motifs. Classification of proteins into this class is difficult because of overlaps to the other three classes and therefore is not used in the CATH domain database.
- α/β domains are made from a combination of β - α - β motifs that predominantly form a parallel β -sheet surrounded by amphipathic α -helices. The secondary structures are arranged in layers or barrels.

Domains have limits on size

Domains have limits on size. The size of individual structural domains varies from 36 residues in E-selectin to 692 residues in lipoxygenase-1, but the majority, 90%, have less than 200 residues with an average of approximately 100 residues. Very short domains, less than 40 residues, are often stabilised by metal ions or disulfide bonds. Larger domains, greater than 300 residues, are likely to consist of multiple hydrophobic cores.

Domains and quaternary structure

Many proteins have a quaternary structure, which consists of several polypeptide chains that associate into an oligomeric molecule. Each polypeptide chain in such a protein is called a subunit. Hemoglobin, for example, consists of two α and two β subunits. Each of the four chains has an all- α globin fold with a heme pocket.

Domain swapping is a mechanism for forming oligomeric assemblies. In domain swapping, a secondary or tertiary element of a monomeric protein is replaced by the same element of another protein. Domain swapping can range from secondary structure elements to whole structural domains. It also represents a model of evolution for functional adaptation by oligomerisation, e.g. oligomeric enzymes that have their active site at subunit interfaces.

Domains as evolutionary modules

Nature is a tinkerer and not an inventor, new sequences are adapted from pre-existing sequences rather than invented. Domains are the common material used by nature to

generate new sequences, they can be thought of as genetically mobile units, referred to as 'modules'. Often, the C and N termini of domains are close together in space, allowing them to easily be "slotted into" parent structures during the process of evolution. Many domain families are found in all three forms of life, Archaea, Bacteria and Eukarya. Domains that are repeatedly found in diverse proteins are often referred to as modules, examples can be found among extracellular proteins associated with clotting, fibrinolysis, complement, the extracellular matrix, cell surface adhesion molecules and cytokine receptors.

Molecular evolution gives rise to families of related proteins with similar sequence and structure. However, sequence similarities can be extremely low between proteins that share the same structure. Protein structures may be similar because proteins have diverged from a common ancestor. Alternatively, some folds may be more favored than others as they represent stable arrangements of secondary structures and some proteins may converge towards these folds over the course of evolution. There are currently about 45,000 experimentally determined protein 3D structures deposited within the Protein Data Bank (PDB). However this set contains a lot of identical or very similar structures. All proteins should be classified to structural families to understand their evolutionary relationships. Structural comparisons are best achieved at the domain level. For this reason many algorithms have been developed to automatically assign domains in proteins with known 3D structure.

The CATH domain database classifies domains into approximately 800 fold families, ten of these folds are highly populated and are referred to as 'super-folds'. Super-folds are defined as folds for which there are at least three structures without significant sequence similarity. The most populated is the α/β -barrel super-fold as described previously.

Multidomain proteins

The majority of genomic proteins, two-thirds in unicellular organisms and more than 80% in metazoa, are multidomain proteins created as a result of gene duplication events. Many domains in multidomain structures could have once existed as independent proteins. More and more domains in eukaryotic multidomain proteins can be found as independent proteins in prokaryotes. For example, vertebrates have a multi-enzyme polypeptide containing the GAR synthetase, AIR synthetase and GAR transformylase modules (GARs-AIRs-GARt; GAR: glycinamide ribonucleotide synthetase/transferase; AIR: aminoimidazole ribonucleotide synthetase). In insects, the polypeptide appears as GARs-(AIRs)₂-GARt, in yeast GARs-AIRs is encoded separately from GARt, and in bacteria each domain is encoded separately.

Origin

Multidomain proteins are likely to have emerged from a selective pressure during evolution to create new functions. Various proteins have diverged from common ancestors by different combinations and associations of domains. Modular units

frequently move about, within and between biological systems through mechanisms of genetic shuffling:

- transposition of mobile elements including horizontal transfers (between species);
- gross rearrangements such as inversions, translocations, deletions and duplications;
- homologous recombination;
- slippage of DNA polymerase during replication.

Types of organisation

The simplest multidomain organisation seen in proteins is that of a single domain repeated in tandem. The domains may interact with each other or remain isolated, like beads on string. The giant 30,000 residue muscle protein titin comprises about 120 fibronectin-III-type and Ig-type domains. In the serine proteases, a gene duplication event has led to the formation of a two β -barrel domain enzyme. The repeats have diverged so widely that there is no obvious sequence similarity between them. The active site is located at a cleft between the two β -barrel domains, in which functionally important residues are contributed from each domain. Genetically engineered mutants of the chymotrypsin serine protease were shown to have some proteinase activity even though their active site residues were abolished and it has therefore been postulated that the duplication event enhanced the enzyme's activity.

Modules frequently display different connectivity relationships, as illustrated by the kinesins and ABC transporters. The kinesin motor domain can be at either end of a polypeptide chain that includes a coiled-coil region and a cargo domain. ABC transporters are built with up to four domains consisting of two unrelated modules, ATP-binding cassette and an integral membrane module, arranged in various combinations.

Not only do domains recombine, but there are many examples of a domain having been inserted into another. Sequence or structural similarities to other domains demonstrate that homologues of inserted and parent domains can exist independently. An example is that of the 'fingers' inserted into the 'palm' domain within the polymerases of the Pol I family. Since a domain can be inserted into another, there should always be at least one continuous domain in a multidomain protein. This is the main difference between definitions of structural domains and evolutionary/functional domains. An evolutionary domain will be limited to one or two connections between domains, whereas structural domains can have unlimited connections, within a given criterion of the existence of a common core. Several structural domains could be assigned to an evolutionary domain.

Domains are autonomous folding units

Folding

Protein folding - the unsolved problem : Since the seminal work of Anfinsen over forty years ago, the goal to completely understand the mechanism by which a polypeptide

rapidly folds into its stable native conformation remains elusive. Many experimental folding studies have contributed much to our understanding, but the principles that govern protein folding are still based on those discovered in the very first studies of folding. Anfinsen showed that the native state of a protein is thermodynamically stable, the conformation being at a global minimum of its free energy.

Folding is a directed search of conformational space allowing the protein to fold on a biologically feasible time scale. The Levinthal paradox states that if an averaged sized protein would sample all possible conformations before finding the one with the lowest energy, the whole process would take billions of years. Proteins typically fold within 0.1 and 1000 seconds, therefore the protein folding process must be directed some way through a specific folding pathway. The forces that direct this search are likely to be a combination of local and global influences whose effects are felt at various stages of the reaction.

Advances in experimental and theoretical studies have shown that folding can be viewed in terms of energy landscapes, where folding kinetics is considered as a progressive organisation of an ensemble of partially folded structures through which a protein passes on its way to the folded structure. This has been described in terms of a folding funnel, in which an unfolded protein has a large number of conformational states available and there are fewer states available to the folded protein. A funnel implies that for protein folding there is a decrease in energy and loss of entropy with increasing tertiary structure formation. The local roughness of the funnel reflects kinetic traps, corresponding to the accumulation of misfolded intermediates. A folding chain progresses toward lower intra-chain free-energies by increasing its compactness. The chains conformational options become increasingly narrowed ultimately toward one native structure.

Advantage of domains in protein folding

The organisation of large proteins by structural domains represents an advantage for protein folding, with each domain being able to individually fold, accelerating the folding process and reducing a potentially large combination of residue interactions. Furthermore, given the observed random distribution of hydrophobic residues in proteins, domain formation appears to be the optimal solution for a large protein to bury its hydrophobic residues while keeping the hydrophilic residues at the surface.

However, the role of inter-domain interactions in protein folding and in energetics of stabilisation of the native structure, probably differs for each protein. In T4 lysozyme, the influence of one domain on the other is so strong that the entire molecule is resistant to proteolytic cleavage. In this case, folding is a sequential process where the C-terminal domain is required to fold independently in an early step, and the other domain requires the presence of the folded C-terminal domain for folding and stabilisation.

It has been found that the folding of an isolated domain can take place at the same rate or sometimes faster than that of the integrated domain. Suggesting that unfavourable interactions with the rest of the protein can occur during folding. Several arguments

suggest that the slowest step in the folding of large proteins is the pairing of the folded domains. This is either because the domains are not folded entirely correctly or because the small adjustments required for their interaction are energetically unfavourable, such as the removal of water from the domain interface.

Domains and protein flexibility

The presence of multiple domains in proteins gives rise to a great deal of flexibility and mobility, leading to **protein domain dynamics**. Domain motions can be inferred by comparing different structures of a protein, or they can be directly observed using spectra measured by neutron spin echo spectroscopy. They can also be suggested by sampling in extensive molecular dynamics trajectories. Domain motions are important for:

- catalysis;
- regulatory activity;
- transport of metabolites;
- formation of protein assemblies; and
- cellular locomotion.

One of the largest observed domain motions is the 'swivelling' mechanism in pyruvate phosphate dikinase. The phosphoinositide domain swivels between two states in order to bring a phosphate group from the active site of the nucleotide binding domain to that of the phosphoenolpyruvate/pyruvate domain. The phosphate group is moved over a distance of 45Å involving a domain motion of about 100 degrees around a single residue. In enzymes, the closure of one domain onto another captures a substrate by an induced fit, allowing the reaction to take place in a controlled way. A detailed analysis by Gerstein led to the classification of two basic types of domain motion; hinge and shear. Only a relatively small portion of the chain, namely the inter-domain linker and side chains undergo significant conformational changes upon domain rearrangement.

Hinges by secondary structures

A study by Hayward found that the termini of α -helices and β -sheets form hinges in a large number of cases. Many hinges were found to involve two secondary structure elements acting like hinges of a door, allowing an opening and closing motion to occur. This can arise when two neighbouring strands within a β -sheet situated in one domain, diverge apart as they join the other domain. The two resulting termini then form the bending regions between the two domains. α -helices that preserve their hydrogen bonding network when bent are found to behave as mechanical hinges, storing 'elastic energy' that drives the closure of domains for rapid capture of a substrate.

Helical to extended conformation

The interconversion of helical and extended conformations at the site of a domain boundary is not uncommon. In calmodulin, torsion angles change for five residues in the

middle of a domain linking α -helix. The helix is split into two, almost perpendicular, smaller helices separated by four residues of an extended strand.

Shear motions

Shear motions involve a small sliding movement of domain interfaces, controlled by the amino acid side chains within the interface. Proteins displaying shear motions often have a layered architecture: stacking of secondary structures. The interdomain linker has merely the role of keeping the domains in close proximity.

Domain motion and functional dynamics in enzymes

The analysis of the internal dynamics of structurally different, but functionally similar enzymes has highlighted a common relationship between the positioning of the active site and the two principal protein sub-domains. In fact, for several members of the hydrolase superfamily, the catalytic site is located close to the interface separating the two principal quasi-rigid domains. Such positioning appears instrumental for maintaining the precise geometry of the active site, while allowing for an appreciable functionally-oriented modulation of the flanking regions resulting from the relative motion of the two sub-domains.

Domain definition from structural co-ordinates

The importance of domains as structural building blocks and elements of evolution has brought about many automated methods for their identification and classification in proteins of known structure. Automatic procedures for reliable domain assignment is essential for the generation of the domain databases, especially as the number of protein structures is increasing. Although the boundaries of a domain can be determined by visual inspection, construction of an automated method is not straightforward. Problems occur when faced with domains that are discontinuous or highly associated. The fact that there is no standard definition of what a domain really is has meant that domain assignments have varied enormously, with each researcher using a unique set of criteria.

A structural domain is a compact, globular sub-structure with more interactions within it than with the rest of the protein. Therefore, a structural domain can be determined by two visual characteristics; its compactness and its extent of isolation. Measures of local compactness in proteins have been used in many of the early methods of domain assignment and in several of the more recent methods.

Methods

One of the first algorithms used a $C\alpha$ - $C\alpha$ distance map together with a hierarchical clustering routine that considered proteins as several small segments, 10 residues in length. The initial segments were clustered one after another based on inter-segment distances; segments with the shortest distances were clustered and considered as single segments thereafter. The stepwise clustering finally included the full protein. Go also

exploited the fact that inter-domain distances are normally larger than intra-domain distances; all possible C α -C α distances were represented as diagonal plots in which there were distinct patterns for helices, extended strands and combinations of secondary structures.

The method by Sowdhamini and Blundell clusters secondary structures in a protein based on their C α -C α distances and identifies domains from the pattern in their dendrograms. As the procedure does not consider the protein as a continuous chain of amino acids there are no problems in treating discontinuous domains. Specific nodes in these dendrograms are identified as tertiary structural clusters of the protein, these include both super-secondary structures and domains. The DOMAK algorithm is used to create the 3Dee domain database. It calculates a 'split value' from the number of each type of contact when the protein is divided arbitrarily into two parts. This split value is large when the two parts of the structure are distinct.

The method of Wodak and Janin was based on the calculated interface areas between two chain segments repeatedly cleaved at various residue positions. Interface areas were calculated by comparing surface areas of the cleaved segments with that of the native structure. Potential domain boundaries can be identified at a site where the interface area was at a minimum. Other methods have used measures of solvent accessibility to calculate compactness.

The PUU algorithm incorporates a harmonic model used to approximate inter-domain dynamics. The underlying physical concept is that many rigid interactions will occur within each domain and loose interactions will occur between domains. This algorithm is used to define domains in the FSSP domain database.

Swindells (1995) developed a method, DETECTIVE, for identification of domains in protein structures based on the idea that domains have a hydrophobic interior. Deficiencies were found to occur when hydrophobic cores from different domains continue through the interface region.

RigidFinder is a novel method for identification of protein rigid blocks (domains and loops) from two different conformations. Rigid blocks are defined as blocks where all inter residue distances are conserved across conformations.

A general method to identify *dynamical domains*, that is protein regions that behave approximately as rigid units in the course of structural fluctuations, has been introduced by Potestio et al. and, among other applications was also used to compare the consistency of the dynamics-based domain subdivisions with standard structure-based ones. The method, termed PiSQRD, is publicly available in the form of a webserver. The latter allows users to optimally subdivide single-chain or multimeric proteins into quasi-rigid domains based on the collective modes of fluctuation of the system. By default the latter are calculated through an elastic network model; alternatively pre-calculated essential dynamical spaces can be uploaded by the user.

Chapter 13

Proteomics



Robotic preparation of MALDI mass spectrometry samples on a sample carrier.

Proteomics is the large-scale study of proteins, particularly their structures and functions. Proteins are vital parts of living organisms, as they are the main components of the physiological metabolic pathways of cells. The term "proteomics" was first coined in 1997 to make an analogy with genomics, the study of the genes. The word "proteome" is a blend of "**protein**" and "**genome**", and was coined by Marc Wilkins in 1994 while working on the concept as a PhD student. The proteome is the entire complement of

proteins, including the modifications made to a particular set of proteins, produced by an organism or system. This will vary with time and distinct requirements, or stresses, that a cell or organism undergoes.

Complexity of the problem

After genomics, proteomics is considered the next step in the study of biological systems. It is much more complicated than genomics mostly because while an organism's genome is more or less constant, the proteome differs from cell to cell and from time to time. This is because distinct genes are expressed in distinct cell types. This means that even the basic set of proteins which are produced in a cell needs to be determined.

In the past this was done by mRNA analysis, but this was found not to correlate with protein content. It is now known that mRNA is not always translated into protein, and the amount of protein produced for a given amount of mRNA depends on the gene it is transcribed from and on the current physiological state of the cell. Proteomics confirms the presence of the protein and provides a direct measure of the quantity present.

Post-translational modifications

Not only does the translation from mRNA cause differences, many proteins are also subjected to a wide variety of chemical modifications after translation. Many of these post-translational modifications are critical to the protein's function.

Phosphorylation

One such modification is phosphorylation, which happens to many enzymes and structural proteins in the process of cell signaling. The addition of a phosphate to particular amino acids—most commonly serine and threonine mediated by serine/threonine kinases, or more rarely tyrosine mediated by tyrosine kinases—causes a protein to become a target for binding or interacting with a distinct set of other proteins that recognize the phosphorylated domain.

Because protein phosphorylation is one of the most-studied protein modifications, many "proteomic" efforts are geared to determining the set of phosphorylated proteins in a particular cell or tissue-type under particular circumstances. This alerts the scientist to the signaling pathways that may be active in that instance.

Ubiquitination

Ubiquitin is a small protein that can be affixed to certain protein substrates by enzymes called E3 ubiquitin ligases. Determining which proteins are poly-ubiquitinated can be helpful in understanding how protein pathways are regulated. This is therefore an additional legitimate "proteomic" study. Similarly, once it is determined what substrates are ubiquitinated by each ligase, determining the set of ligases expressed in a particular cell type will be helpful.

Additional modifications

Listing all the protein modifications that might be studied in a "Proteomics" project would require a discussion of most of biochemistry; therefore, a short list will serve here to illustrate the complexity of the problem. In addition to phosphorylation and ubiquitination, proteins can be subjected to (among others) methylation, acetylation, glycosylation, oxidation and nitrosylation. Some proteins undergo ALL of these modifications, often in time-dependent combinations, aptly illustrating the potential complexity one has to deal with when studying protein structure and function.

Distinct proteins are made under distinct settings

Even if one is studying a particular cell type, that cell may make different sets of proteins at different times, or under different conditions. Furthermore, as mentioned, any one protein can undergo a wide range of post-translational modifications.

Therefore a "proteomics" study can become quite complex very quickly, even if the object of the study is very restricted. In more ambitious settings, such as when a biomarker for a tumor is sought - when the proteomics scientist is obliged to study sera samples from multiple cancer patients - the amount of complexity that must be dealt with is as great as in any modern biological project.

Limitations to genomic study

Scientists are very interested in proteomics because it gives a much better understanding of an organism than genomics. First, the level of transcription of a gene gives only a rough estimate of its level of expression into a protein. An mRNA produced in abundance may be degraded rapidly or translated inefficiently, resulting in a small amount of protein. Second, as mentioned above many proteins experience post-translational modifications that profoundly affect their activities; for example some proteins are not active until they become phosphorylated. Methods such as phosphoproteomics and glycoproteomics are used to study post-translational modifications. Third, many transcripts give rise to more than one protein, through alternative splicing or alternative post-translational modifications. Fourth, many proteins form complexes with other proteins or RNA molecules, and only function in the presence of these other molecules. Finally, protein degradation rate plays an important role in protein content.

Methods of studying proteins

Determining proteins which are post-translationally modified

One way in which a particular protein can be studied is to develop an antibody which is specific to that modification. For example, there are antibodies which only recognize certain proteins when they are tyrosine-phosphorylated, known as phospho-specific antibodies; also, there are antibodies specific to other modifications. These can be used to determine the set of proteins that have undergone the modification of interest.

For sugar modifications, such as glycosylation of proteins, certain lectins have been discovered which bind sugars. These too can be used.

A more common way to determine post-translational modification of interest is to subject a complex mixture of proteins to electrophoresis in "two-dimensions", which simply means that the proteins are electrophoresed first in one direction, and then in another... this allows small differences in a protein to be visualized by separating a modified protein from its unmodified form. This methodology is known as "two-dimensional gel electrophoresis".

Recently, another approach has been developed called PROTOMAP which combines SDS-PAGE with shotgun proteomics to enable detection of changes in gel-migration such as those caused by proteolysis or post translational modification.

Determining the existence of proteins in complex mixtures

Classically, antibodies to particular proteins or to their modified forms have been used in biochemistry and cell biology studies. These are among the most common tools used by practicing biologists today.

For more quantitative determinations of protein amounts, techniques such as ELISAs can be used.

For proteomic study, more recent techniques such as matrix-assisted laser desorption/ionization (MALDI) have been employed for rapid determination of proteins in particular mixtures and increasingly electrospray ionization (ESI).

Computational methods in studying protein biomarkers

Computational predictive models have shown that extensive and diverse fetomaternal protein trafficking occurs during pregnancy and can be readily detected non-invasively in maternal whole blood. This computational approach circumvented a major limitation, the abundance of maternal proteins interfering with the detection of fetal proteins, to fetal proteomic analysis of maternal blood. Computational models can use fetal gene transcripts previously identified in maternal whole blood to create a comprehensive proteomic network of the term neonate. Such work shows that the fetal proteins detected in pregnant woman's blood originate from a diverse group of tissues and organs from the developing fetus. The proteomic networks contain many biomarkers that are proxies for development and illustrate the potential clinical application of this technology as a way to monitor normal and abnormal fetal development.

An information theoretic framework has also been introduced for biomarker discovery, integrating biofluid and tissue information. This new approach takes advantage of functional synergy between certain biofluids and tissues with the potential for clinically significant findings not possible if tissues and biofluids were considered individually. By conceptualizing tissue-biofluid as information channels, significant biofluid proxies can

be identified and then used for guided development of clinical diagnostics. Candidate biomarkers are then predicted based on information transfer criteria across the tissue-biofluid channels. Significant biofluid-tissue relationships can be used to prioritize clinical validation of biomarkers.

Establishing protein-protein interactions

Most proteins function in collaboration with other proteins, and one goal of proteomics is to identify which proteins interact. This is especially useful in determining potential partners in cell signaling cascades.

Several methods are available to probe protein-protein interactions. The traditional method is yeast two-hybrid analysis. New methods include protein microarrays, immunoaffinity chromatography followed by mass spectrometry, dual polarisation interferometry, Microscale Thermophoresis and experimental methods such as phage display and computational methods

Practical applications of proteomics

One of the most promising developments to come from the study of human genes and proteins has been the identification of potential new drugs for the treatment of disease. This relies on genome and proteome information to identify proteins associated with a disease, which computer software can then use as targets for new drugs. For example, if a certain protein is implicated in a disease, its 3D structure provides the information to design drugs to interfere with the action of the protein. A molecule that fits the active site of an enzyme, but cannot be released by the enzyme, will inactivate the enzyme. This is the basis of new drug-discovery tools, which aim to find new drugs to inactivate proteins involved in disease. As genetic differences among individuals are found, researchers expect to use these techniques to develop personalized drugs that are more effective for the individual.

