**1**
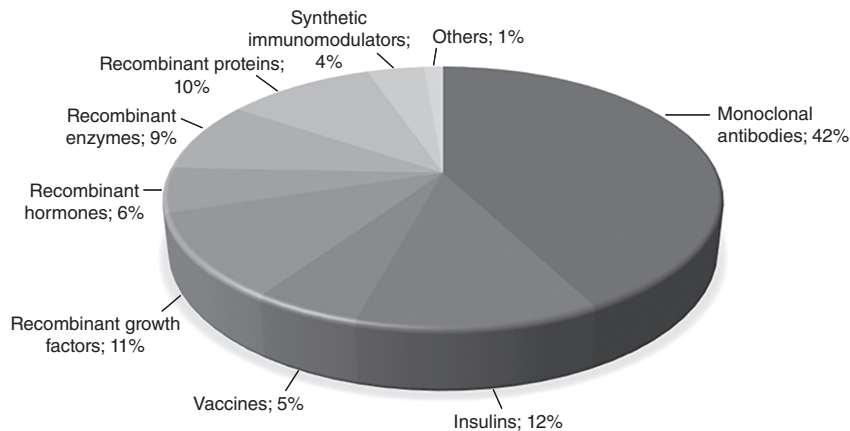
# Downstream Processing of Biotechnology Products

## 1.1 Introduction

Biological products are important for many applications including biotransformations, diagnostics, research and development, and in the food, pharmaceutical, and cosmetics industries. For certain applications, biological products can be used as crude extracts with little or no purification. However, biopharmaceuticals typically require exquisite purity, making downstream processing a critical component of the overall process. From the regulatory viewpoint, the production process and its subsequent clinical testing define the biopharmaceutical product. As a result, proper definition of effective and efficient downstream processing steps is crucial early in process development. Currently, proteins and, in particular, antibodies are the most important biopharmaceuticals [1]. The history of their development as industrial products goes back more than half a century. Blood plasma fractionation was the first full-scale biopharmaceutical industry with a current annual production in the 100-ton scale [2, 3]. Precipitation with organic solvents has been and continues to be the principal purification tool in plasma fractionation; although, recently, chromatographic separation processes have also been integrated in this industry. Antivenom antibodies and other antitoxins extracted from animal sources are other examples of early biopharmaceuticals that are also purified by a combination of precipitation, filtration, and chromatography. In contrast, current biopharmaceuticals are almost exclusively produced with recombinant DNA technology. For these products, chromatography and membrane filtration serve as the main tools for purification.

Figure 1.1 depicts the 2018 market share of different biopharmaceuticals. Approximately 42% are antibodies or antibody fragments, nearly 11% are erythropoietins, and about 12% are insulins. The rest are enzymes, growth factors, and cytokines. Although many nonproteinaceous biomolecules such as plasmids, virus, extracellular vehicles, and cells for therapy are currently being developed, it is likely that proteins will continue to dominate as biopharmaceuticals. Proteins are well tolerated, can be highly potent, and often possess a long half-life after administration, making them effective therapeutics. Some of these properties also make proteins potentially useful in cosmetics, although applications in this field are complicated in part by the US and European legal frameworks that do not allow application of pharmacologically active compounds in cosmetics.

**Figure 1.1** Biopharmaceuticals market share in 2018. By this year, 285 biological proteins have gained regulatory approval in the United States and the European Union. Source: Adapted from Walch 2018 [1].

Currently, only a few proteins are used in this area, such as botulinum toxin, Botox®, human growth factor, or fibroblast growth factor for skin care [4]. Pharmaceutical compounds must be exclusively administered by physicians and thus are not considered to be cosmetics. Often, cosmetic products with low amount of bioactive compounds are referred to as "cosmeceuticals." However, according to the US Food and Drug Administration (FDA), the Food, Drug, and Cosmetic Act does not recognize any such category as cosmeceuticals. A product can be a drug, a cosmetic, or a combination of both, but the term cosmeceutical is not recognized in the current regulatory framework.

## 1.2 Bioproducts and Their Contaminants

This section gives an overview of the chemical and biophysical properties of proteins and their main contaminants such as DNA and endotoxins. The description is not comprehensive; only biophysical properties relevant to chromatographic purification will be considered. A detailed description of the chemistry of proteins and DNA would exceed the scope of this book and can be found in a number of biochemistry or molecular biology texts [5, 6].