Biomarkers

The FDA defines a biomarker as, “A characteristic that is objectively measured and evaluated as an indicator of normal biologic processes, pathogenic processes, or pharmacologic responses to a therapeutic intervention”.

Understanding the proteome, the structure and function of each protein and the complexities of protein-protein interactions will be critical for developing the most effective diagnostic techniques and disease treatments in the future.

An interesting use of proteomics is using specific protein biomarkers to diagnose disease. A number of techniques allow to test for proteins produced during a particular disease, which helps to diagnose the disease quickly. Techniques include western blot, immunohistochemical staining, enzyme linked immunosorbent assay (ELISA) or mass spectrometry.

Current research methodologies

There are many approaches to attempting to characterize the human proteome, which is estimated to exceed 100,000 unique forms, 25,000 genes plus post-translational modifications.

In addition, first promising attempts to decipher the proteom of animal tumors have recently been reported.

Chapter 14

Proteinogenic Amino Acid

Proteinogenic amino acids are those amino acids that can be found in proteins and require cellular machinery coded for in the genetic code of any organism for their isolated production. There are 22 standard amino acids, but only 21 are found in eukaryotes. Of the 22, 20 are directly encoded by the universal genetic code. Humans can synthesize 11 of these 20 from each other or from other molecules of intermediary metabolism. The other 9 must be consumed in the diet, and so are called *essential amino acids*; those are histidine, isoleucine, leucine, lysine, methionine, phenylalanine, threonine, tryptophan, and valine. The remaining two, selenocysteine and pyrrolysine, are incorporated into proteins by unique synthetic mechanisms.

The word *proteinogenic* means "protein building". Proteinogenic amino acids can be assembled into a polypeptide (the subunit of a protein) through a process called translation (the second stage of protein biosynthesis, part of the overall process of gene expression).

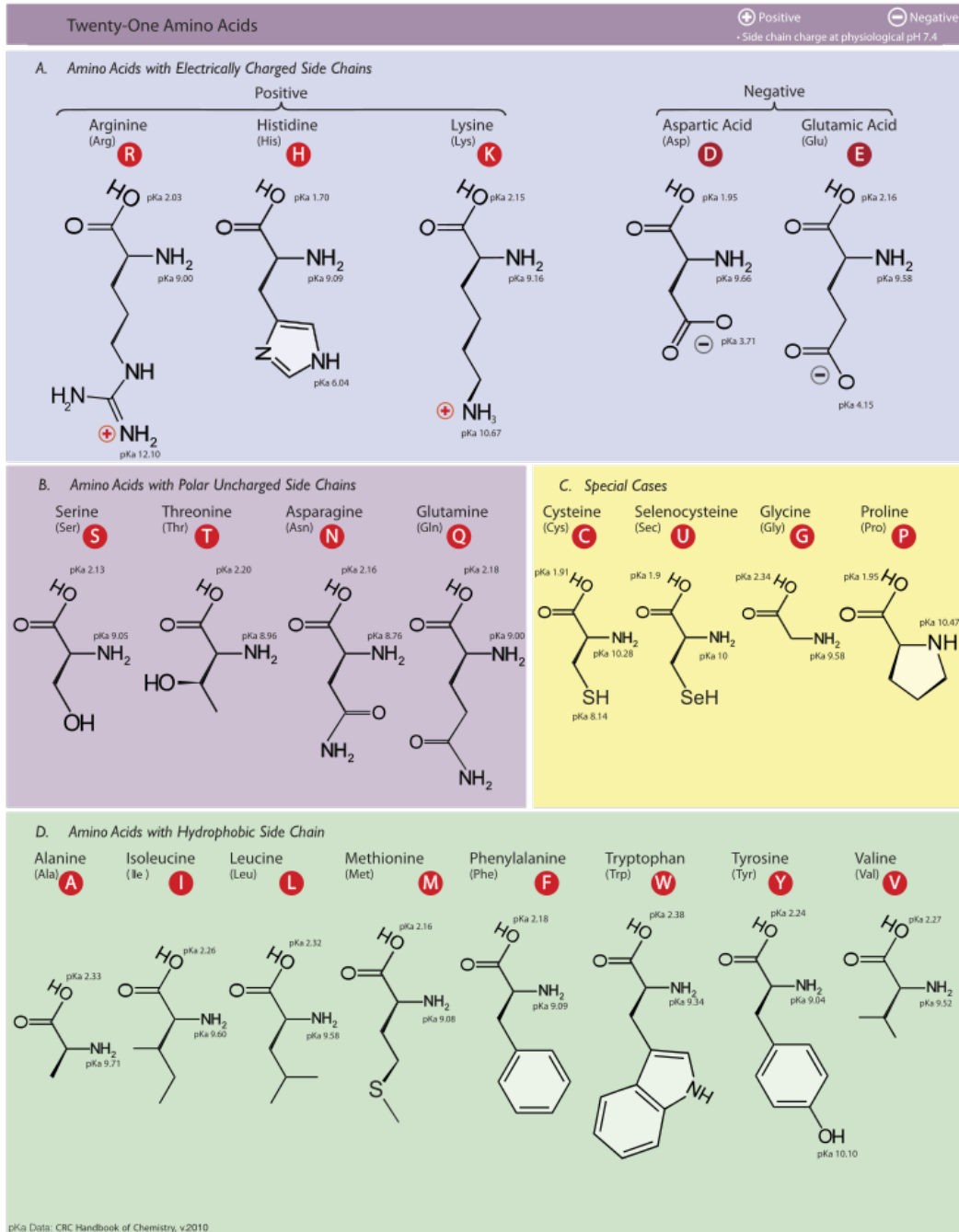
Non-proteinogenic amino acids are either not found in proteins (like carnitine, GABA, or L-DOPA), or are not produced directly and in isolation by standard cellular machinery (like hydroxyproline and selenomethionine). The latter often results from posttranslational modification of proteins.

There are clear reasons why organisms have not evolved to incorporate certain non-proteinogenic amino acids into proteins: for example, ornithine and homoserine cyclize against the peptide backbone and fragment the protein with relatively short half-lives, while others are toxic because they can be mistakenly incorporated into proteins, such as the arginine analog canavanine.

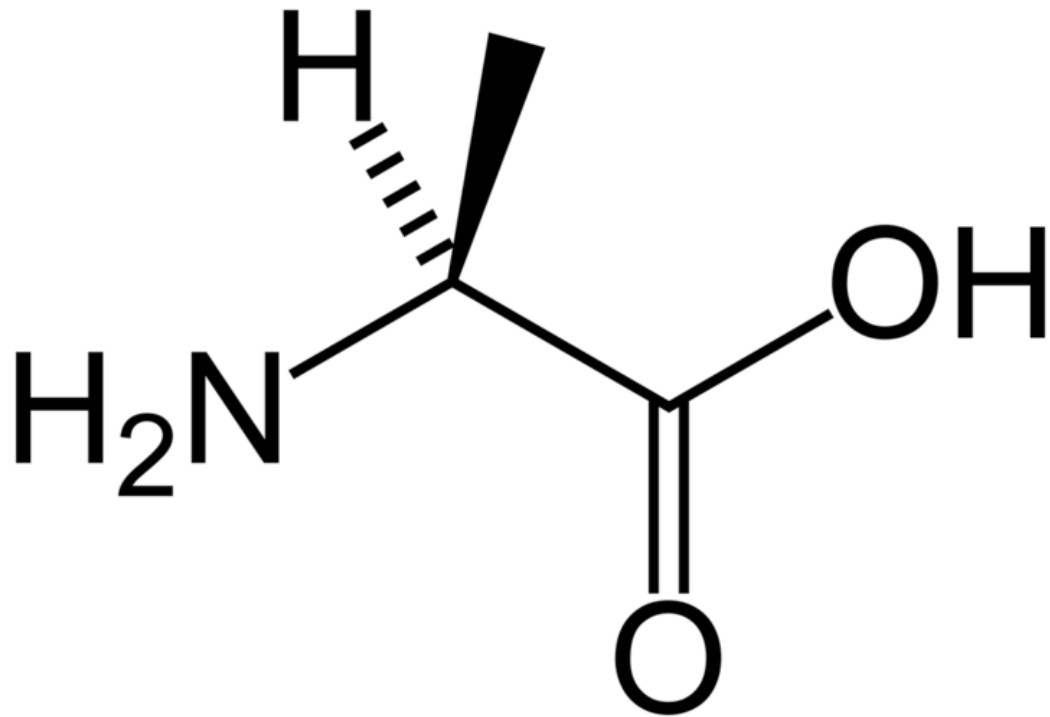
Non-proteinogenic amino acids are found in nonribosomal peptides, which are not produced by the ribosome during translation.

Structures

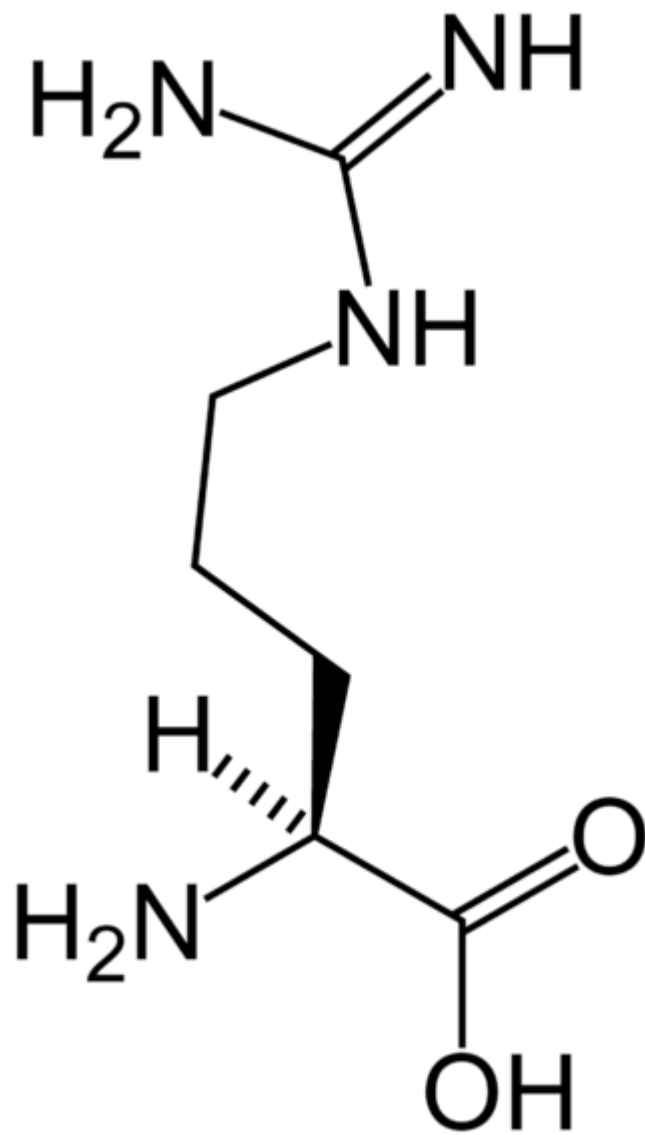
The following illustrates the structures and abbreviations of the 21 amino acids that are directly encoded for protein synthesis by the genetic code of eukaryotes. The structures given below are standard chemical structures, not the typical zwitterion forms that exist in aqueous solutions.



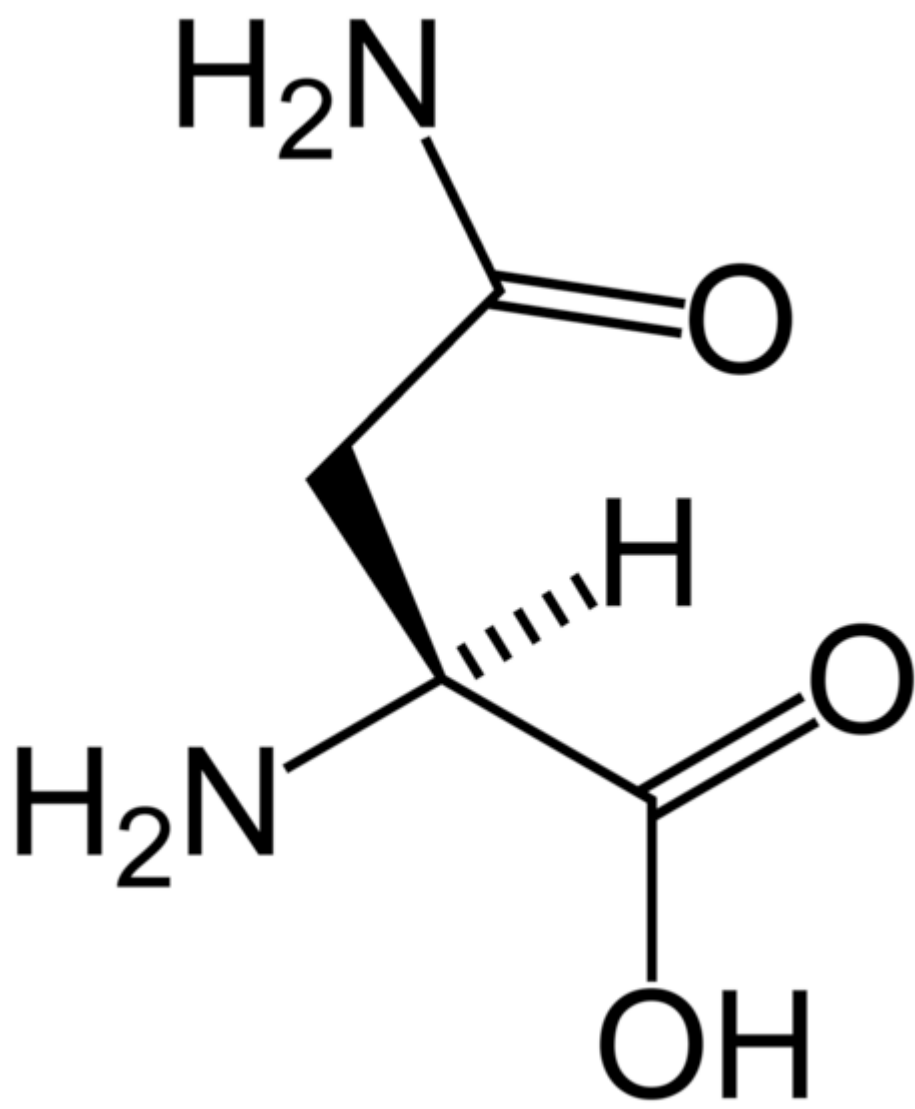
Grouped table of 21 amino acids' structures, nomenclature, and their side groups' pKa's.



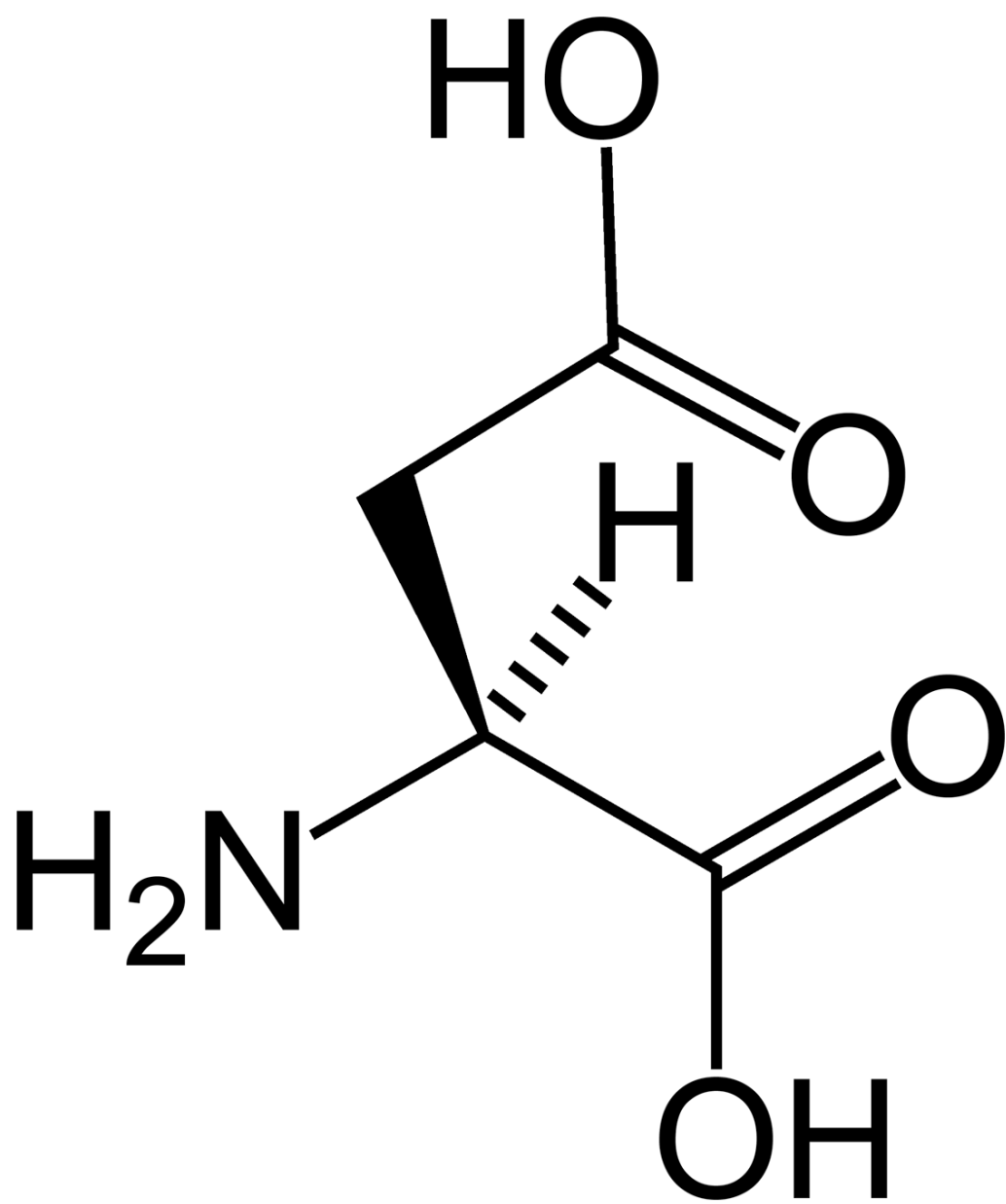
L-Alanine
(Ala / A)



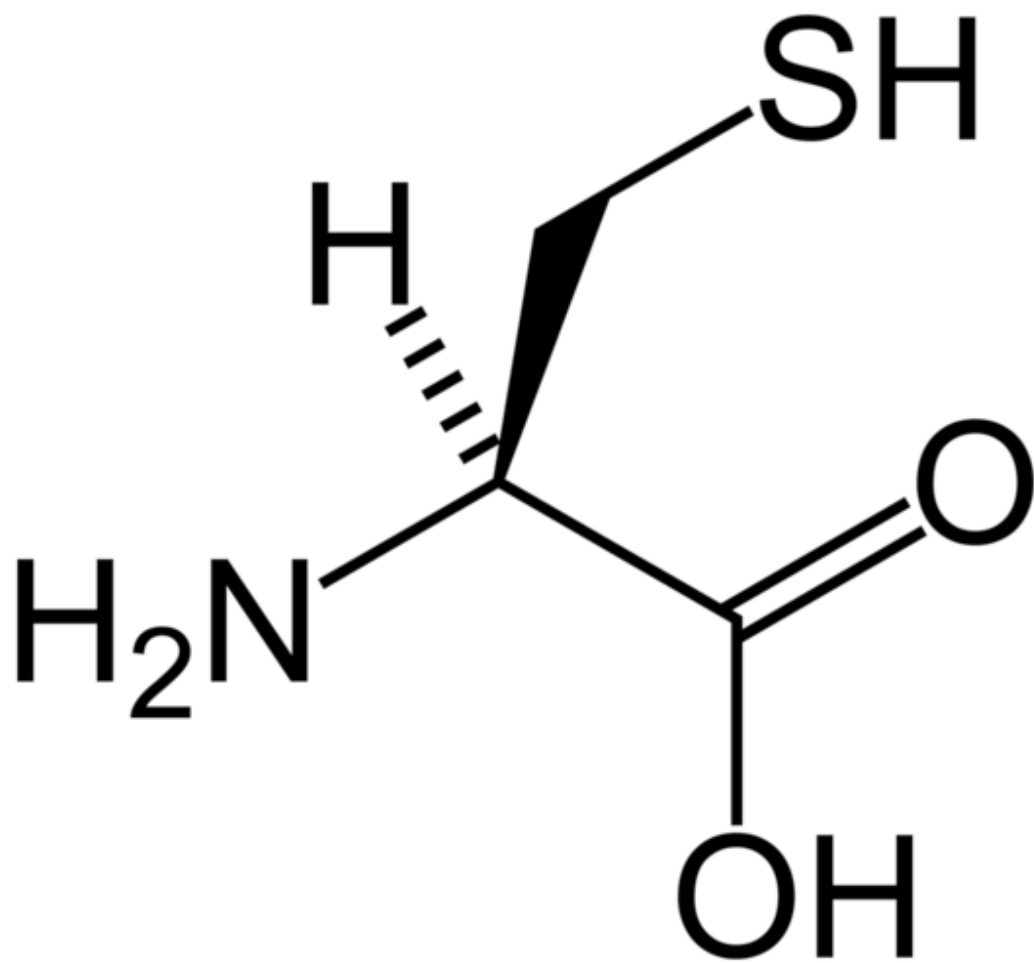
L-Arginine
(Arg / R)



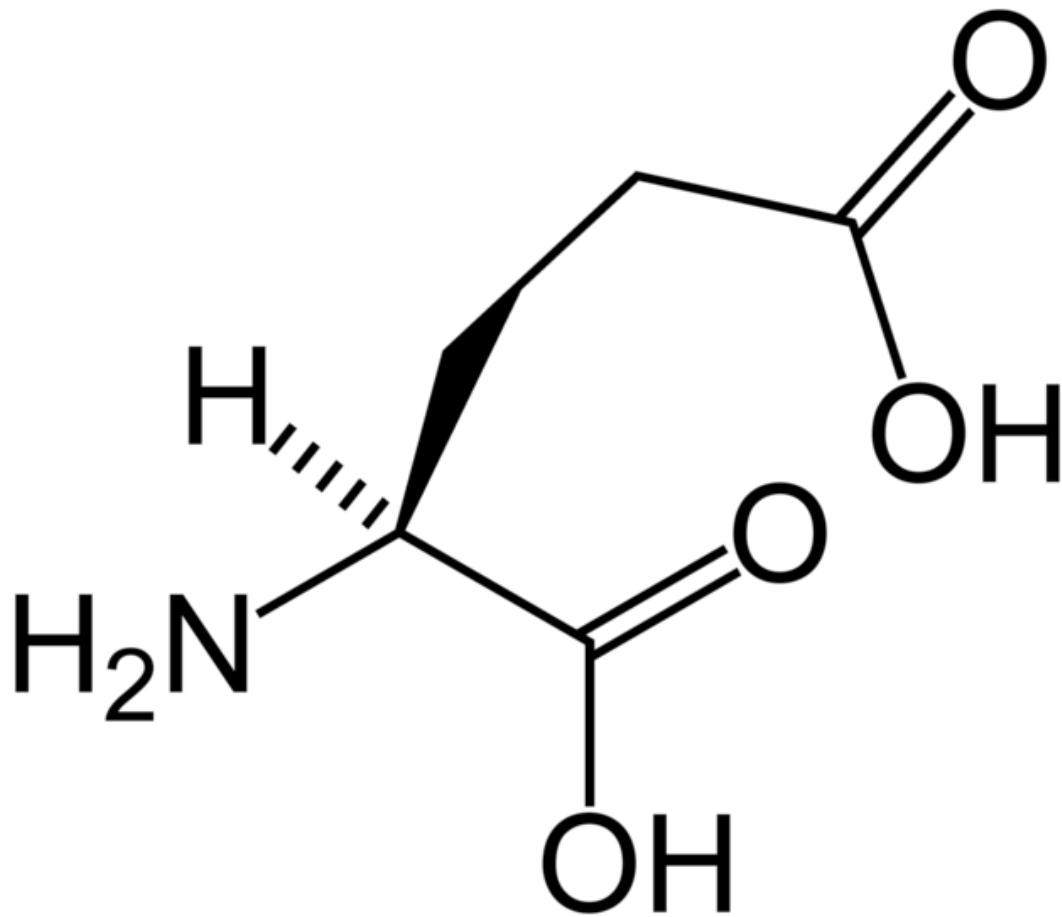
L-Asparagine
(Asn / N)



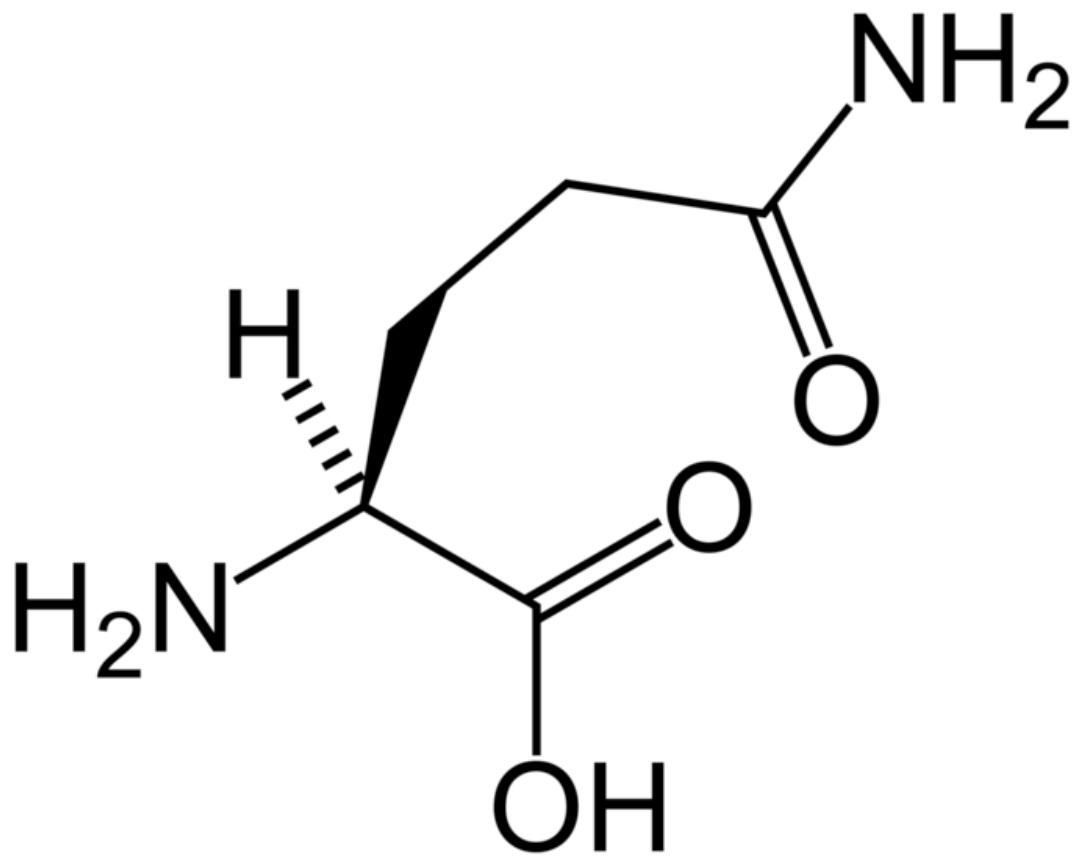
L-Aspartic acid
(Asp / D)



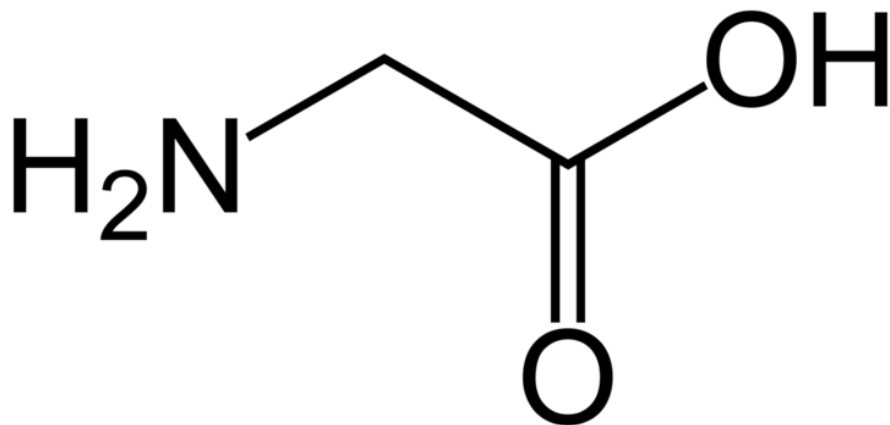
L-Cysteine
(Cys / C)



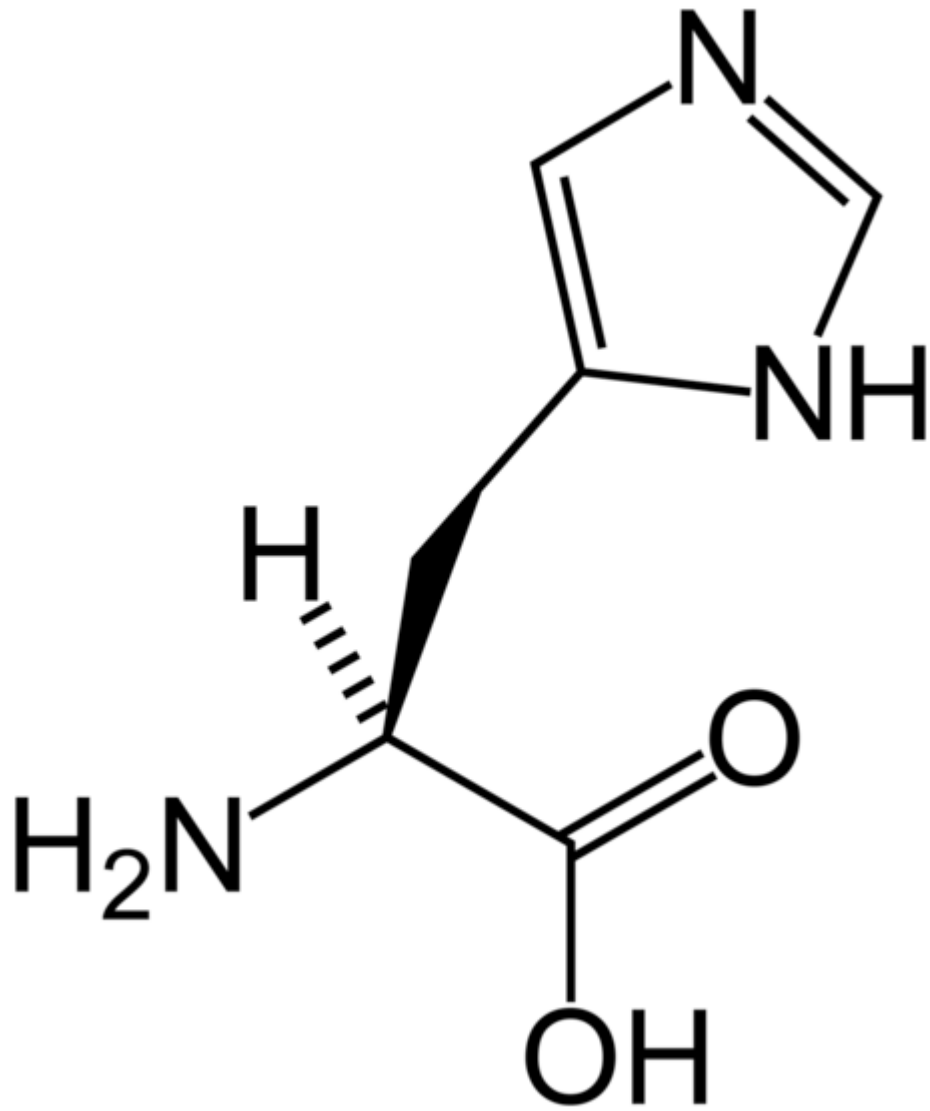
L-Glutamic acid
(Glu / E)



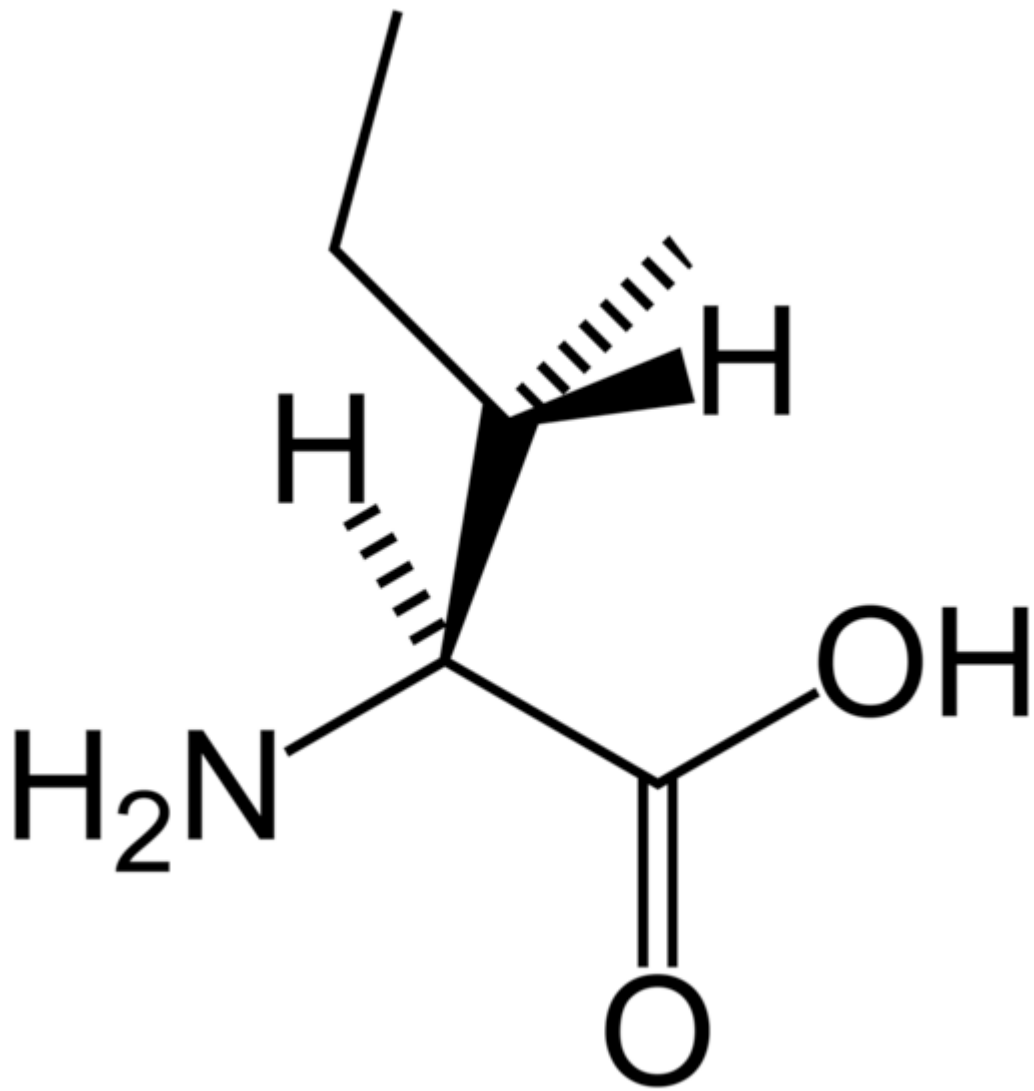
L-Glutamine
(Gln / Q)



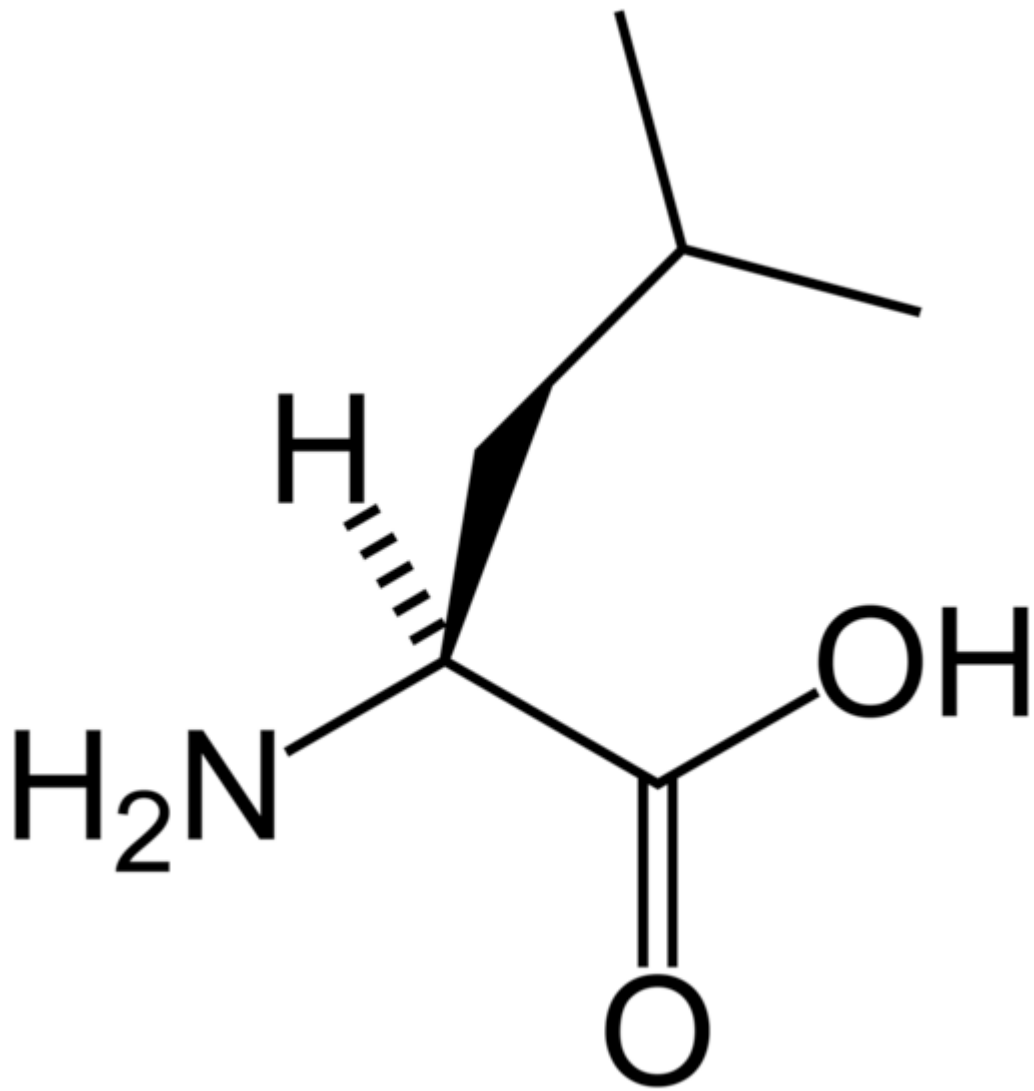
Glycine
(Gly / G)



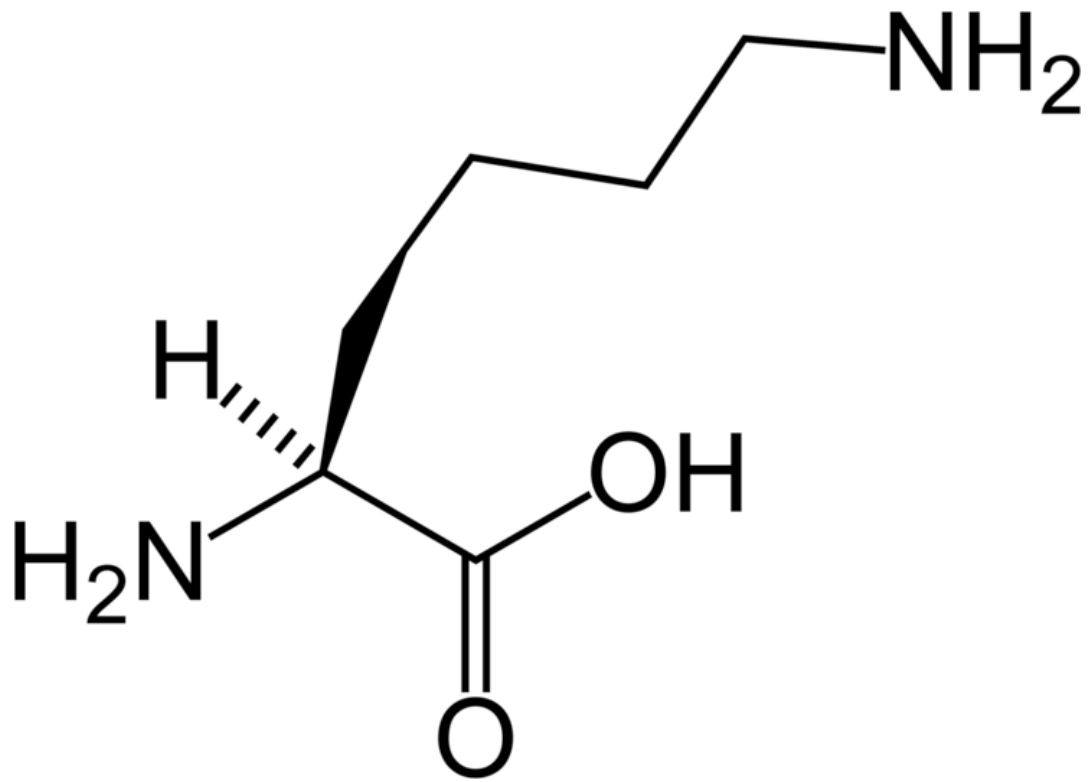
L-Histidine
(His / H)



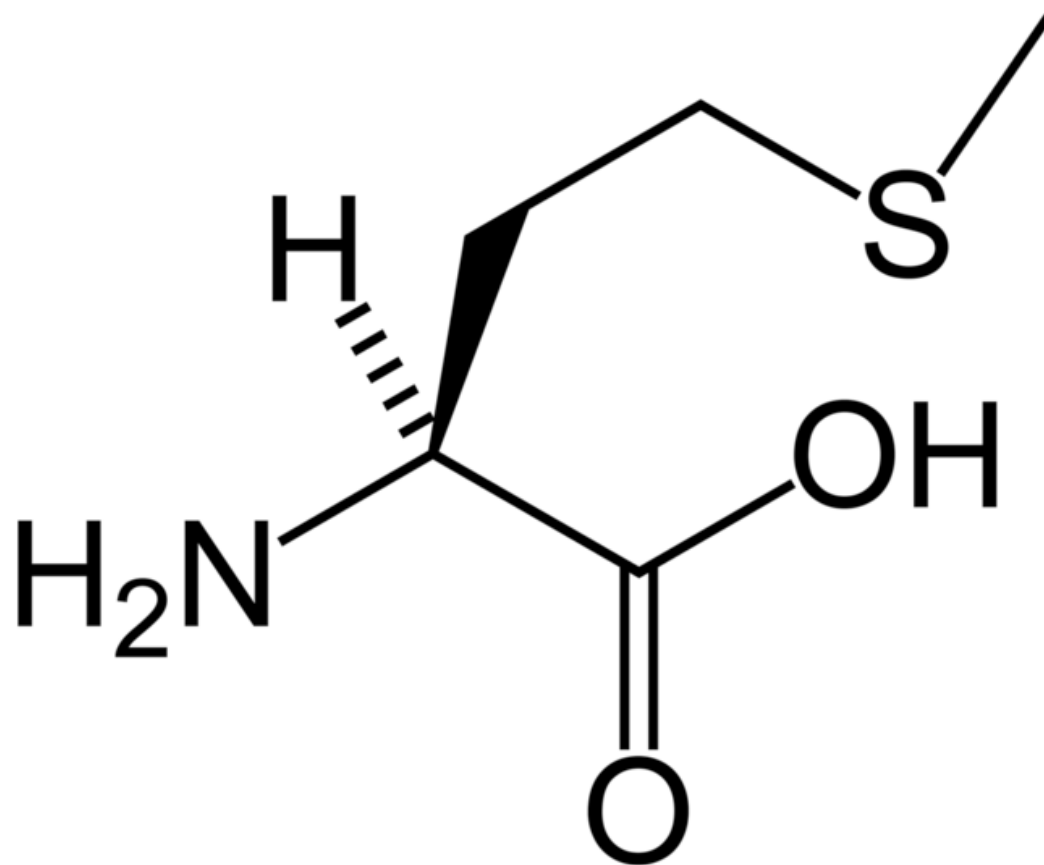
L-Isoleucine
(Ile / I)