### 1.2.1 Biomolecular Chemistry and Structure

#### 1.2.1.1 Proteins

Proteins constitute a large class of amphoteric biopolymers with molecular mass ranging from 5 to 20 000 kDa, which are based on amino acids as building blocks. Enormous variations in structure and properties exist within this class. Insulin, for example, a peptide with a molecular mass of 5808 Da, has a relatively simple and well-defined structure. On the other hand, human von Willebrand factor, a large multimeric glycoprotein with a molecular mass of 20 000 kDa, has

an extremely complex structure consisting of up to 80 subunits, each 250 kDa in mass. Most proteins have molecular masses well within these two extremes, typically between 15 and 200 kDa. Proteins are generally rather compact molecules; yet they are flexible enough to undergo substantial conformational change in different environments, at interphases, upon binding of substrates, or upon adsorption on surfaces.

Proteins are highly structured molecules, and their structure is generally critical to their biological function. This structure is organized in four different levels: *primary*, *secondary*, *tertiary*, and *quaternary*. The primary structure is determined by the amino acid sequence, and the secondary structure is determined by folded elements or domains in the polypeptide chain. The tertiary structure is defined by the association of multiple secondary structure domains. Finally, the quaternary structure is defined by the association of multiple folded polypeptide chains. The final result is a complex three-dimensional superstructure linked by various intra- and intermolecular interactions. Often, nonamino acid elements are incorporated into a protein. Well-known examples are prosthetic groups in enzymes and iron-carrying heme groups in oxygen transport or storage proteins such as hemoglobin or myoglobin.

### 1.2.1.2 Primary Structure

The building blocks of proteins are amino acids. During biosynthesis, following transcription and translation, these molecules are linked together via peptide bonds forming a polypeptide chain in a sequence that is uniquely determined by the genetic code. The general structure of amino acids and the formation of a peptide bond are depicted in Figure 1.2. The order in which the amino acids are arranged in the polypeptide chain defines the protein's *primary structure.* Note that although amino acids are chiral molecules with L- and D-isomers, only the L-isomer is found in natural proteins. The 20 amino acids naturally found in proteins are listed in Table 1.1. In typical proteins, the average molecular mass of the amino acid components is 109 Da. Thus, the approximate molecular mass of a protein can be easily estimated from the number of amino acids in the polypeptide chain.

The peptide bond formed when amino acids are linked together has a partial double-bond character and is thus planar. This structure restricts rotation in the peptide chain, making free rotation possible only in two out of three bonds. As a consequence, unique structures are formed depending on the particular sequence of amino acids. Certain conformations are not allowed owing to the restricted rotation, while others are energetically favored owing to the formation of hydrogen bonds and other intramolecular interactions. The amino acid

**Figure 1.2** General structure of amino acids and formation of a peptide bond.

**Table 1.1** The proteinogenic amino acids, including 3- and 1-letter code, the structure of their R-group, relative abundance in *E. coli*, molecular mass, and $pK_a$ of the R-group.
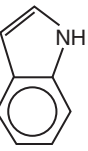
| Name | 3-Letter code | 1-Letter code | R-group | Abundance in *E. coli* (%) | Molecular mass | $pK_a$ of R-group |
|------|------|------|------|------|------|------|
| *Hydrophobic R-groups* | | | | | | |
| Alanine | Ala | A | $-CH_3$ | 13.0 | 89 | |
| Valine | Val | V | $-CH$ with $CH_3$, $CH_3$ | 6.0 | 117 | |
| Proline | Pro | P | ring structure | 4.6 | 115 | |
| Leucine | Leu | L | $-CH_2-CH-CH_3$ with $CH_3$ | 7.8 | 131 | |
| Isoleucine | Ile | I | $-CH-CH_2-CH_3$ with $CH_3$ | 4.4 | 131 | |
| Methionine | Met | M | $-CH_2-CH_2-S-CH_3$ | 3.8 | 149 | |
| Phenylalanine | Phe | F | $-CH_2-$ phenyl | 3.0 | 165 | |
| Tryptophan | Trp | W | $-CH_2-$ indole | 1.0 | 204 | |
| *Polar but uncharged R-groups* | | | | | | |
| Glycine | Gly | G | $-H$ | 7.8 | 75 | |
| Serine | Ser | S | $-CH_2OH$ | 6.0 | 105 | |
| Threonine | Thr | T | $-CH-CH_3$ with $OH$ | 4.6 | 119 | |
| Cysteine | Cys | C | $-CH_2-SH$ | 1.8 | 121 | 8.5 |
| Asparagine | Asn | N | $-CH_2-C-NH_2$ with $O$ | 11.4 | 132 | |
| Glutamine | Gln | Q | $-CH_2-CH_2-C-NH_2$ with $O$ | 10.8 | 146 | |
| Tyrosine | Tyr | Y | $-CH_2-$ phenyl $-OH$ | 2.2 | 181 | 10.0 |