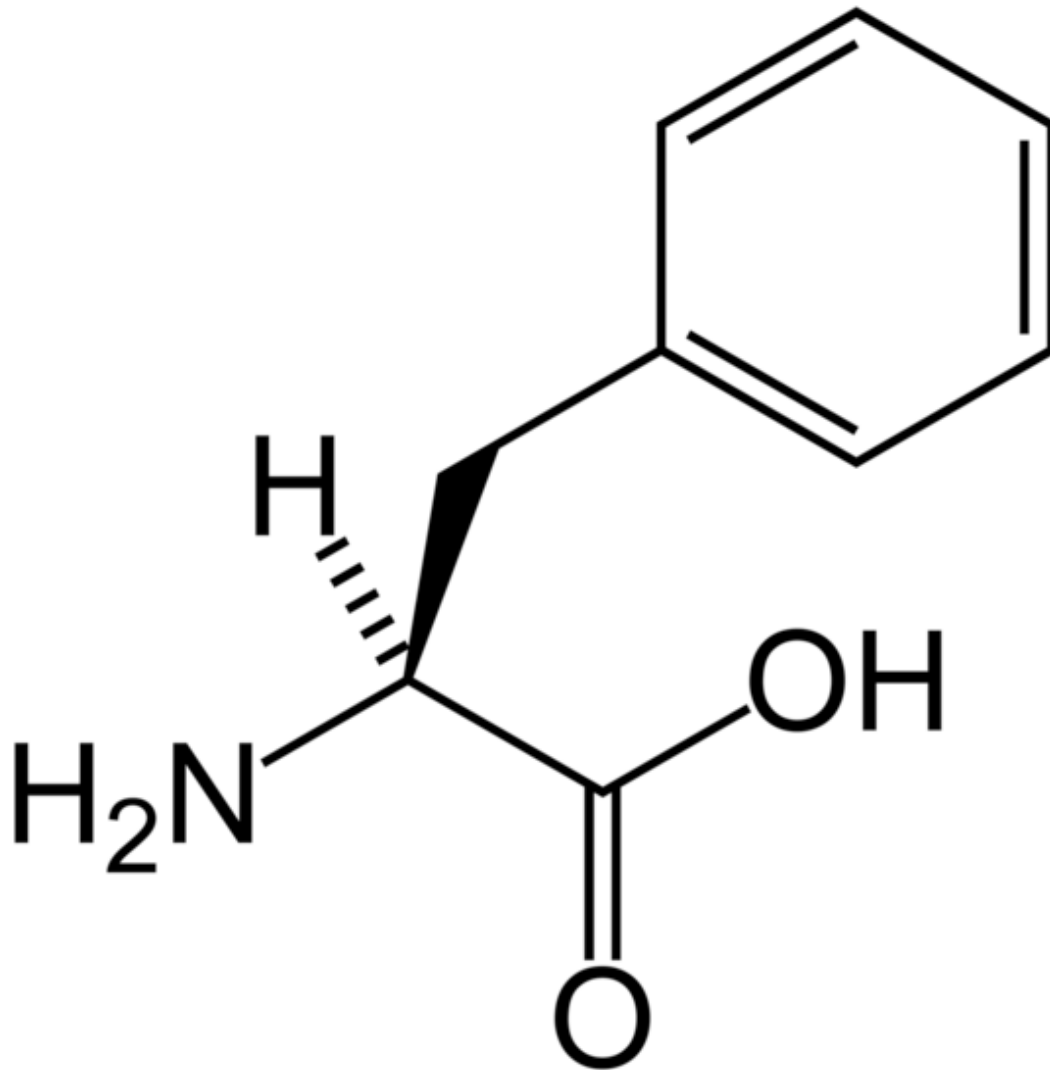
L-Leucine
(Leu / L)



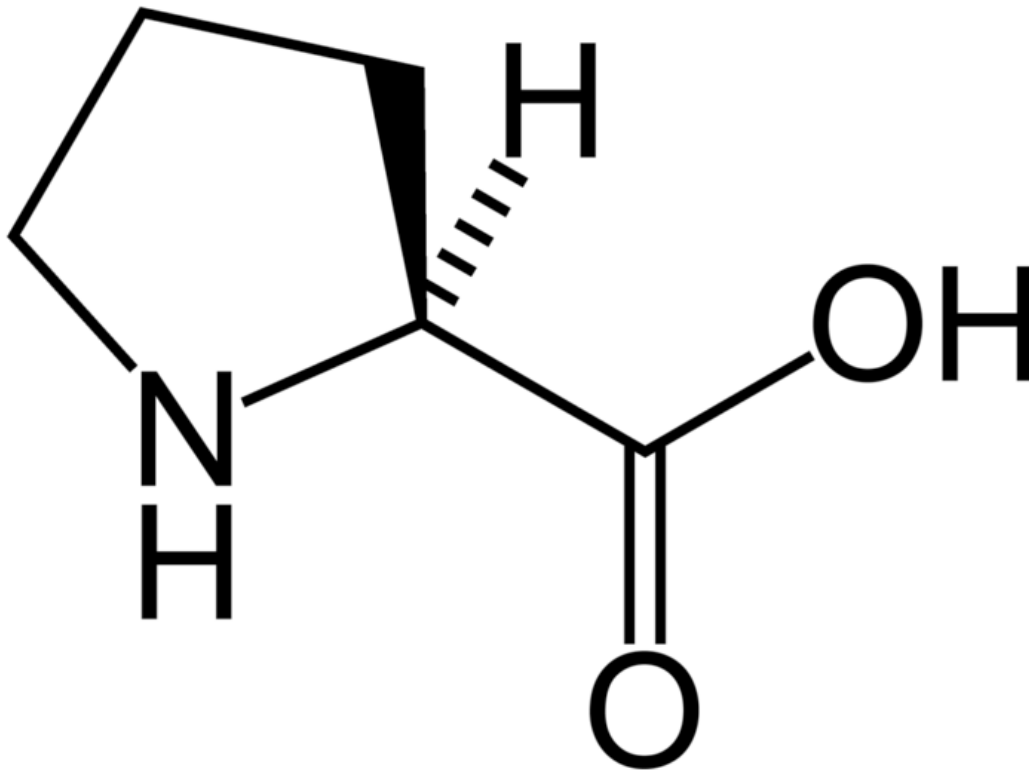
L-Lysine
(Lys / K)



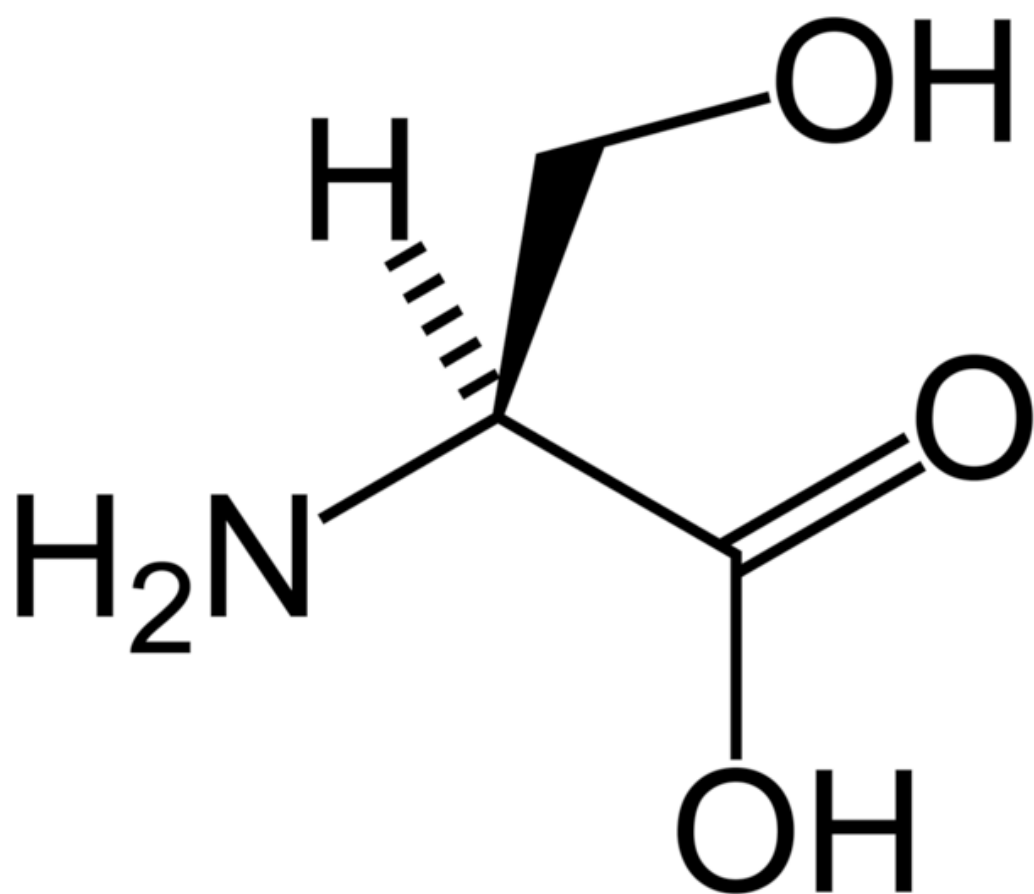
L-Methionine
(Met / M)



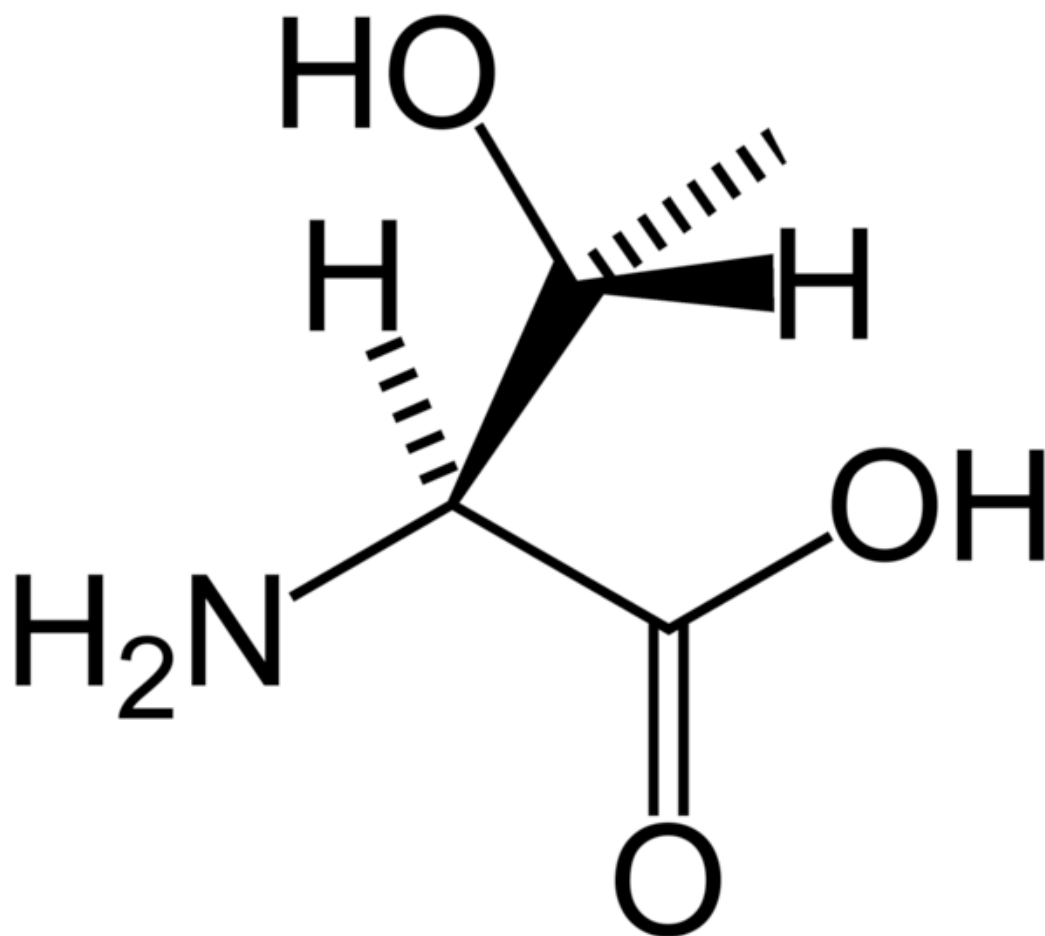
L-Phenylalanine
(Phe / F)



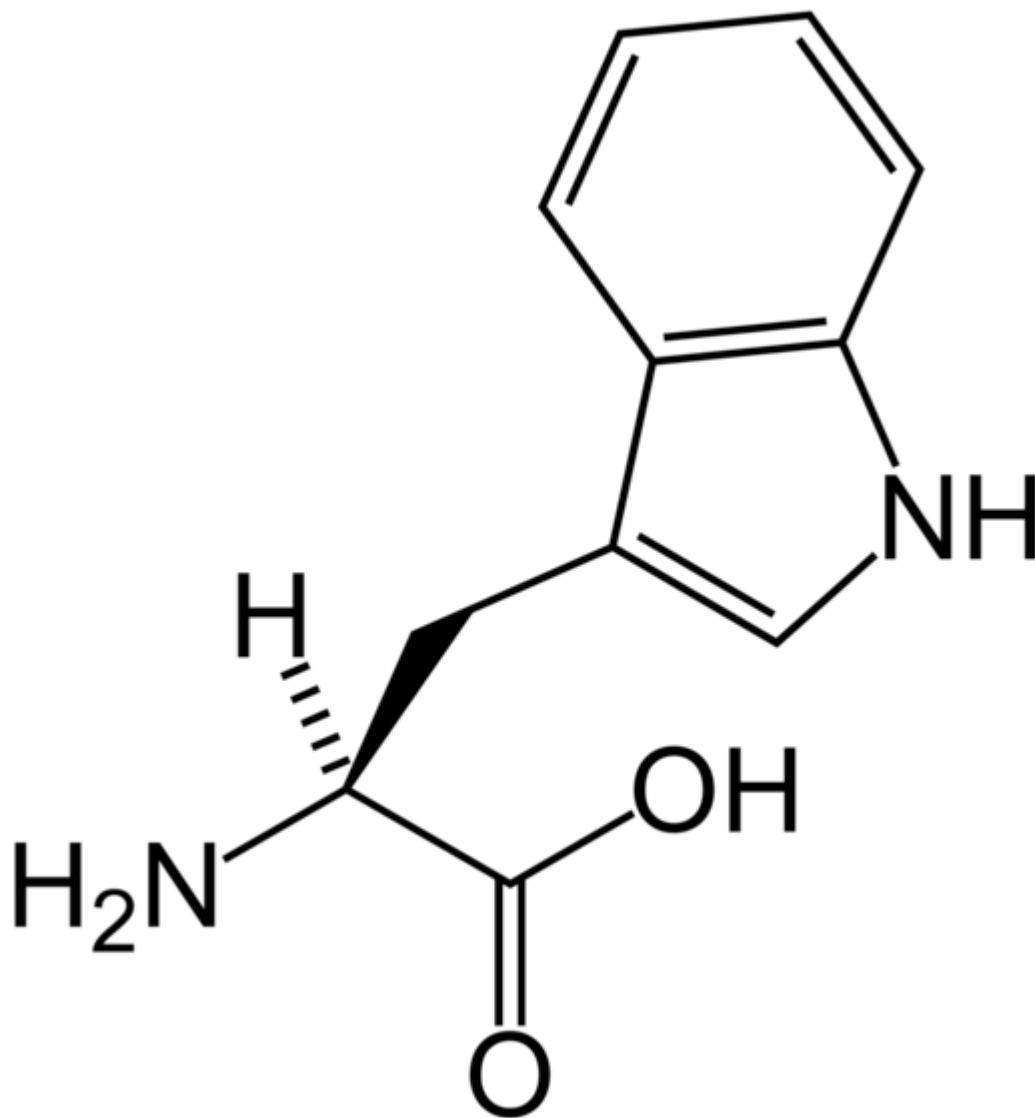
L-Proline
(Pro / P)



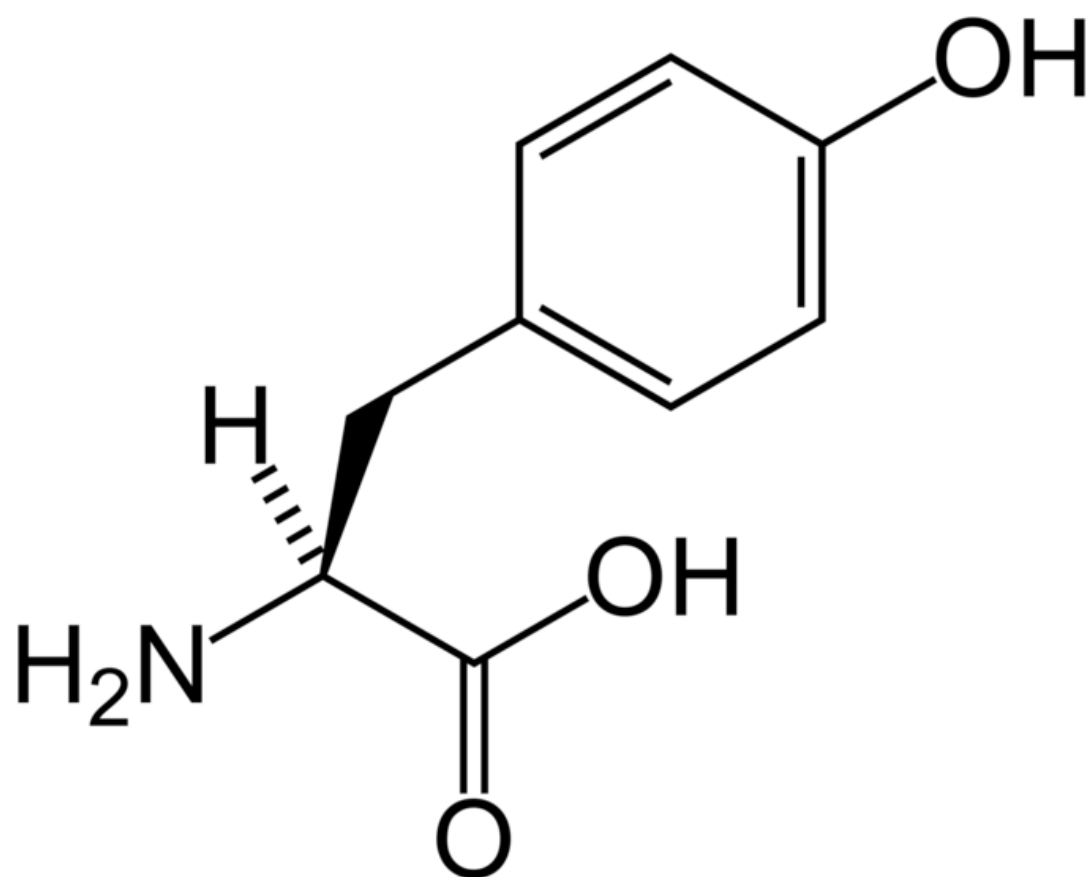
L-Serine
(Ser / S)



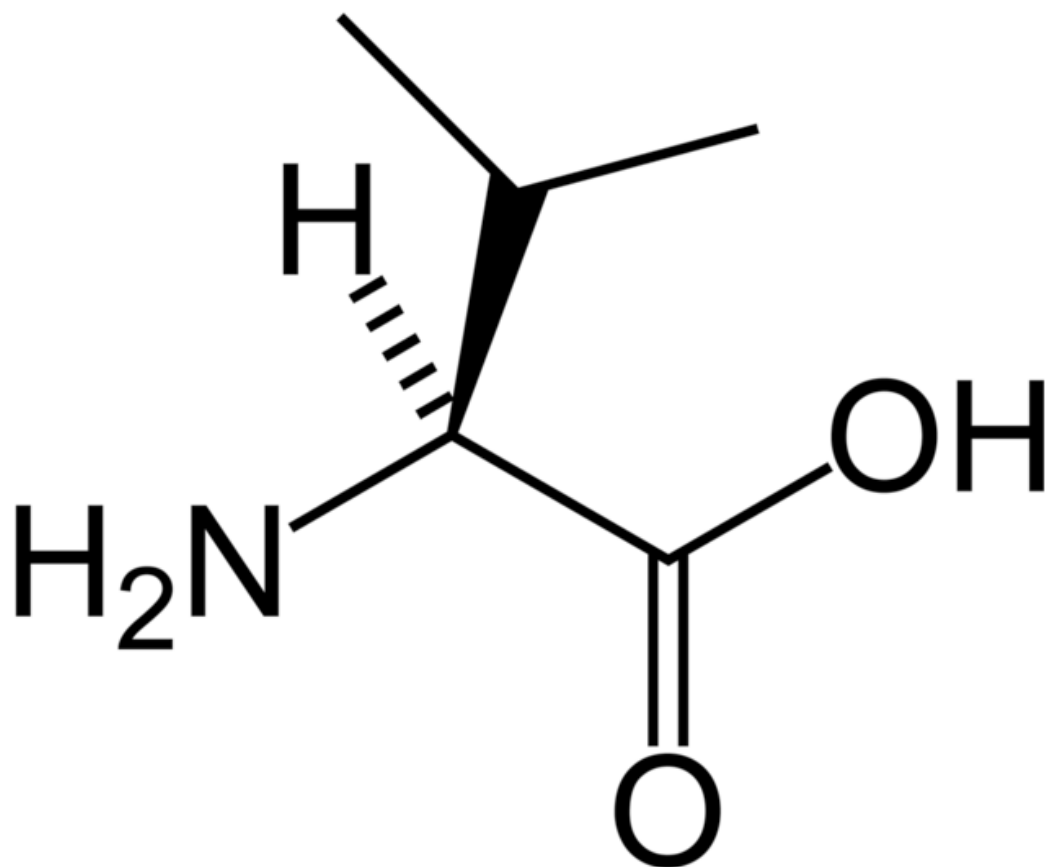
L-Threonine
(Thr / T)



L-Tryptophan
(Trp / W)

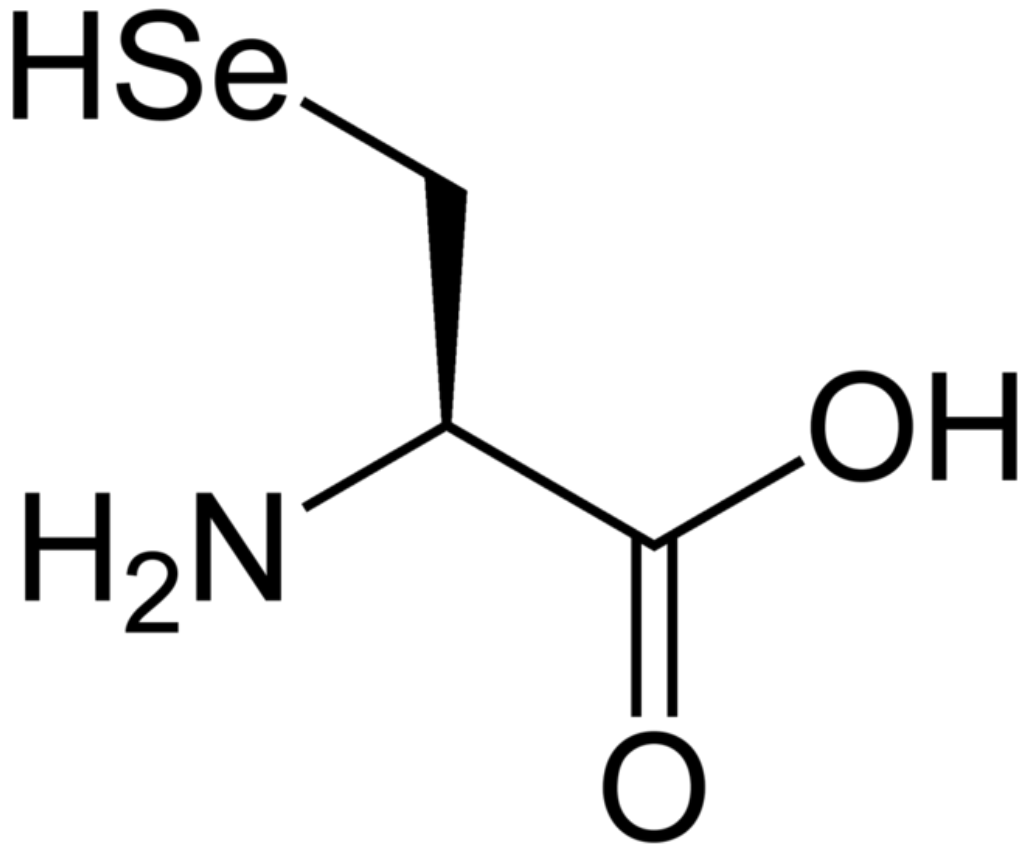


L-Tyrosine
(Tyr / Y)

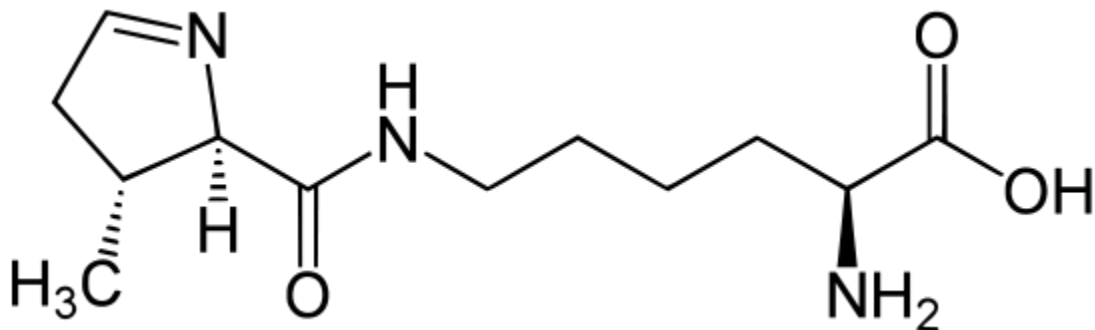


L-Valine
(Val / V)

IUPAC/IUBMB now also recommends standard abbreviations for the following two amino acids:



L-Selenocysteine



L-Pyrrolysine
(Pyl / O)

Non-specific abbreviations

Sometimes the specific identity of an amino acid cannot be determined unambiguously. Certain protein sequencing techniques do not distinguish among certain pairs. Thus, the following codes are used:

- *Asx* (*B*) is "asparagine or aspartic acid"
- *Glx* (*Z*) is "glutamic acid or glutamine"
- *Xle* (*J*) is "leucine or isoleucine"

In addition, the symbol *X* is used to indicate an amino acid that is completely unidentified.

Chemical properties

Following is a table listing the one-letter symbols, the three-letter symbols, and the chemical properties of the side-chains of the standard amino acids. The masses listed are based on weighted averages of the elemental isotopes at their natural abundances. Note that forming a peptide bond results in elimination of a molecule of water, so the mass of an amino acid unit within a protein chain is reduced by 18.01524 Da.

General chemical properties

Amino Acid	Short	Abbrev.	Avg. Mass (Da)	pI	pK ₁ (α -COOH)	pK ₂ (α -NH ₃ ⁺)
Alanine	A	Ala	89.09404	6.01	2.35	9.87
Cysteine	C	Cys	121.15404	5.05	1.92	10.70
Aspartic acid	D	Asp	133.10384	2.85	1.99	9.90
Glutamic acid	E	Glu	147.13074	3.15	2.10	9.47
Phenylalanine	F	Phe	165.19184	5.49	2.20	9.31
Glycine	G	Gly	75.06714	6.06	2.35	9.78
Histidine	H	His	155.15634	7.60	1.80	9.33
Isoleucine	I	Ile	131.17464	6.05	2.32	9.76
Lysine	K	Lys	146.18934	9.60	2.16	9.06
Leucine	L	Leu	131.17464	6.01	2.33	9.74
Methionine	M	Met	149.20784	5.74	2.13	9.28
Asparagine	N	Asn	132.11904	5.41	2.14	8.72
Pyrrolysine	O	Pyl				
Proline	P	Pro	115.13194	6.30	1.95	10.64
Glutamine	Q	Gln	146.14594	5.65	2.17	9.13
Arginine	R	Arg	174.20274	10.76	1.82	8.99
Serine	S	Ser	105.09344	5.68	2.19	9.21

Threonine	T	Thr	119.12034	5.60	2.09	9.10
Selenocysteine	U	Sec	168.053			
Valine	V	Val	117.14784	6.00	2.39	9.74
Tryptophan	W	Trp	204.22844	5.89	2.46	9.41
Tyrosine	Y	Tyr	181.19124	5.64	2.20	9.21

Side chain properties

Amino Acid	Short	Abbrev.	Side chain	Hydro- phobic	pKa	Polar	pH	Small	Tiny	Aromatic or Aliphatic	van der Waals volume
Alanine	A	Ala	-CH ₃	X	-	-	-	X	X	-	67
Cysteine	C	Cys	-CH ₂ SH	X	8.18	-	acidic	X	-	-	86
Aspartic acid	D	Asp	-CH ₂ COOH	-	3.90	X	acidic	X	-	-	91
Glutamic acid	E	Glu	-CH ₂ CH ₂ COOH	-	4.07	X	acidic	-	-	-	109
Phenylalanine	F	Phe	-CH ₂ C ₆ H ₅	X	-	-	-	-	-	Aromatic	135
Glycine	G	Gly	-H	X	-	-	-	X	X	-	48
Histidine	H	His	-CH ₂ -C ₃ H ₃ N ₂	-	6.04	X	weak basic	-	-	Aromatic	118
Isoleucine	I	Ile	-CH(CH ₃)CH ₂ CH ₃	X	-	-	-	-	-	Aliphatic	124
Lysine	K	Lys	-(CH ₂) ₄ NH ₂	-	10.54	X	basic	-	-	-	135
Leucine	L	Leu	-CH ₂ CH(CH ₃) ₂	X	-	-	-	-	-	Aliphatic	124
Methionine	M	Met	-CH ₂ CH ₂ SCH ₃	X	-	-	-	-	-	-	124
Asparagine	N	Asn	-CH ₂ CONH ₂	-	-	X	-	X	-	-	96
Pyrrolysine	O	Pyl									
Proline	P	Pro	-CH ₂ CH ₂ CH ₂ -	X	-	-	-	X	-	-	90
Glutamine	Q	Gln	-CH ₂ CH ₂ CONH ₂	-	-	X	-	-	-	-	114
Arginine	R	Arg	-(CH ₂) ₃ NH- C(NH)NH ₂	-	12.48	X	strongly basic	-	-	-	148
Serine	S	Ser	-CH ₂ OH	-	-	X	-	X	X	-	73
Threonine	T	Thr	-CH(OH)CH ₃	-	-	X	weak acidic	X	-	-	93
Selenocysteine	U	Sec	-CH ₂ SeH	X	5.73	-	-	X	-	-	
Valine	V	Val	-CH(CH ₃) ₂	X	-	-	-	X	-	Aliphatic	105
Tryptophan	W	Trp	-CH ₂ C ₈ H ₆ N	X	-	-	-	-	-	Aromatic	163
Tyrosine	Y	Tyr	-CH ₂ -C ₆ H ₄ OH	-	10.46	X	-	-	-	Aromatic	141

Note: The pKa values of amino acids are typically slightly different when the amino acid is inside a protein. Protein pKa calculations are sometimes used to calculate the change in the pKa value of an amino acid in this situation.

Gene expression and biochemistry

Amino Acid	Short Abbrev.	Codon(s)	Occurrence in human proteins (%)	Essential‡ in humans
Alanine	A Ala	GCU, GCC, GCA, GCG	7.8	-
Cysteine	C Cys	UGU, UGC	1.9	Conditionally
Aspartic acid	D Asp	GAU, GAC	5.3	-
Glutamic acid	E Glu	GAA, GAG	6.3	Conditionally
Phenylalanine	F Phe	UUU, UUC	3.9	Yes
Glycine	G Gly	GGU, GGC, GGA, GGG	7.2	Conditionally
Histidine	H His	CAU, CAC	2.3	Yes
Isoleucine	I Ile	AUU, AUC, AUA	5.3	Yes
Lysine	K Lys	AAA, AAG	5.9	Yes
Leucine	L Leu	UUA, UUG, CUU, CUC, CUA, CUG	9.1	Yes
Methionine	M Met	AUG	2.3	Yes
Asparagine	N Asn	AAU, AAC	4.3	-
Pyrrolysine	O Pyl	UAG*		-
Proline	P Pro	CCU, CCC, CCA, CCG	5.2	-
Glutamine	Q Gln	CAA, CAG	4.2	-
Arginine	R Arg	CGU, CGC, CGA, CGG, AGA, AGG	5.1	Conditionally
Serine	S Ser	UCU, UCC, UCA, UCG, AGU, AGC	6.8	-
Threonine	T Thr	ACU, ACC, ACA, ACG	5.9	Yes
Selenocysteine	U Sec	UGA**		-
Valine	V Val	GUU, GUC, GUA, GUG	6.6	Yes
Tryptophan	W Trp	UGG	1.4	Yes
Tyrosine	Y Tyr	UAU, UAC	3.2	Conditionally
Stop codon†	- Term	UAA, UAG, UGA	-	-

* UAG is normally the amber stop codon, but encodes pyrrolysine if a PYLIS element is present.

** UGA is normally the opal (or umber) stop codon, but encodes selenocysteine if a SECIS element is present.

† The stop codon is not an amino acid, but is included for completeness.

‡ An essential amino acid cannot be synthesized in humans and must, therefore, be supplied in the diet. Conditionally essential amino acids are not normally required in the diet, but must be supplied exogenously to specific populations that do not synthesize it in adequate amounts.

Mass spectrometry

In mass spectrometry of peptides and proteins, it is useful to know the masses of the residues. The mass of the peptide or protein is the sum of the residue masses plus the mass of water.

Amino Acid	Short	Abbrev.	Formula	Mon. Mass§ (Da)	Avg. Mass (Da)
Alanine	A	Ala	C ₃ H ₅ NO	71.03711	71.0788
Cysteine	C	Cys	C ₃ H ₅ NOS	103.00919	103.1388
Aspartic acid	D	Asp	C ₄ H ₅ NO ₃	115.02694	115.0886
Glutamic acid	E	Glu	C ₅ H ₇ NO ₃	129.04259	129.1155
Phenylalanine	F	Phe	C ₉ H ₉ NO	147.06841	147.1766
Glycine	G	Gly	C ₂ H ₃ NO	57.02146	57.0519
Histidine	H	His	C ₆ H ₇ N ₃ O	137.05891	137.1411
Isoleucine	I	Ile	C ₆ H ₁₁ NO	113.08406	113.1594
Lysine	K	Lys	C ₆ H ₁₂ N ₂ O	128.09496	128.1741
Leucine	L	Leu	C ₆ H ₁₁ NO	113.08406	113.1594
Methionine	M	Met	C ₅ H ₉ NOS	131.04049	131.1986
Asparagine	N	Asn	C ₄ H ₆ N ₂ O ₂	114.04293	114.1039
Pyrrolysine	O	Pyl	C ₁₂ H ₂₁ N ₃ O ₃	255.15829	255.3172
Proline	P	Pro	C ₅ H ₇ NO	97.05276	97.1167
Glutamine	Q	Gln	C ₅ H ₈ N ₂ O ₂	128.05858	128.1307
Arginine	R	Arg	C ₆ H ₁₂ N ₄ O	156.10111	156.1875
Serine	S	Ser	C ₃ H ₅ NO ₂	87.03203	87.0782
Threonine	T	Thr	C ₄ H ₇ NO ₂	101.04768	101.1051
Selenocysteine	U	Sec	C ₃ H ₅ NOSe	150.95364	150.0388
Valine	V	Val	C ₅ H ₉ NO	99.06841	99.1326
Tryptophan	W	Trp	C ₁₁ H ₁₀ N ₂ O	186.07931	186.2132
Tyrosine	Y	Tyr	C ₉ H ₉ NO ₂	163.06333	163.1760

§ Monoisotopic mass

Stoichiometry and metabolic cost in cell

Following table lists the abundance of amino acids in E.coli cell and the metabolic cost (ATP) for synthesis the amino acids. Negative numbers indicate the metabolic processes are energy favorable and do not cost net ATP of the cell. Note that the abundance of amino acids include amino acids in free-form and in polymerization form (proteins).

Amino acid	Abundance (# of molecules ($\times 10^8$) per <i>E. coli</i> cell)	ATP cost in synthesis under aerobic condition	ATP cost in synthesis under anaerobic condition
Alanine	2.9	-1	1
Cysteine	0.52	11	15
Aspartic acid	1.4	0	2
Glutamic acid	1.5	-7	-1
Phenylalanine	1.1	-6	2
Glycine	3.5	-2	2
Histidine	0.54	1	7
Isoleucine	1.7	7	11
Lysine	2.0	5	9
Leucine	2.6	-9	1
Methionine	0.88	21	23
Asparagine	1.4	3	5
Proline	1.3	-2	4
Glutamine	1.5	-6	0
Arginine	1.7	5	13
Serine	1.2	-2	2
Threonine	1.5	6	8
Tryptophan	0.33	-7	7
Tyrosine	0.79	-8	2
Valine	2.4	-2	2

Remarks

Amino Acid	Abbrev.	Remarks
Alanine	A Ala	Very abundant, very versatile. More stiff than glycine, but small enough to pose only small steric limits for the protein conformation. It behaves fairly neutrally, and can be located in both hydrophilic regions on the protein outside and the

			hydrophobic areas inside.
Asparagine or aspartic acid	B	Asx	A placeholder when either amino acid may occupy a position.
Cysteine	C	Cys	The sulfur atom bonds readily to heavy metal ions. Under oxidizing conditions, two cysteines can join together in a disulfide bond to form the amino acid cystine. When cystines are part of a protein, insulin for example, the tertiary structure is stabilized, which makes the protein more resistant to denaturation; therefore, disulfide bonds are common in proteins that have to function in harsh environments including digestive enzymes (e.g., pepsin and chymotrypsin) and structural proteins (e.g., keratin). Disulfides are also found in peptides too small to hold a stable shape on their own (eg. insulin).
Aspartic acid	D	Asp	Behaves similarly to glutamic acid. Carries a hydrophilic acidic group with strong negative charge. Usually is located on the outer surface of the protein, making it water-soluble. Binds to positively-charged molecules and ions, often used in enzymes to fix the metal ion. When located inside of the protein, aspartate and glutamate are usually paired with arginine and lysine.
Glutamic acid	E	Glu	Behaves similar to aspartic acid. Has longer, slightly more flexible side chain.
Phenylalanine	F	Phe	Essential for humans. Phenylalanine, tyrosine, and tryptophan contain large rigid aromatic group on the side-chain. These are the biggest amino acids. Like isoleucine, leucine and valine, these are hydrophobic and tend to orient towards the interior of the folded protein molecule. Phenylalanine can be converted into Tyrosine.
Glycine	G	Gly	Because of the two hydrogen atoms at the α carbon, glycine is not optically active. It is the smallest amino acid, rotates easily, adds flexibility to the protein chain. It is able to fit into the tightest spaces, e.g., the triple helix of collagen. As too much flexibility is usually not desired, as a structural component it is less common than alanine.
Histidine	H	His	In even slightly acidic conditions protonation of the nitrogen occurs, changing the properties of histidine and the polypeptide as a whole. It is used by many proteins as a regulatory mechanism, changing the conformation and behavior of the polypeptide in acidic regions such as the late endosome or lysosome, enforcing conformation change in enzymes. However only a few histidines are needed for this, so it is comparatively scarce.
Isoleucine	I	Ile	Essential for humans. Isoleucine, leucine and valine have large aliphatic hydrophobic side chains. Their molecules are

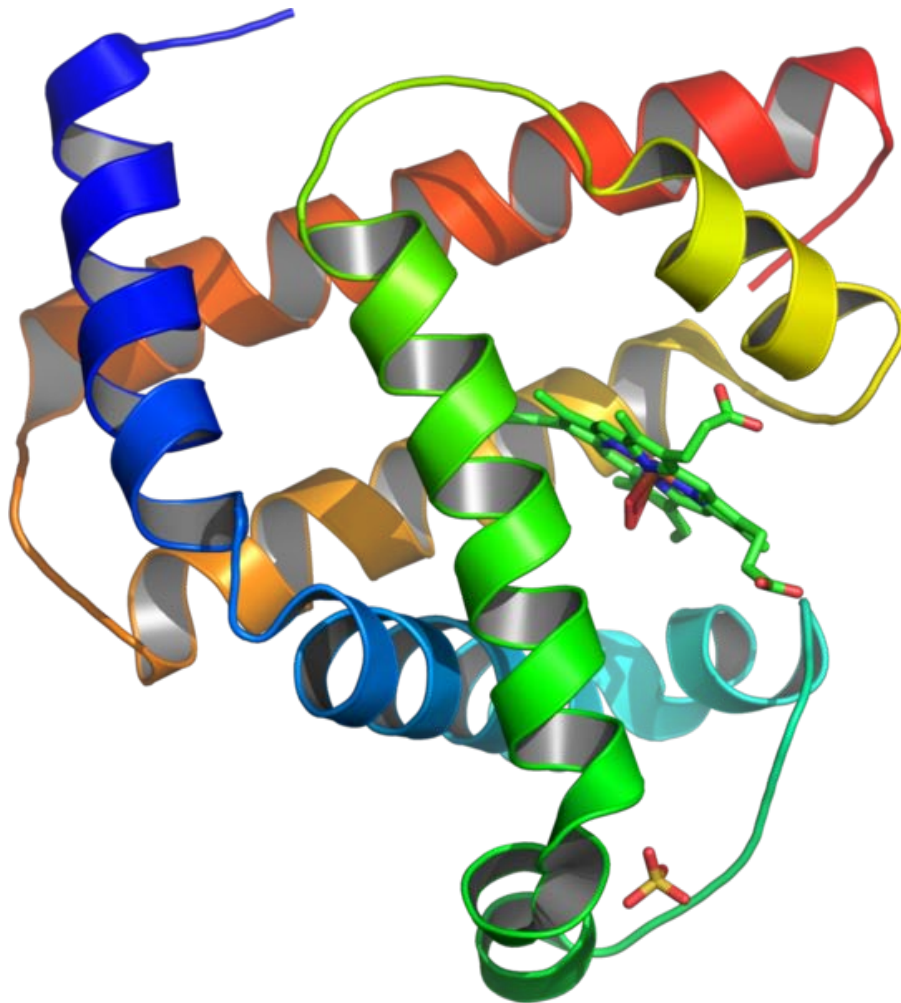
rigid, and their mutual hydrophobic interactions are important for the correct folding of proteins, as these chains tend to be located inside of the protein molecule.