**Table 1.1** (Continued)

| Name | 3-Letter code | 1-Letter code | R-group | Abundance in *E. coli* (%) | Molecular mass | p$K_a$ of R-group |
|---|---|---|---|---|---|---|
| *Acidic R-groups (negatively charged at pH ~ 6)* | | | | | | |
| Aspartic acid | Asp | D | $-CH_2-C-O^-$ with $\|$ $O$ below | 9.9 | 133 | 3.7 |
| Glutamic acid | Glu | E | $-CH_2-CH_2-C-O^-$ with $\|$ $O$ below | 12.8 | 147 | 4.2 |
| *Basic R-groups (positively charged at pH ~ 6)* | | | | | | |
| Lysine | Lys | K | $-CH_2-CH_2-CH_2-CH_2-NH_2^+$ | 7.0 | 146 | 10.5 |
| Histidine | His | H | $-CH_2-C=CH$, $HN$ $NH^+$, $C$ $H$ | 0.7 | 155 | 6.1 |
| Arginine | Arg | R | $-CH_2-CH_2-CH_2-NH-C-NH_2$ with $\|$ $NH_2^+$ below | 5.3 | 174 | 12.5 |

Note that proline is a cyclic imino acid and its structure is shown in its entirety.

side chains can be charged, polar, or hydrophobic (see Table 1.1), thereby determining the biophysical properties of a protein. The charged groups are acids and bases of differing strength or p$K_a$. Thus, these groups will determine the protein net charge as a function of pH. Hydrophobic side chains, on the other hand, will determine the hydrophobic character of the primary structure, greatly influencing folding. The amino acid cysteine and proline play special roles. By oxidation, free cysteines can form disulfide bonds or bridges yielding cystine as shown in Figure 1.3. When cysteines are parts of polypeptide chains, these bridges can be either intramolecular (along the same polypeptide chain) or between different polypeptide chains. On the one hand, these bridges contribute to the stabilization of a protein's folded structure. On the other hand, they can lead to the formation of covalently bonded multimeric protein structures.

The formation of disulfide bridges is generally reversible. Bonds formed in an oxidative environment can be broken in a reductive one destabilizing the protein's folded structure and disrupting associated forms. This property is exploited, for example, in high-resolution analytical protein separation methods such as SDS polyacrylamide gel electrophoresis (SDS-PAGE) that are often performed under reductive conditions. In this case, loss of structure ensues and the resulting elimination of associated forms allows a precise determination of the protein's molecular mass.

Proline plays a special role in defining the protein structure. Proline is a cyclic imino acid and can exist in cis and trans forms. In turn, these forms influence the conformation of the polypeptide chain. In free solution, these isomeric forms are
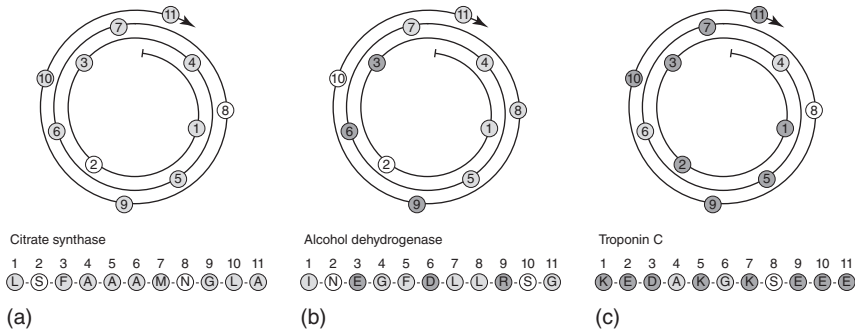
**Figure 1.3** Formation of a disulfide bond upon oxidation of two cysteines.

in equilibrium. However, in a polypeptide, the interconversion of these isomeric forms is often slow and can be the rate-limiting step in the establishment of folded protein structures.

### 1.2.1.3 Secondary Structure

The polypeptide chains found in proteins do not form knots or rings and are not β-branched. However, these chains can form α-helices, β-sheets, and loops, which define the protein's *secondary structure*. α-Helices consist of a spiral arrangement of the polypeptide chain comprising 3.6 amino acid residues per turn. The helix is stabilized by intramolecular hydrogen bonds and may have hydrophobic, amphipathic, or hydrophilic character, depending on the particular sequence of amino acids in the primary structure. Examples are given in Figure 1.4. In each case, the character of the α-helix can be predicted by placing each amino acid residue in a spiral at 100° intervals so that there will be 3.6 residues per turn. As can be seen in Figure 1.4, for citrate synthase, the hydrophobic residues are dominant and uniformly distributed so that the α-helix will be hydrophobic. In the last case, for troponin C, the charged residues are dominant but also uniformly distributed so that the resulting helix will be hydrophilic. Finally, for alcohol dehydrogenase, hydrophobic and charged residues will be distributed nonuniformly resulting in an amphipathic helix that has hydrophilic character on one side and hydrophobic character on the other.