Leucine or isoleucine	J	Xle	A placeholder when either amino acid may occupy a position
Lysine	K	Lys	Essential for humans. Behaves similarly to arginine. Contains a long flexible side-chain with a positively-charged end. The flexibility of the chain makes lysine and arginine suitable for binding to molecules with many negative charges on their surfaces. E.g., DNA-binding proteins have their active regions rich with arginine and lysine. The strong charge makes these two amino acids prone to be located on the outer hydrophilic surfaces of the proteins; when they are found inside, they are usually paired with a corresponding negatively-charged amino acid, e.g., aspartate or glutamate.
Leucine	L	Leu	Essential for humans. Behaves similar to isoleucine and valine.
Methionine	M	Met	Essential for humans. Always the first amino acid to be incorporated into a protein; sometimes removed after translation. Like cysteine, contains sulfur, but with a methyl group instead of hydrogen. This methyl group can be activated, and is used in many reactions where a new carbon atom is being added to another molecule.
Asparagine	N	Asn	Similar to aspartic acid. Asn contains an amide group where Asp has a carboxyl.
Pyrrolysine	O	Pyl	Similar to lysine, with a pyrroline ring attached.
Proline	P	Pro	Contains an unusual ring to the N-end amine group, which forces the CO-NH amide sequence into a fixed conformation. Can disrupt protein folding structures like α helix or β sheet, forcing the desired kink in the protein chain. Common in collagen, where it often undergoes a posttranslational modification to hydroxyproline.
Glutamine	Q	Gln	Similar to glutamic acid. Gln contains an amide group where Glu has a carboxyl. Used in proteins and as a storage for ammonia. The most abundant Amino Acid in the body.
Arginine	R	Arg	Functionally similar to lysine.
Serine	S	Ser	Serine and threonine have a short group ended with a hydroxyl group. Its hydrogen is easy to remove, so serine and threonine often act as hydrogen donors in enzymes. Both are very hydrophilic, therefore the outer regions of soluble proteins tend to be rich with them.
Threonine	T	Thr	Essential for humans. Behaves similarly to serine.
Selenocysteine	U	Sec	Selenated form of cysteine, which replaces sulfur.

Valine	V	Val	Essential for humans. Behaves similarly to isoleucine and leucine.
Tryptophan	W	Trp	Essential for humans. Behaves similarly to phenylalanine and tyrosine. Precursor of serotonin. Naturally fluorescent.
Unknown	X	Xaa	Placeholder when the amino acid is unknown or unimportant.
Tyrosine	Y	Tyr	Behaves similarly to phenylalanine (precursor to Tyrosine) and tryptophan. Precursor of melanin, epinephrine, and thyroid hormones. Naturally fluorescent, although fluorescence is usually quenched by energy transfer to tryptophans.
Glutamic acid or glutamine	Z	Glx	A placeholder when either amino acid may occupy a position.

Chapter 15

Protein



A representation of the 3D structure of the protein myoglobin showing colored alpha helices. This protein was the first to have its structure solved by X-ray crystallography. Towards the right-center among the coils, a prosthetic group called a heme group is shown colored largely in green.

Proteins are biochemical compounds consisting of one or more polypeptides typically folded into a globular or fibrous form in a biologically functional way. A polypeptide is a single linear polymer chain of amino acids bonded together by peptide bonds between the carboxyl and amino groups of adjacent amino acid residues. The sequence of amino acids in a protein is defined by the sequence of a gene, which is encoded in the genetic code. In general, the genetic code specifies 20 standard amino acids; however, in certain organisms the genetic code can include selenocysteine—and in certain archaea—pyrrolysine. Shortly after or even during synthesis, the residues in a protein are often chemically modified by post-translational modification, which alters the physical and chemical properties, folding, stability, activity, and ultimately, the function of the proteins. Sometimes proteins have non-peptide groups attached, which can be called prosthetic groups or cofactors. Proteins can also work together to achieve a particular function, and they often associate to form stable complexes.

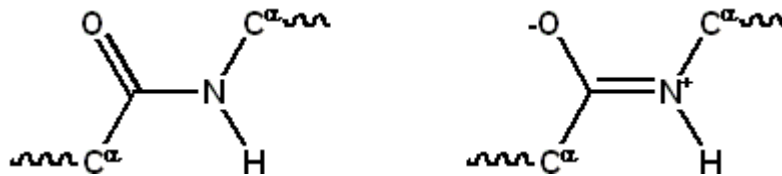
One of the most distinguishing features of polypeptides is their ability to fold into a globular state, or "structure". The extent to which proteins fold into a defined structure varies widely. Some proteins fold into a highly rigid structure with small fluctuations and are therefore considered to be single structure. Other proteins undergo large rearrangements from one conformation to another. This conformational change is often associated with a signaling event. Thus, the structure of a protein serves as a medium through which to regulate either the function of a protein or activity of an enzyme. Not all proteins requiring a folding process in order to function, as some function in an unfolded state.

Like other biological macromolecules such as polysaccharides and nucleic acids, proteins are essential parts of organisms and participate in virtually every process within cells. Many proteins are enzymes that catalyze biochemical reactions and are vital to metabolism. Proteins also have structural or mechanical functions, such as actin and myosin in muscle and the proteins in the cytoskeleton, which form a system of scaffolding that maintains cell shape. Other proteins are important in cell signaling, immune responses, cell adhesion, and the cell cycle. Proteins are also necessary in animals' diets, since animals cannot synthesize all the amino acids they need and must obtain essential amino acids from food. Through the process of digestion, animals break down ingested protein into free amino acids that are then used in metabolism.

Proteins were first described by the Dutch chemist Gerhardus Johannes Mulder and named by the Swedish chemist Jöns Jakob Berzelius in 1838. Early nutritional scientists such as the German Carl von Voit believed that protein was the most important nutrient for maintaining the structure of the body, because it was generally believed that "flesh makes flesh." The central role of proteins as enzymes in living organisms was however not fully appreciated until 1926, when James B. Sumner showed that the enzyme urease was in fact a protein. The first protein to be sequenced was insulin, by Frederick Sanger, who won the Nobel Prize for this achievement in 1958. The first protein structures to be solved were hemoglobin and myoglobin, by Max Perutz and Sir John Cowdery Kendrew, respectively, in 1958. The three-dimensional structures of both proteins were first determined by X-ray diffraction analysis; Perutz and Kendrew shared the 1962 Nobel

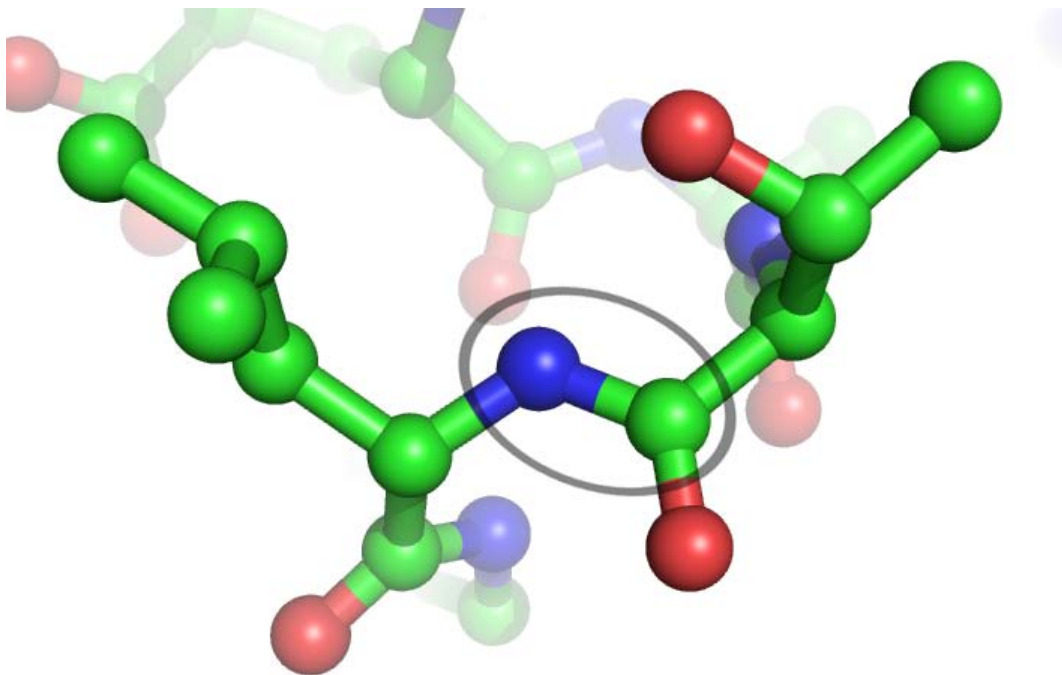
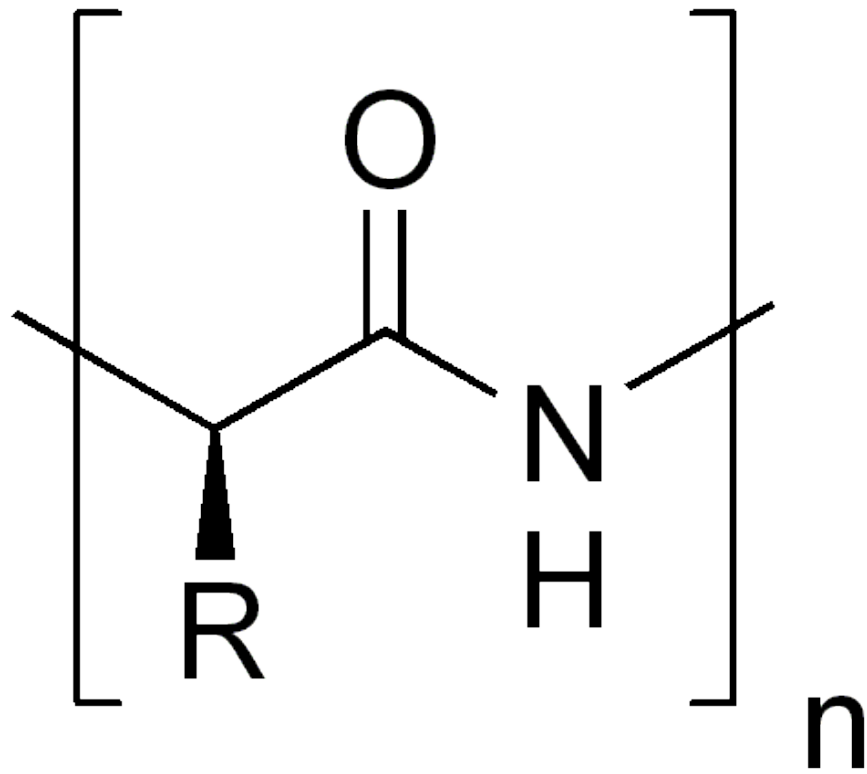
Prize in Chemistry for these discoveries. Proteins may be purified from other cellular components using a variety of techniques such as ultracentrifugation, precipitation, electrophoresis, and chromatography; the advent of genetic engineering has made possible a number of methods to facilitate purification. Methods commonly used to study protein structure and function include immunohistochemistry, site-directed mutagenesis, nuclear magnetic resonance and mass spectrometry. Distributed computing is a relatively new tool researchers are using to examine the infamously complex interactions that govern protein folding; the statistical analysis techniques employed to calculate a protein's probable tertiary structure from its amino acid sequence (primary structure) are well-suited for the distributed computing environment, which has made this otherwise prohibitively expensive and time consuming problem significantly more manageable.

Biochemistry



Resonance structures of the peptide bond that links individual amino acids to form a protein polymer

Most proteins consist of linear polymers built from series of up to 20 different L- α -amino acids. All proteinogenic amino acids possess common structural features, including an α -carbon to which an amino group, a carboxyl group, and a variable side chain are bonded. Only proline differs from this basic structure as it contains an unusual ring to the N-end amine group, which forces the CO–NH amide moiety into a fixed conformation. The side chains of the standard amino acids, detailed in the list of standard amino acids, have a great variety of chemical structures and properties; it is the combined effect of all of the amino acid side chains in a protein that ultimately determines its three-dimensional structure and its chemical reactivity.

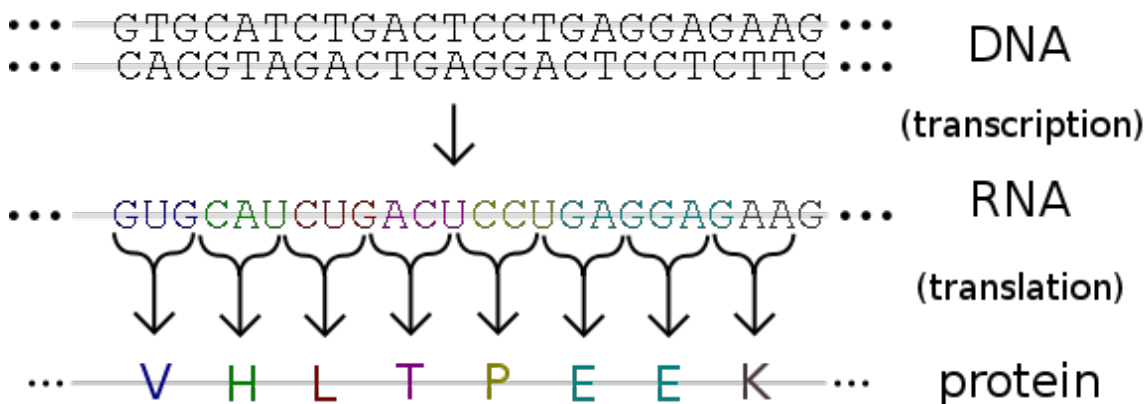


Chemical structure of the peptide bond (above) and a peptide bond between leucine and threonine (below)

The amino acids in a polypeptide chain are linked by peptide bonds. Once linked in the protein chain, an individual amino acid is called a *residue*, and the linked series of carbon, nitrogen, and oxygen atoms are known as the *main chain* or *protein backbone*. The peptide bond has two resonance forms that contribute some double-bond character and inhibit rotation around its axis, so that the alpha carbons are roughly coplanar. The other two dihedral angles in the peptide bond determine the local shape assumed by the protein backbone. The end of the protein with a free carboxyl group is known as the C-terminus or carboxy terminus, whereas the end with a free amino group is known as the N-terminus or amino terminus.

The words *protein*, *polypeptide*, and *peptide* are a little ambiguous and can overlap in meaning. *Protein* is generally used to refer to the complete biological molecule in a stable conformation, whereas *peptide* is generally reserved for a short amino acid oligomers often lacking a stable three-dimensional structure. However, the boundary between the two is not well defined and usually lies near 20–30 residues. *Polypeptide* can refer to any single linear chain of amino acids, usually regardless of length, but often implies an absence of a defined conformation.

Synthesis



The DNA sequence of a gene encodes the amino acid sequence of a protein.

Proteins are assembled from amino acids using information encoded in genes. Each protein has its own unique amino acid sequence that is specified by the nucleotide sequence of the gene encoding this protein. The genetic code is a set of three-nucleotide sets called codons and each three-nucleotide combination designates an amino acid, for example AUG (adenine-uracil-guanine) is the code for methionine. Because DNA contains four nucleotides, the total number of possible codons is 64; hence, there is some redundancy in the genetic code, with some amino acids specified by more than one codon. Genes encoded in DNA are first transcribed into pre-messenger RNA (mRNA) by proteins such as RNA polymerase. Most organisms then process the pre-mRNA (also known as a *primary transcript*) using various forms of post-transcriptional modification to form the mature mRNA, which is then used as a template for protein synthesis by the ribosome. In prokaryotes the mRNA may either be used as soon as it is produced, or be

bound by a ribosome after having moved away from the nucleoid. In contrast, eukaryotes make mRNA in the cell nucleus and then translocate it across the nuclear membrane into the cytoplasm, where protein synthesis then takes place. The rate of protein synthesis is higher in prokaryotes than eukaryotes and can reach up to 20 amino acids per second.

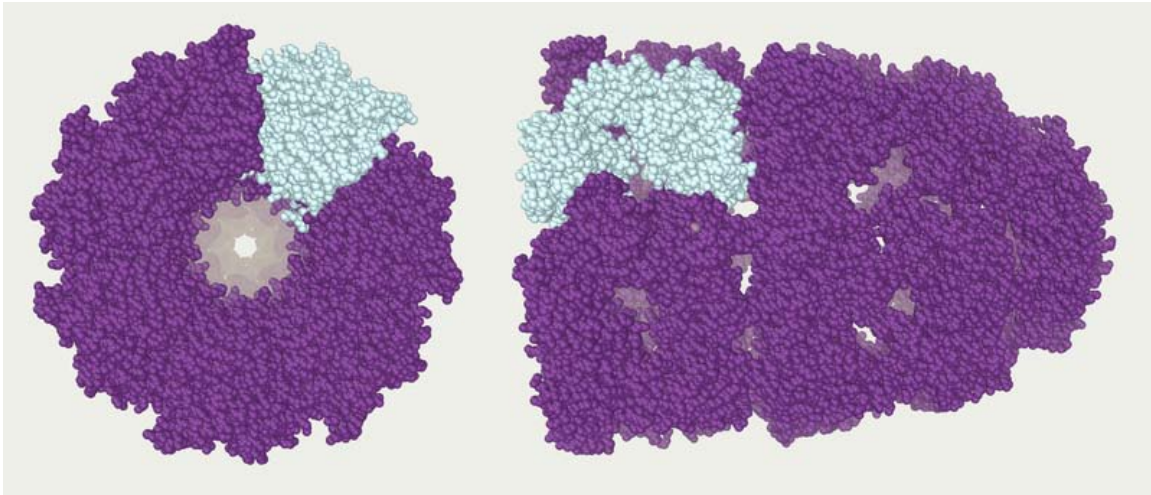
The process of synthesizing a protein from an mRNA template is known as translation. The mRNA is loaded onto the ribosome and is read three nucleotides at a time by matching each codon to its base pairing anticodon located on a transfer RNA molecule, which carries the amino acid corresponding to the codon it recognizes. The enzyme aminoacyl tRNA synthetase "charges" the tRNA molecules with the correct amino acids. The growing polypeptide is often termed the *nascent chain*. Proteins are always biosynthesized from N-terminus to C-terminus.

The size of a synthesized protein can be measured by the number of amino acids it contains and by its total molecular mass, which is normally reported in units of *daltons* (synonymous with atomic mass units), or the derivative unit kilodalton (kDa). Yeast proteins are on average 466 amino acids long and 53 kDa in mass. The largest known proteins are the titins, a component of the muscle sarcomere, with a molecular mass of almost 3,000 kDa and a total length of almost 27,000 amino acids.

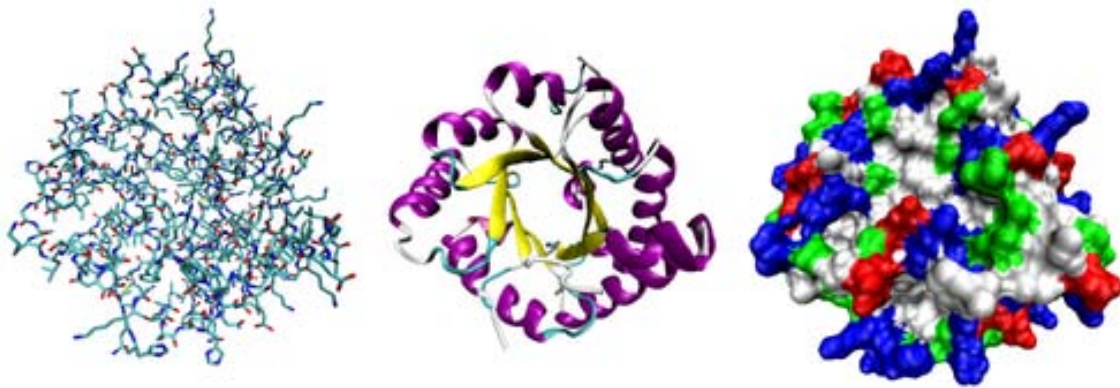
Chemical synthesis

Short proteins can also be synthesized chemically by a family of methods known as peptide synthesis, which rely on organic synthesis techniques such as chemical ligation to produce peptides in high yield. Chemical synthesis allows for the introduction of non-natural amino acids into polypeptide chains, such as attachment of fluorescent probes to amino acid side chains. These methods are useful in laboratory biochemistry and cell biology, though generally not for commercial applications. Chemical synthesis is inefficient for polypeptides longer than about 300 amino acids, and the synthesized proteins may not readily assume their native tertiary structure. Most chemical synthesis methods proceed from C-terminus to N-terminus, opposite the biological reaction.

Structure



The crystal structure of the chaperonin. Chaperonins assist protein folding.