β-Sheets are very stable secondary structure elements that also occur as a result of hydrogen bonding. Although one hydrogen bond makes up a free energy of binding ($\Delta G$) of only about 1 kJ/mol, the large number of such bonds in β-sheets make them highly stable. As can be seen in Figure 1.5, β-sheets have a planar structure, which can be parallel, antiparallel, or mixed depending on the direction of the polypeptide chains that form these structures. Formation of β-sheets is often observed during irreversible protein aggregation. Because of the strong intermolecular forces in these structures, vigorous denaturing agents are needed to break the resulting aggregates. Urea, a strong hydrogen bond breaker, can be used for this purpose. However, the high urea concentrations needed will often

**Figure 1.4** Schematic structures of (a) hydrophobic, (b) amphipathic, and (c) hydrophilic α-helices. Hydrophobic amino acid residues are shown in light gray, polar in white, and charged in dark gray. Source: Adapted from Branden and Tooze 1991 [7].

result in a complete destabilization and unfolding of the whole protein structure. Amyloid proteins and fibers contain large amounts of β-sheets, which explain in part the properties of these classes of aggregation-prone proteins.

Finally, loops are very flexible parts of the protein often connecting other secondary structure elements with each other. For example, loops often connect the portions of a polypeptide chain that form antiparallel of parallel β-sheets or form the links between different α-helical and β-sheet domains. Several types of loops have been described such as α and ω loops. Loops also play a critical role when different proteins are fused together in an artificial way as in the case of single-chain antibodies. These artificial antibodies are connected with loops that significantly contribute to the stability of the protein.

The relative amount of secondary structure elements present in a protein can be measured by several spectroscopic methods including circular dichroism (CD) and infrared spectroscopy. CD spectroscopy is based on the anisotropic nature of the protein. In circularly polarized light, the electric field vector has a constant length but rotates about its propagation direction. Hence, the light forms a helix in space while propagating. If this is a left-handed helix, the light is referred to as left circularly polarized and vice versa for a right-handed helix. Because of the interaction with the molecule, the electric field vector of the light traces out an elliptical path while propagating. At a given wavelength, the difference between absorbance of left circularly polarized ($A_L$) and right circularly polarized ($A_R$) light is

$$\Delta A = A_L - A_R \tag{1.1}$$

Although $\Delta A$ is the measured quantity, the results are usually reported in degrees of ellipticity [$\theta$]. Molar CD ($\varepsilon$) and molar ellipticity, [$\theta$], are readily interconverted by the equation
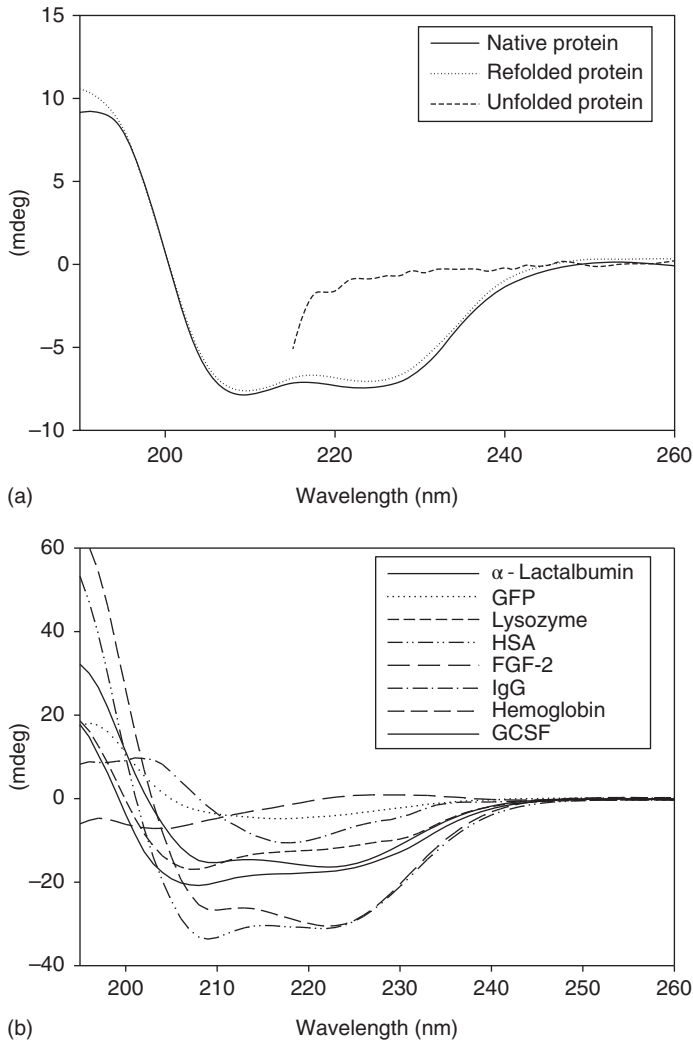
$$[\theta] = 3298.2 \cdot \Delta\varepsilon \tag{1.2}$$

The wavelength scan provides information about the secondary structure content of a protein and is a useful measure of integrity. Such measurement is often used either to follow protein refolding or to confirm the native structure of a protein

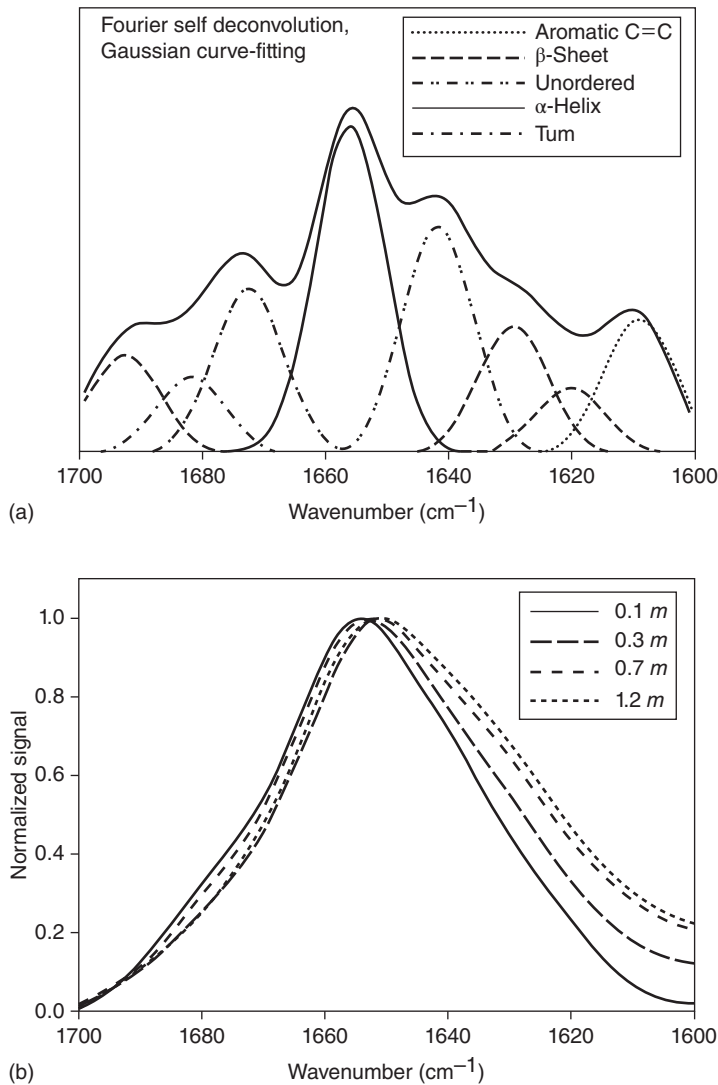**Figure 1.5** Schematic structure of parallel (a) and antiparallel (b) β-sheets in proteins.

**Figure 1.6** (a) CD-spectrum of native, refolded, and unfolded α-lactalbumin. (b) A selection of CD spectra of various proteins (GFP, green fluorescent protein; HSA, human serum albumin; FGF-2, fibroblast growthfactor 2; GCSF, granulocyte colony stimulating factor). Source: (Panel b) Reproduced from Duerkop et al. 2018 [8].
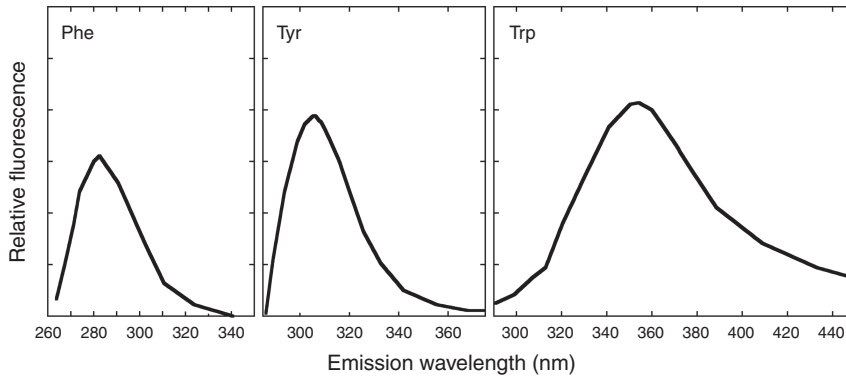
(see Figure 1.6). Different algorithms have been applied to determine the content of secondary structure elements based on these measurements and quantification is highly dependent on the particular algorithm used. Although CD spectroscopy is not sensitive enough to trace residual unfolded protein in a protein preparation, the method is well suited and accepted to study thermally or chemically induced unfolding of proteins.