Three possible representations of the three-dimensional structure of the protein triose phosphate isomerase. Left: all-atom representation colored by atom type. Middle: Simplified representation illustrating the backbone conformation, colored by secondary structure. Right: Solvent-accessible surface representation colored by residue type (acidic residues red, basic residues blue, polar residues green, nonpolar residues white)

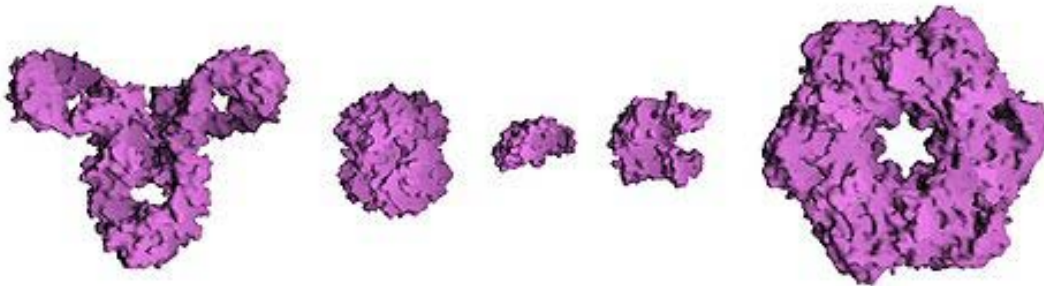
Most proteins fold into unique 3-dimensional structures. The shape into which a protein naturally folds is known as its native conformation. Although many proteins can fold unassisted, simply through the chemical properties of their amino acids, others require the aid of molecular chaperones to fold into their native states. Biochemists often refer to four distinct aspects of a protein's structure:

- *Primary structure*: the amino acid sequence.
- *Secondary structure*: regularly repeating local structures stabilized by hydrogen bonds. The most common examples are the alpha helix, beta sheet and turns.

Because secondary structures are local, many regions of different secondary structure can be present in the same protein molecule.

- *Tertiary structure*: the overall shape of a single protein molecule; the spatial relationship of the secondary structures to one another. Tertiary structure is generally stabilized by nonlocal interactions, most commonly the formation of a hydrophobic core, but also through salt bridges, hydrogen bonds, disulfide bonds, and even post-translational modifications. The term "tertiary structure" is often used as synonymous with the term *fold*. The tertiary structure is what controls the basic function of the protein.
- *Quaternary structure*: the structure formed by several protein molecules (polypeptide chains), usually called *protein subunits* in this context, which function as a single protein complex.

Proteins are not entirely rigid molecules. In addition to these levels of structure, proteins may shift between several related structures while they perform their functions. In the context of these functional rearrangements, these tertiary or quaternary structures are usually referred to as "conformations", and transitions between them are called *conformational changes*. Such changes are often induced by the binding of a substrate molecule to an enzyme's active site, or the physical region of the protein that participates in chemical catalysis. In solution proteins also undergo variation in structure through thermal vibration and the collision with other molecules.



Molecular surface of several proteins showing their comparative sizes. From left to right are: immunoglobulin G (IgG, an antibody), hemoglobin, insulin (a hormone), adenylate kinase (an enzyme), and glutamine synthetase (an enzyme).

Proteins can be informally divided into three main classes, which correlate with typical tertiary structures: globular proteins, fibrous proteins, and membrane proteins. Almost all globular proteins are soluble and many are enzymes. Fibrous proteins are often structural, such as collagen, the major component of connective tissue, or keratin, the protein component of hair and nails. Membrane proteins often serve as receptors or provide channels for polar or charged molecules to pass through the cell membrane.

A special case of intramolecular hydrogen bonds within proteins, poorly shielded from water attack and hence promoting their own dehydration, are called dehydrons.

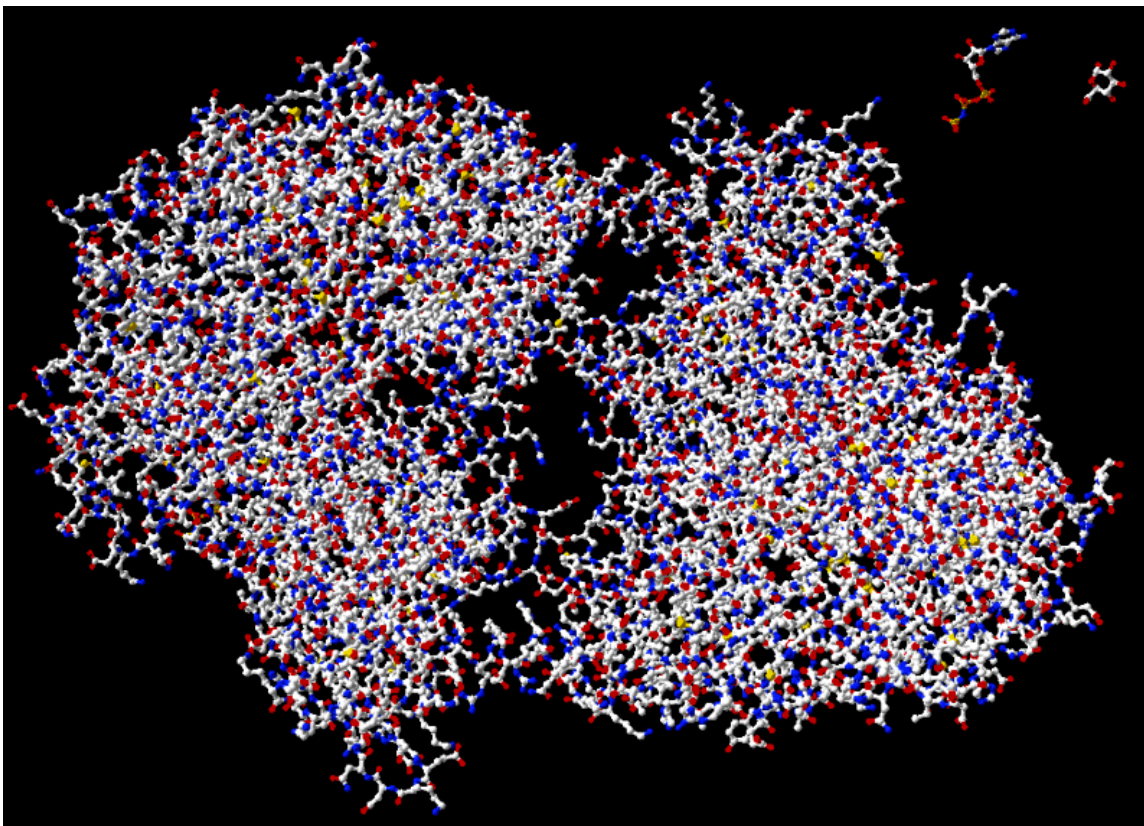
Structure determination

Discovering the tertiary structure of a protein, or the quaternary structure of its complexes, can provide important clues about how the protein performs its function. Common experimental methods of structure determination include X-ray crystallography and NMR spectroscopy, both of which can produce information at atomic resolution. However, NMR experiments are able to provide information from which a subset of distances between pairs of atoms can be estimated, and the final possible conformations for a protein are determined by solving a distance geometry problem. Dual polarisation interferometry is a quantitative analytical method for measuring the overall protein conformation and conformational changes due to interactions or other stimulus. Circular dichroism is another laboratory technique for determining internal beta sheet/ helical composition of proteins. Cryoelectron microscopy is used to produce lower-resolution structural information about very large protein complexes, including assembled viruses; a variant known as electron crystallography can also produce high-resolution information in some cases, especially for two-dimensional crystals of membrane proteins. Solved structures are usually deposited in the Protein Data Bank (PDB), a freely available resource from which structural data about thousands of proteins can be obtained in the form of Cartesian coordinates for each atom in the protein.

Many more gene sequences are known than protein structures. Further, the set of solved structures is biased toward proteins that can be easily subjected to the conditions required in X-ray crystallography, one of the major structure determination methods. In particular, globular proteins are comparatively easy to crystallize in preparation for X-ray crystallography. Membrane proteins, by contrast, are difficult to crystallize and are underrepresented in the PDB. Structural genomics initiatives have attempted to remedy these deficiencies by systematically solving representative structures of major fold classes. Protein structure prediction methods attempt to provide a means of generating a plausible structure for proteins whose structures have not been experimentally determined.

Cellular functions

Proteins are the chief actors within the cell, said to be carrying out the duties specified by the information encoded in genes. With the exception of certain types of RNA, most other biological molecules are relatively inert elements upon which proteins act. Proteins make up half the dry weight of an *Escherichia coli* cell, whereas other macromolecules such as DNA and RNA make up only 3% and 20%, respectively. The set of proteins expressed in a particular cell or cell type is known as its proteome.



The enzyme hexokinase is shown as a simple ball-and-stick molecular model. To scale in the top right-hand corner are two of its substrates, ATP and glucose.

The chief characteristic of proteins that also allows their diverse set of functions is their ability to bind other molecules specifically and tightly. The region of the protein responsible for binding another molecule is known as the binding site and is often a depression or "pocket" on the molecular surface. This binding ability is mediated by the tertiary structure of the protein, which defines the binding site pocket, and by the chemical properties of the surrounding amino acids' side chains. Protein binding can be extraordinarily tight and specific; for example, the ribonuclease inhibitor protein binds to human angiogenin with a sub-femtomolar dissociation constant ($<10^{-15}$ M) but does not bind at all to its amphibian homolog onconase (>1 M). Extremely minor chemical changes such as the addition of a single methyl group to a binding partner can sometimes suffice to nearly eliminate binding; for example, the aminoacyl tRNA synthetase specific to the amino acid valine discriminates against the very similar side chain of the amino acid isoleucine.

Proteins can bind to other proteins as well as to small-molecule substrates. When proteins bind specifically to other copies of the same molecule, they can oligomerize to form fibrils; this process occurs often in structural proteins that consist of globular monomers that self-associate to form rigid fibers. Protein-protein interactions also regulate enzymatic activity, control progression through the cell cycle, and allow the assembly of large protein complexes that carry out many closely related reactions with a common

biological function. Proteins can also bind to, or even be integrated into, cell membranes. The ability of binding partners to induce conformational changes in proteins allows the construction of enormously complex signaling networks. Importantly, as interactions between proteins are reversible, and depend heavily on the availability of different groups of partner proteins to form aggregates that are capable to carry out discrete sets of function, study of the interactions between specific proteins is a key to understand important aspects of cellular function, and ultimately the properties that distinguish particular cell types.

Enzymes

The best-known role of proteins in the cell is as enzymes, which catalyze chemical reactions. Enzymes are usually highly specific and accelerate only one or a few chemical reactions. Enzymes carry out most of the reactions involved in metabolism, as well as manipulating DNA in processes such as DNA replication, DNA repair, and transcription. Some enzymes act on other proteins to add or remove chemical groups in a process known as post-translational modification. About 4,000 reactions are known to be catalyzed by enzymes. The rate acceleration conferred by enzymatic catalysis is often enormous—as much as 10^{17} -fold increase in rate over the uncatalyzed reaction in the case of orotate decarboxylase (78 million years without the enzyme, 18 milliseconds with the enzyme).

The molecules bound and acted upon by enzymes are called substrates. Although enzymes can consist of hundreds of amino acids, it is usually only a small fraction of the residues that come in contact with the substrate, and an even smaller fraction—three to four residues on average—that are directly involved in catalysis. The region of the enzyme that binds the substrate and contains the catalytic residues is known as the active site.

Cell signaling and ligand binding



Ribbon diagram of a mouse antibody against cholera that binds a carbohydrate antigen

Many proteins are involved in the process of cell signaling and signal transduction. Some proteins, such as insulin, are extracellular proteins that transmit a signal from the cell in which they were synthesized to other cells in distant tissues. Others are membrane proteins that act as receptors whose main function is to bind a signaling molecule and induce a biochemical response in the cell. Many receptors have a binding site exposed on the cell surface and an effector domain within the cell, which may have enzymatic activity or may undergo a conformational change detected by other proteins within the cell.

Antibodies are protein components of adaptive immune system whose main function is to bind antigens, or foreign substances in the body, and target them for destruction. Antibodies can be secreted into the extracellular environment or anchored in the membranes of specialized B cells known as plasma cells. Whereas enzymes are limited in their binding affinity for their substrates by the necessity of conducting their reaction, antibodies have no such constraints. An antibody's binding affinity to its target is extraordinarily high.

Many ligand transport proteins bind particular small biomolecules and transport them to other locations in the body of a multicellular organism. These proteins must have a high binding affinity when their ligand is present in high concentrations, but must also release the ligand when it is present at low concentrations in the target tissues. The canonical example of a ligand-binding protein is haemoglobin, which transports oxygen from the lungs to other organs and tissues in all vertebrates and has close homologs in every biological kingdom. Lectins are sugar-binding proteins which are highly specific for their sugar moieties. Lectins typically play a role in biological recognition phenomena involving cells and proteins. Receptors and hormones are highly specific binding proteins.

Transmembrane proteins can also serve as ligand transport proteins that alter the permeability of the cell membrane to small molecules and ions. The membrane alone has a hydrophobic core through which polar or charged molecules cannot diffuse. Membrane proteins contain internal channels that allow such molecules to enter and exit the cell. Many ion channel proteins are specialized to select for only a particular ion; for example, potassium and sodium channels often discriminate for only one of the two ions.

Structural proteins

Structural proteins confer stiffness and rigidity to otherwise-fluid biological components. Most structural proteins are fibrous proteins; for example, actin and tubulin are globular and soluble as monomers, but polymerize to form long, stiff fibers that comprise the cytoskeleton, which allows the cell to maintain its shape and size. Collagen and elastin are critical components of connective tissue such as cartilage, and keratin is found in hard or filamentous structures such as hair, nails, feathers, hooves, and some animal shells.

Other proteins that serve structural functions are motor proteins such as myosin, kinesin, and dynein, which are capable of generating mechanical forces. These proteins are crucial for cellular motility of single celled organisms and the sperm of many multicellular organisms which reproduce sexually. They also generate the forces exerted by contracting muscles.

Methods of study

As some of the most commonly studied biological molecules, the activities and structures of proteins are examined both *in vitro* and *in vivo*. *In vitro* studies of purified proteins in controlled environments are useful for learning how a protein carries out its function: for

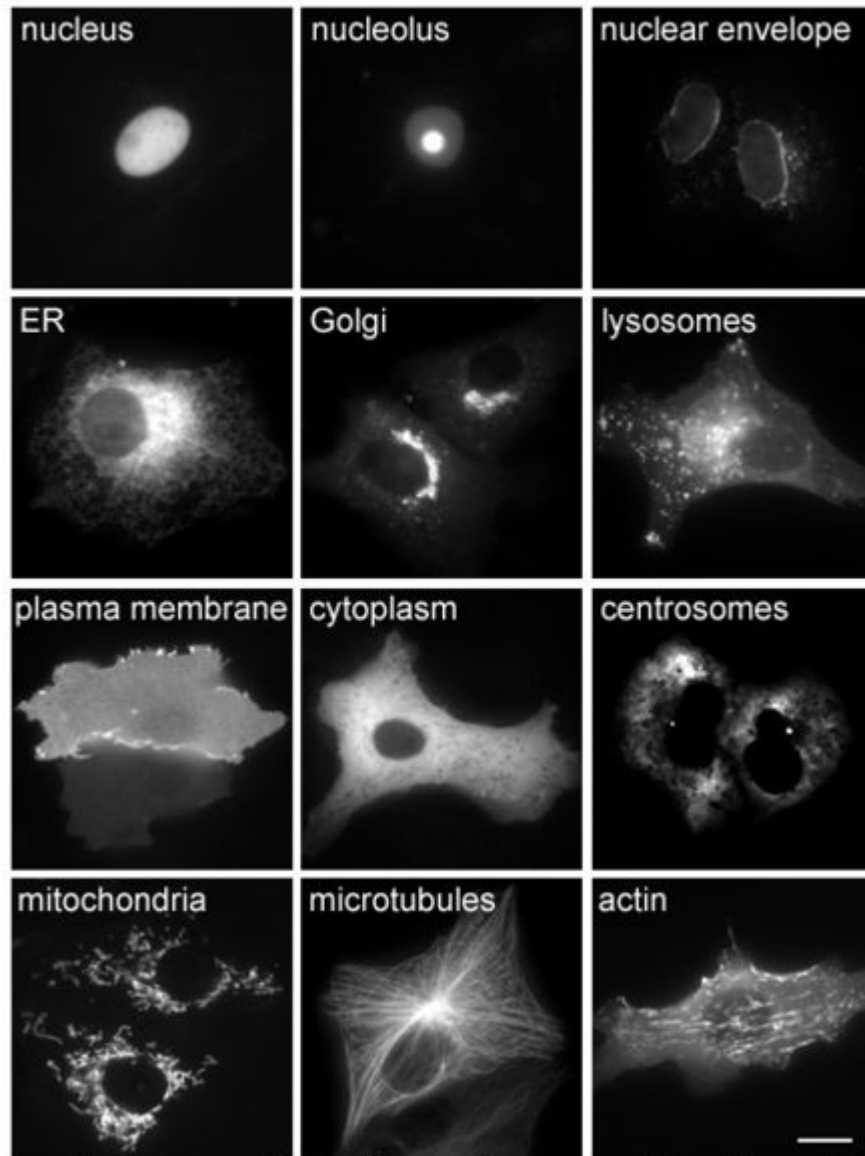
example, enzyme kinetics studies explore the chemical mechanism of an enzyme's catalytic activity and its relative affinity for various possible substrate molecules. By contrast, *in vivo* experiments on proteins' activities within cells or even within whole organisms can provide complementary information about where a protein functions and how it is regulated.

Protein purification

In order to perform *in vitro* analysis, a protein must be purified away from other cellular components. This process usually begins with cell lysis, in which a cell's membrane is disrupted and its internal contents released into a solution known as a crude lysate. The resulting mixture can be purified using ultracentrifugation, which fractionates the various cellular components into fractions containing soluble proteins; membrane lipids and proteins; cellular organelles, and nucleic acids. Precipitation by a method known as salting out can concentrate the proteins from this lysate. Various types of chromatography are then used to isolate the protein or proteins of interest based on properties such as molecular weight, net charge and binding affinity. The level of purification can be monitored using various types of gel electrophoresis if the desired protein's molecular weight and isoelectric point are known, by spectroscopy if the protein has distinguishable spectroscopic features, or by enzyme assays if the protein has enzymatic activity. Additionally, proteins can be isolated according their charge using electrofocusing.

For natural proteins, a series of purification steps may be necessary to obtain protein sufficiently pure for laboratory applications. To simplify this process, genetic engineering is often used to add chemical features to proteins that make them easier to purify without affecting their structure or activity. Here, a "tag" consisting of a specific amino acid sequence, often a series of histidine residues (a "His-tag"), is attached to one terminus of the protein. As a result, when the lysate is passed over a chromatography column containing nickel, the histidine residues ligate the nickel and attach to the column while the untagged components of the lysate pass unimpeded. A number of different tags have been developed to help researchers purify specific proteins from complex mixtures.

Cellular localization



with friendly permission of Jeremy Simpson and Rainer Pepperkok

Proteins in different cellular compartments and structures tagged with green fluorescent protein (here, white)

The study of proteins *in vivo* is often concerned with the synthesis and localization of the protein within the cell. Although many intracellular proteins are synthesized in the cytoplasm and membrane-bound or secreted proteins in the endoplasmic reticulum, the specifics of how proteins are targeted to specific organelles or cellular structures is often unclear. A useful technique for assessing cellular localization uses genetic engineering to express in a cell a fusion protein or chimera consisting of the natural protein of interest linked to a "reporter" such as green fluorescent protein (GFP). The fused protein's position within the cell can be cleanly and efficiently visualized using microscopy, as shown in the figure opposite.

Other methods for elucidating the cellular location of proteins requires the use of known compartmental markers for regions such as the ER, the Golgi, lysosomes/vacuoles, mitochondria, chloroplasts, plasma membrane, etc. With the use of fluorescently tagged versions of these markers or of antibodies to known markers, it becomes much simpler to identify the localization of a protein of interest. For example, indirect immunofluorescence will allow for fluorescence colocalization and demonstration of location. Fluorescent dyes are used to label cellular compartments for a similar purpose.

Other possibilities exist, as well. For example, immunohistochemistry usually utilizes an antibody to one or more proteins of interest that are conjugated to enzymes yielding either luminescent or chromogenic signals that can be compared between samples, allowing for localization information. Another applicable technique is cofractionation in sucrose (or other material) gradients using isopycnic centrifugation. While this technique does not prove colocalization of a compartment of known density and the protein of interest, it does increase the likelihood, and is more amenable to large-scale studies.

Finally, the gold-standard method of cellular localization is immunoelectron microscopy. This technique also uses an antibody to the protein of interest, along with classical electron microscopy techniques. The sample is prepared for normal electron microscopic examination, and then treated with an antibody to the protein of interest that is conjugated to an extremely electro-dense material, usually gold. This allows for the localization of both ultrastructural details as well as the protein of interest.

Through another genetic engineering application known as site-directed mutagenesis, researchers can alter the protein sequence and hence its structure, cellular localization, and susceptibility to regulation. This technique even allows the incorporation of unnatural amino acids into proteins, using modified tRNAs, and may allow the rational design of new proteins with novel properties.

Proteomics and bioinformatics

The total complement of proteins present at a time in a cell or cell type is known as its proteome, and the study of such large-scale data sets defines the field of proteomics, named by analogy to the related field of genomics. Key experimental techniques in proteomics include 2D electrophoresis, which allows the separation of a large number of proteins, mass spectrometry, which allows rapid high-throughput identification of proteins and sequencing of peptides (most often after in-gel digestion), protein microarrays, which allow the detection of the relative levels of a large number of proteins present in a cell, and two-hybrid screening, which allows the systematic exploration of protein–protein interactions. The total complement of biologically possible such interactions is known as the interactome. A systematic attempt to determine the structures of proteins representing every possible fold is known as structural genomics.

The large amount of genomic and proteomic data available for a variety of organisms, including the human genome, allows researchers to efficiently identify homologous proteins in distantly related organisms by sequence alignment. Sequence profiling tools

can perform more specific sequence manipulations such as restriction enzyme maps, open reading frame analyses for nucleotide sequences, and secondary structure prediction. From this data phylogenetic trees can be constructed and evolutionary hypotheses developed using special software like ClustalW regarding the ancestry of modern organisms and the genes they express. The field of bioinformatics seeks to assemble, annotate, and analyze genomic and proteomic data, applying computational techniques to biological problems such as gene finding and cladistics.