Attenuated total reflectance Fourier transform infrared (ATR FT-IR) spectroscopy is also used to study conformational changes of the protein 3D structure *in situ*. With ATR FT-IR, changes of the secondary structure elements can be

assessed. ATR FT-IR allows the measurement of secondary structure elements of proteins also in suspensions and turbid solutions. The amide I band in the spectral region from 1600 to 1700 cm$^{-1}$ is used to evaluate structural changes (see Figure 1.7). As in the case of CD, application of certain algorithms allows the extraction of the secondary structure content of a protein, although, again,

**Figure 1.7** Panel (a) shows the infrared spectrum of the amide I band of a protein and its deconvolution according to the contribution of the secondary structure domains. Panel (b) shows the shift of the amide I band of the protein upon adsorption to a HIC resin at different concentrations of ammonium sulfate. The spectral shift indicates a significant change in secondary structure content at higher ammonium sulfate concentrations. Source: Ueberbacher et al. 2008 [9]. Reproduced with permission of Elsevier.

**Figure 1.8** Relative fluorescence of the amino acids Phe, Tyr, and Trp with excitation at 257, 274, and 278 nm, respectively. Source: Adapted from Schmid 1997 [10].

the extracted structure content is highly dependent on the applied algorithm. An advantage of this method is that the structure can be determined when the protein is adsorbed.
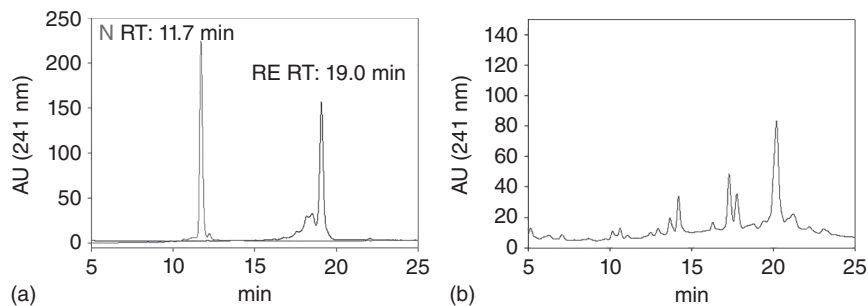
### 1.2.1.4 Tertiary Structure

The *tertiary structure* is formed when secondary structure elements (α-helices, β-sheets, and loops) are combined together in a three-dimensional arrangement. Hydrophobic interactions and disulfide bridges are primarily responsible for the stabilization of the tertiary structure. An example is the packing of amphipathic α-helices into a four-helix bundle. In this structure, the hydrophobic residues are tightly packed in its core, shielded from the surrounding aqueous environment, while polar and charged residues remain exposed on its surface.

Fluorescence spectroscopy can provide information about the location of highly hydrophobic residues, tryptophan, phenylalanine, and tyrosine, and their degree of surface exposure in such folded structures. As can be seen in Figure 1.8, these residues have characteristic fluorescence spectra, which are dependent on where they are located in the protein structure and what the folding state of the protein is. When these residues are exposed at the protein surface, the fluorescence maximum shifts providing an indication that unfolding has occurred. Thus, the extent of unfolding can be calculated when the fluorescence spectra of native and unfolded forms are known.

### 1.2.1.5 Quaternary Structure

The *quaternary structure* is established when two or more polypeptide chains are associated to form a superstructure, which, in many cases, is essential for the biological function. One of the best-known examples is hemoglobin, which consists of four polypeptide units held together by hydrogen bonding and hydrophobic interactions. In this case, the flexibility of the quaternary structure in response to oxygen binding is critical for favorable oxygen uptake in the lung where the oxygen tension is high and favorable release behavior in capillary environments where the oxygen tension is low. Antibodies are another example of proteins with

**Figure 1.9** (a) Retention of native and fully folded α-lactalbumin on a Vydac C4 reversed-phase column with 5 μm particles and a pore size of 30 nm. Mobile phase was a water–acetonitrile mixture. (b) Separation of folding intermediates of α-lactalbumin with the same column and same conditions.

quaternary structure. These molecules consist of four polypeptide chains (two light and two heavy) linked together by disulfide bridges. The resulting structure is generally quite stable, allowing antibodies to circulate freely in plasma.

#### 1.2.1.6 Folding

Although individual steps in the folding pathway can be extremely fast, the overall process of protein folding can be relatively slow. For instance, the helix–coil transition and the diffusion-limited collapse of proteins occur on time scales on the order of microseconds. On the other hand, the cis–trans isomerization of imidic bonds is a slow reaction occurring on time scales of up to several hours. As a result, in some instances, folding and chromatography occur on similar time scales so that structural rearrangements can take place during separation. When folding processes are particularly slow, chromatography can be used to resolve intermediate folding variants. For example, as can be seen in Figure 1.9, partially unfolded proteins show different retention in reversed-phase chromatography, which can be used either to analyze protein solutions during an industrial refolding process or for preparative separation of partially unfolded forms.

Protein structures are classified into several hierarchies with protein superfamilies and families. Dayhoff [11] introduced the term *protein superfamily* in 1974. Currently, the term "*folds*" is more commonly used to describe broad classes of protein structures. Table 1.2 shows the relative abundances of protein folds found in the PIR-International Protein Sequence Database. An excellent description of structural hierarchies of proteins can be found in the web site: http://supfam.mrc-lmb.cam.ac.uk/SUPERFAMILY/description.html.

Proteins have been classified into classes and folds to allow searching for common origins and evolutionary patterns. However, it should be noted that even proteins belonging to the same class may behave differently because even the exchange of a single amino acid can result in large variations in biophysical properties.

#### 1.2.1.7 Post-translational Modifications

Post-translational modifications are often critical to a protein's biological function and can dramatically impact downstream processing. Post-translational

**Table 1.2** Classes of folds found in protein databases.

| Class of protein fold | Relative abundance (%) |
|---|---|
| All α | 20–30 |
| All β | 10–20 |
| α and β with mainly parallel β-sheets (α/β) | 15–25 |
| α and β with mainly antiparallel β-sheets with segregated α- and β-regions (α + β) | 20–30 |
| Multidomain | <10 |
| Membrane and cell surface proteins | <10 |
| Small proteins (dominated by cofactors or disulfide bridges) | 5–15 |

Source: Nötling 1999 [12]. Reproduced with permission of Springer.

modifications occur after the primary structure is formed and are highly cell-specific. They can also vary with the physiological status of the cell, which, in turn, can vary, for example, when the cells are deprived of certain nutrients or are present in a low oxygen environment. Further modifications can also occur following expression. As a result, many protein-based biopharmaceuticals are highly heterogeneous and their biological and pharmacological activity is often highly influenced by the production process. Difficulties encountered in fully characterizing the corresponding broad range of molecular diversity often require that protein pharmaceuticals be defined by the process by which they are produced rather than as uniquely defined molecular entities. Although considerable effort is being devoted to develop "well-characterized biologicals," for which molecular qualities rather than processing define the product, current regulations continue to define biologicals strictly by their manufacturing process.

Post-translational modifications often represent a productivity bottleneck. At high expression rates, post-translational modifications are often altered or become incomplete when the cell's ability to perform these transformations lags behind the protein translation machinery. The result is the expression of additional protein variants, with potentially varying biological activity, stability, and biophysical properties such as solubility, charge, hydrophobicity, and size. Thus, improved fermentation or cell culture titers often have to be balanced against the increased heterogeneity of the product formed.

More than 200 types of post-translational modifications have been described in the literature. Table 1.3 summarizes the most relevant ones. The individual molecular entities produced by such modifications are called isoforms.

Two especially relevant post-translational modifications are glycosylation and deamidation as both produce changes that can influence the protein chromatographic properties. Thus, chromatography can be used as a tool to separate the corresponding isoforms. Glycosylation is the addition of carbohydrate molecules, either simple sugars or complex oligosaccharides, to the protein molecule. Glycosylation renders the protein more hydrophilic and, thus, more soluble. Additionally, however, as the terminal carbohydrates of such oligosaccharides are often
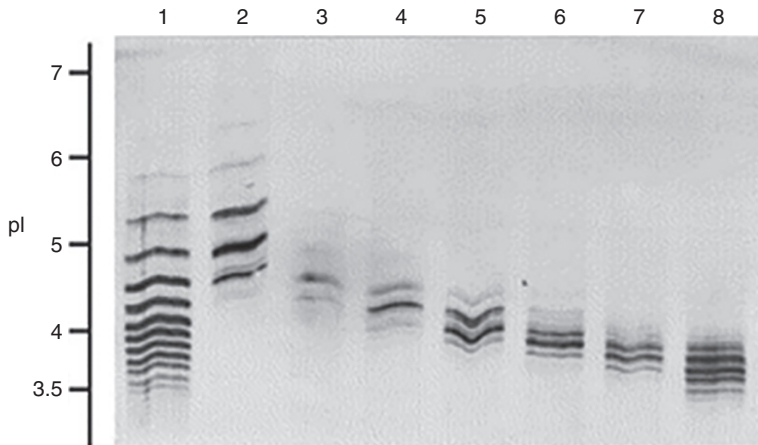
**Table 1.3** Examples of post-translational modifications of proteins.