Structure prediction and simulation

Complementary to the field of structural genomics, protein structure prediction seeks to develop efficient ways to provide plausible models for proteins whose structures have not yet been determined experimentally. The most successful type of structure prediction, known as homology modeling, relies on the existence of a "template" structure with sequence similarity to the protein being modeled; structural genomics' goal is to provide sufficient representation in solved structures to model most of those that remain. Although producing accurate models remains a challenge when only distantly related template structures are available, it has been suggested that sequence alignment is the bottleneck in this process, as quite accurate models can be produced if a "perfect" sequence alignment is known. Many structure prediction methods have served to inform the emerging field of protein engineering, in which novel protein folds have already been designed. A more complex computational problem is the prediction of intermolecular interactions, such as in molecular docking and protein–protein interaction prediction.

The processes of protein folding and binding can be simulated using such technique as molecular mechanics, in particular, molecular dynamics and Monte Carlo, which increasingly take advantage of parallel and distributed computing (Folding@Home project; molecular modeling on GPU). The folding of small alpha-helical protein domains such as the villin headpiece and the HIV accessory protein have been successfully simulated *in silico*, and hybrid methods that combine standard molecular dynamics with quantum mechanics calculations have allowed exploration of the electronic states of rhodopsins.

Nutrition

Most microorganisms and plants can biosynthesize all 20 standard amino acids, while animals (including humans) must obtain some of the amino acids from the diet. The amino acids that an organism cannot synthesize on its own are referred to as essential amino acids. Key enzymes that synthesize certain amino acids are not present in animals — such as aspartokinase, which catalyzes the first step in the synthesis of lysine, methionine, and threonine from aspartate. If amino acids are present in the environment, microorganisms can conserve energy by taking up the amino acids from their surroundings and downregulating their biosynthetic pathways.

In animals, amino acids are obtained through the consumption of foods containing protein. Ingested proteins are then broken down into amino acids through digestion,

which typically involves denaturation of the protein through exposure to acid and hydrolysis by enzymes called proteases. Some ingested amino acids are used for protein biosynthesis, while others are converted to glucose through gluconeogenesis, or fed into the citric acid cycle. This use of protein as a fuel is particularly important under starvation conditions as it allows the body's own proteins to be used to support life, particularly those found in muscle. Amino acids are also an important dietary source of nitrogen.

History and etymology

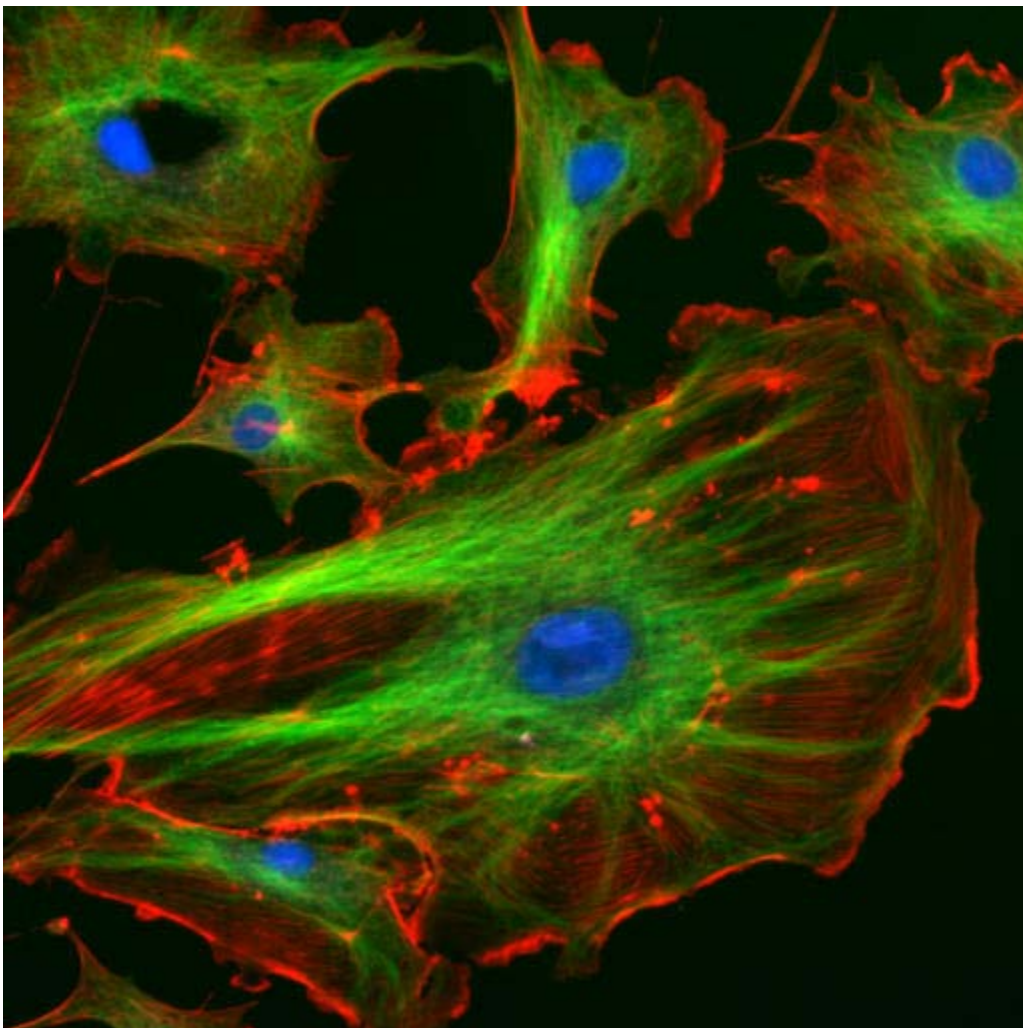
Proteins were recognized as a distinct class of biological molecules in the eighteenth century by Antoine Fourcroy and others, distinguished by the molecules' ability to coagulate or flocculate under treatments with heat or acid. Noted examples at the time included albumin from egg whites, blood serum albumin, fibrin, and wheat gluten. Dutch chemist Gerhardus Johannes Mulder carried out elemental analysis of common proteins and found that nearly all proteins had the same empirical formula, $C_{400}H_{620}N_{100}O_{120}P_1S_1$. He came to the erroneous conclusion that they might be composed of a single type of (very large) molecule. The term "protein" to describe these molecules was proposed in 1838 by Mulder's associate Jöns Jakob Berzelius; protein is derived from the Greek word *πρωτεῖος* (*proteios*), meaning "primary", "in the lead", or "standing in front". Mulder went on to identify the products of protein degradation such as the amino acid leucine for which he found a (nearly correct) molecular weight of 131 Da.

The difficulty in purifying proteins in large quantities made them very difficult for early protein biochemists to study. Hence, early studies focused on proteins that could be purified in large quantities, e.g., those of blood, egg white, various toxins, and digestive/metabolic enzymes obtained from slaughterhouses. In the 1950s, the Armour Hot Dog Co. purified 1 kg of pure bovine pancreatic ribonuclease A and made it freely available to scientists; this gesture helped ribonuclease A become a major target for biochemical study for the following decades.

Linus Pauling is credited with the successful prediction of regular protein secondary structures based on hydrogen bonding, an idea first put forth by William Astbury in 1933. Later work by Walter Kauzmann on denaturation, based partly on previous studies by Kaj Linderstrøm-Lang, contributed an understanding of protein folding and structure mediated by hydrophobic interactions. In 1949 Fred Sanger correctly determined the amino acid sequence of insulin, thus conclusively demonstrating that proteins consisted of linear polymers of amino acids rather than branched chains, colloids, or cyclols. The first atomic-resolution structures of proteins were solved by X-ray crystallography in the 1960s and by NMR in the 1980s. As of 2009, the Protein Data Bank has over 55,000 atomic-resolution structures of proteins. In more recent times, cryo-electron microscopy of large macromolecular assemblies and computational protein structure prediction of small protein domains are two methods approaching atomic resolution.

Chapter 16

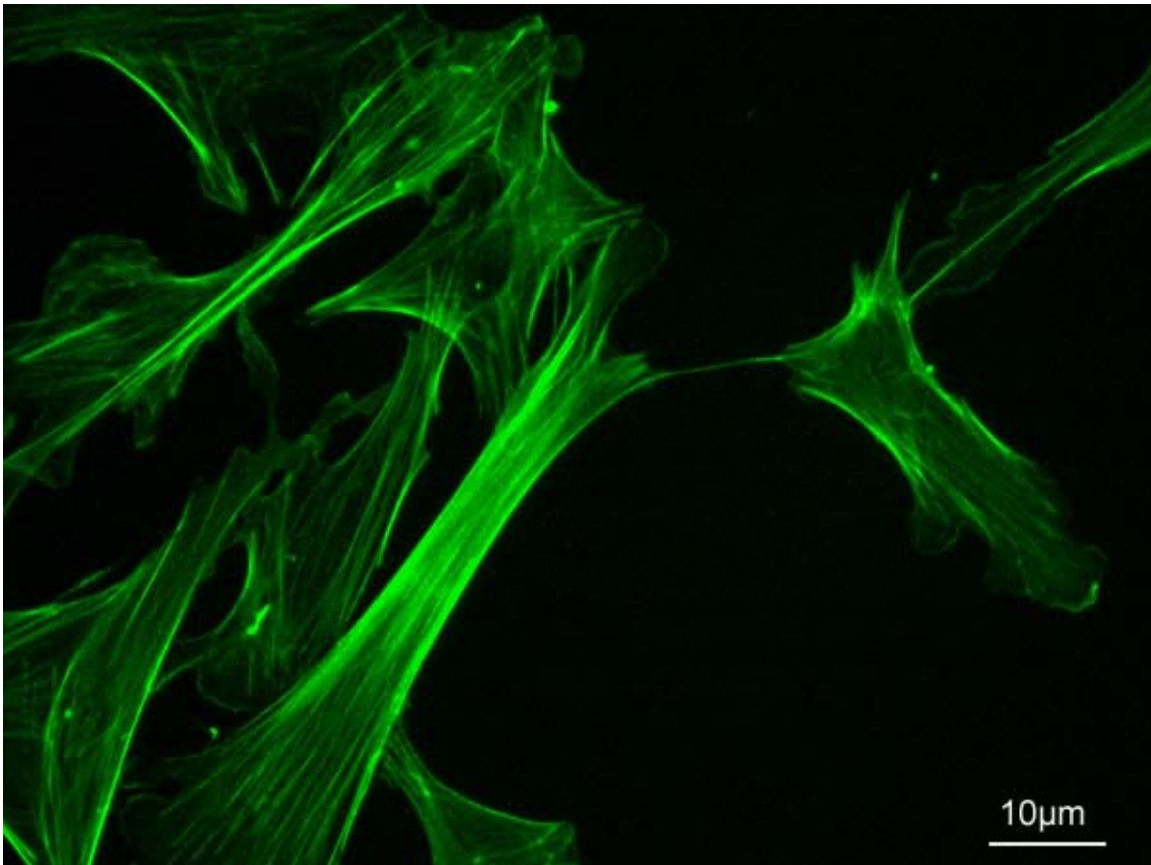
Cytoskeleton



The eukaryotic cytoskeleton. Actin filaments are shown in red, microtubules in green, and the nuclei are in blue.

The **cytoskeleton** (also CSK) is a cellular "scaffolding" or "skeleton" contained within the cytoplasm and is made out of protein. The cytoskeleton is present in all cells; it was once thought to be unique to eukaryotes, but recent research has identified the prokaryotic cytoskeleton. It has structures such as flagella, cilia and lamellipodia and plays important roles in both intracellular transport (the movement of vesicles and organelles, for example) and cellular division. The concept of a protein mosaic that dynamically coordinated cytoplasmic biochemistry was proposed by Rudolph Peters in 1929 while the term (*cytosquelette*, in French) was first introduced by French embryologist Paul Wintrebert in 1931.

The eukaryotic cytoskeleton



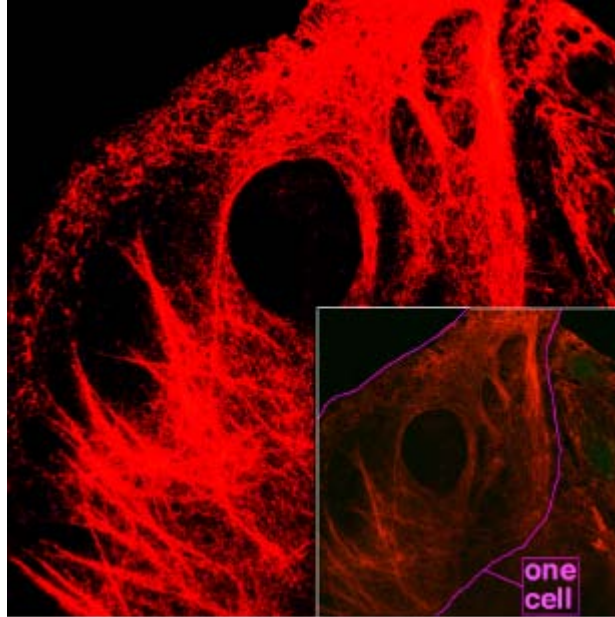
Actin cytoskeleton of mouse embryo fibroblasts, stained with phalloidin.

Eukaryotic cells contain three main kinds of cytoskeletal filaments, which are microfilaments, intermediate filaments, and microtubules. The cytoskeleton provides the cell with structure and shape, and by excluding macromolecules from some of the cytosol it adds to the level of macromolecular crowding in this compartment. Cytoskeletal elements interact extensively and intimately with cellular membranes.

Microfilaments

These are the thinnest filaments of the cytoskeleton. They are composed of linear polymers of actin subunits, and generate force by elongation at one end of the filament coupled with shrinkage at the other, causing net movement of the intervening strand. They also act as tracks for the movement of myosin molecules that attach to the microfilament and "walk" along them.

Intermediate filaments



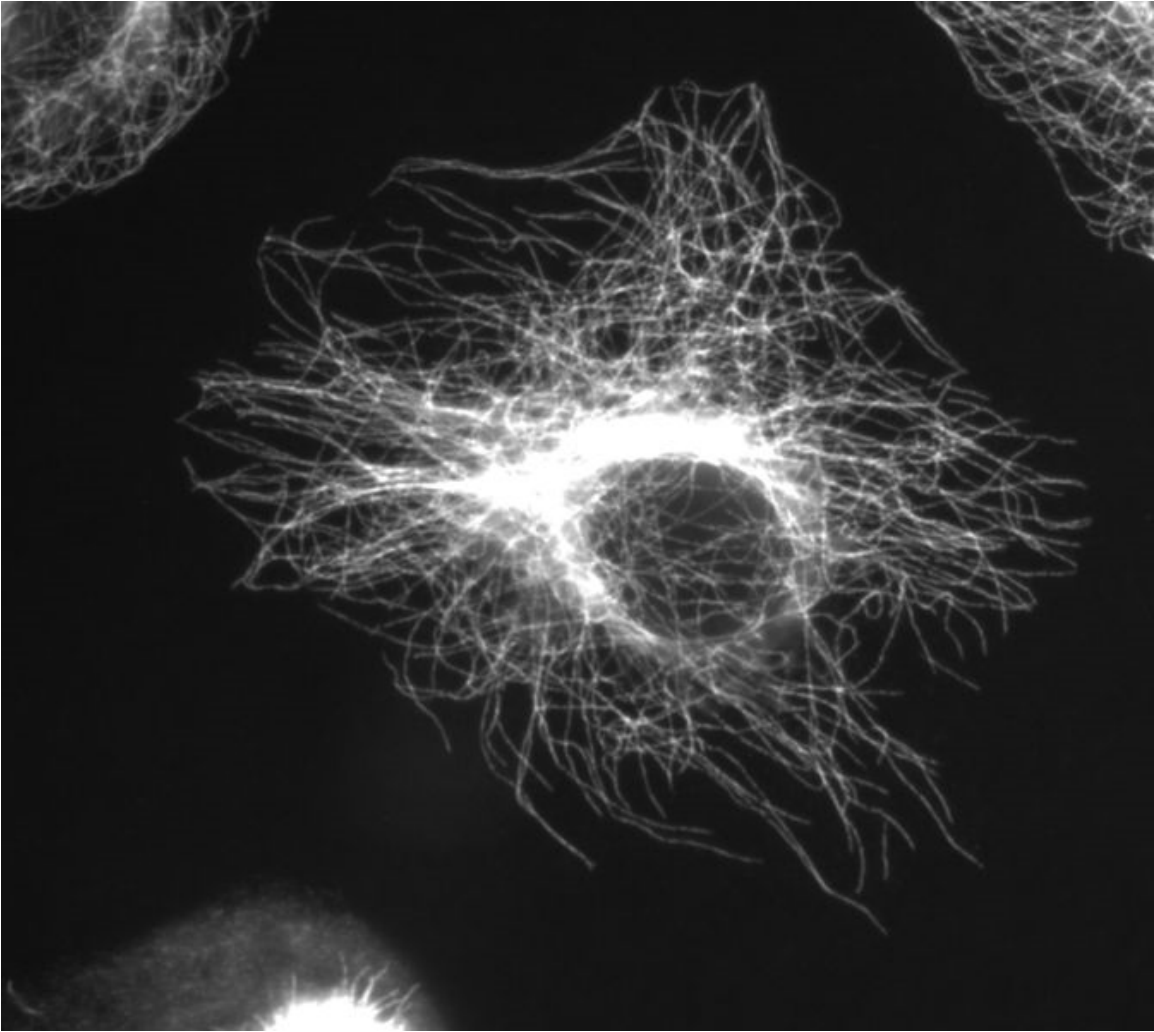
Microscopy of keratin filaments inside cells.

These filaments, around 10 nanometers in diameter, are more stable (strongly bound) than actin filaments, and heterogeneous constituents of the cytoskeleton. Although little work has been done on intermediate filaments in plants, there is some evidence that cytosolic intermediate filaments might be present, and plant nuclear filaments have been detected. Like actin filaments, they function in the maintenance of cell-shape by bearing tension (microtubules, by contrast, resist compression. It may be useful to think of micro- and intermediate filaments as cables, and of microtubules as cellular support beams). Intermediate filaments organize the internal tridimensional structure of the cell, anchoring organelles and serving as structural components of the nuclear lamina and sarcomeres. They also participate in some cell-cell and cell-matrix junctions.

Different intermediate filaments are:

- made of vimentins, being the common structural support of many cells.
- made of keratin, found in skin cells, hair and nails.
- neurofilaments of neural cells.
- made of lamin, giving structural support to the nuclear envelope.

Microtubules



Microtubules in a gel fixated cell.

Microtubules are hollow cylinders about 23 nm in diameter (lumen = approximately 15nm in diameter), most commonly comprising 13 protofilaments which, in turn, are polymers of alpha and beta tubulin. They have a very dynamic behaviour, binding GTP for polymerization. They are commonly organized by the centrosome.

In nine triplet sets (star-shaped), they form the centrioles, and in nine doublets oriented about two additional microtubules (wheel-shaped) they form cilia and flagella. The latter formation is commonly referred to as a "9+2" arrangement, wherein each doublet is connected to another by the protein dynein. As both flagella and cilia are structural components of the cell, and are maintained by microtubules, they can be considered part of the cytoskeleton.

They play key roles in:

- intracellular transport (associated with dyneins and kinesins, they transport organelles like mitochondria or vesicles).
- the axoneme of cilia and flagella.
- the mitotic spindle.
- synthesis of the cell wall in plants.

Comparison

Cytoskeleton type	Diameter (nm)	Structure	Subunit examples
Microfilaments	6	double helix	actin <ul style="list-style-type: none"> • vimentin (mesenchyme) • glial fibrillary acidic protein (glial cells)
Intermediate filaments	10	two anti-parallel helices/dimers, forming tetramers	<ul style="list-style-type: none"> • neurofilament proteins (neuronal processes) • keratins (epithelial cells) • nuclear lamins
Microtubules	23	protofilaments, in turn consisting of tubulin subunits	α - and β -tubulin

The prokaryotic cytoskeleton

The cytoskeleton was previously thought to be a feature only of eukaryotic cells, but homologues to all the major proteins of the eukaryotic cytoskeleton have recently been found in prokaryotes. Although the evolutionary relationships are so distant that they are not obvious from protein sequence comparisons alone, the similarity of their three-dimensional structures and similar functions in maintaining cell shape and polarity provides strong evidence that the eukaryotic and prokaryotic cytoskeletons are truly homologous. However, some structures in the bacterial cytoskeleton may have yet to be identified.

FtsZ

FtsZ was the first protein of the prokaryotic cytoskeleton to be identified. Like tubulin, FtsZ forms filaments in the presence of GTP, but these filaments do not group into tubules. During cell division, FtsZ is the first protein to move to the division site, and is essential for recruiting other proteins that synthesize the new cell wall between the dividing cells.

MreB and ParM

Prokaryotic actin-like proteins, such as MreB, are involved in the maintenance of cell shape. All non-spherical bacteria have genes encoding actin-like proteins, and these proteins form a helical network beneath the cell membrane that guides the proteins involved in cell wall biosynthesis.

Some plasmids encode a partitioning system that involves an actin-like protein ParM. Filaments of ParM exhibit dynamic instability, and may partition plasmid DNA into the dividing daughter cells by a mechanism analogous to that used by microtubules during eukaryotic mitosis.

Crescentin

The bacterium *Caulobacter crescentus* contains a third 3rd protein, crescentin, that is related to the intermediate filaments of eukaryotic cells. Crescentin is also involved in maintaining cell shape, such as helical and vibrioid forms of bacteria, but the mechanism by which it does this is currently unclear.

History

Microtrabeculae

A fourth eukaryotic cytoskeletal element, *microtrabeculae*, was proposed by Keith Porter based on images obtained from high-voltage electron microscopy of whole cells in the 1970s. The images showed short, filamentous structures of unknown molecular composition associated with known cytoplasmic structures. Porter proposed that this microtrabecular structure represented a novel filamentous network distinct from microtubules, filamentous actin, or intermediate filaments. It is now generally accepted that microtrabeculae are nothing more than an artifact of certain types of fixation treatment, although we have yet to fully understand the complexity of the cell's cytoskeleton.