| Modification | Characteristics | Average mass difference |
|---|---|---|
| Glycosylation | O-linked oligosaccharides bound to Ser and Thr; N-linked oligosaccharide bound to Asp | Varies with number of sugar moieties; up to several thousands |
| Phosphorylation | Ser, Thr, and Tyr, a phosphoester is formed, typical modification of allosteric proteins involved in regulation (signal transduction) | 79.9799 |
| Sulfation | Addition of sulfate to Arg and Tyr, a $C—O—SO_3$ bond | 80.0642 |
| Amidation | Addition of $–NH_2$ to C-terminus | −0.9847 |
| Acetylation | Addition of $CH_3CO-$ to N- or C-terminus | 42 |
| Hydroxylation | Addition of –OH to Lys, Pro, and Phe | 16 |
| Cyclization | Formation of pyroglutamate at N-terminal Glu | −0.9847 |
| Complexation of metals | Cys $–CH_2–S–Fe$ complexes in ferredoxins | Varies |
|  | Selenium complexes with Cys and Met |  |
|  | Copper complexes with backbone of peptide bond |  |
| Halogenation | Iodation and bromation of Tyr (3-chloro, 3-bromo) | 34, 78 |
| Desmosin formation | Desmosin is formed by condensation of Lys, frequent in elastin | −58 |
| $\gamma$-Carboxylation | In prothrombin and blood coagulation factor VII | 44.0098 |
| Hydroxyproline | Hydroxyproline formation in collagen responsible for mechanical stability | 15.9994 |
| Adenylation | Tyr residue of glutamine synthetase is adenylated | 209 |
| Methylation | Addition of methyl group to Asp, Gln, His, Lys, and Arg of flagellaprotein | 14.0269 |
| Deamidation | Asn and Gln are susceptible; both biological and processing deamidation are observed | 0.9847 |

Source: Data from http://www.expasy.ch/tools/findmod/findmod_masses.html.

neuraminic acids (generally known as sialic acid), which are negatively charged above pH 3, glycosylation also influences the net charge and isoelectric point of the protein. As a result, chromatographic separations based on the different charge of the glycovariants are possible.

As an example, Figure 1.10 shows the isoelectric focusing (IEF) gel separation of recombinant human erythropoietin (rhEPO), currently still one of the top-selling biopharmaceuticals. As can be seen in Figure 1.10, the starting material contains multiple variants with isoelectric points between 3.5 and 5.5. Loading the starting material on an anion exchange column and eluting with increasing salt concentrations result in eluted fractions that have substantially reduced heterogeneity. Later eluting fractions contain more acidic variants with lower isoelectric

**Figure 1.10** Isoelectric focusing (IEF) of rhEPO. Fractions obtained by DEAE–Sephacel chromatography: (1) starting material; (2) unadsorbed material; (3) material eluted with 0.015 M; (4) 0.03 M; (5) 0.06 M; (6) 0.15 M; (7) 0.35 M; and (8) 1 M NaCl. Source: Gokana et al. 1997 [13]. Reproduced with permission of Elsevier.

points. These variants are more negatively charged and elute only at higher salt concentrations from the positively charged anion exchanger. rhEPO is highly glycosylated and the glycovariants have different bioactivity. Thus, control of the glycosylation pattern and, in some cases, separation of certain undesirable variants are needed to maintain a consistent product quality.

Deamidation can also have dramatic effects both on bioactivity and on chromatographic behavior. Deamidation involves the chemical transformation of asparagine and glutamine, which are uncharged polar amino acids, into aspartic acid and glutamic acid, respectively, both of which are fully deprotonated and, thus, negatively charged at pH values above 5. Deamidation of asparagine residues is observed more frequently than that of glutamine residues, but the process is highly dependent on the location of these residues in the protein structure. Surface-exposed residues tend to be most affected, while those buried within the protein core are usually partially protected. Deamidation is generally facilitated by higher pH values and higher temperatures and occurs via the mechanism illustrated in Figure 1.11. In this process, an amino group is cleaved off from asparagine forming an L-cyclic imide intermediate. This intermediate is generally unstable and is further converted into L-aspartyl and L-*iso*-aspartyl peptides. Both introduce negative charge and lower the isoelectric point of the protein. It should be noted that the unstable L-cyclic amide can also undergo racemization forming a D-cyclic amide, which is further converted into -aspartyl and D-*iso*-aspartyl peptide. The net result is the introduction of D-amino acids in a protein. Removal of deamidated variants is often an important task as these variants can have different bioactivity and is a challenge for downstream processing. Separation by ion exchange chromatography is possible but often difficult as the net charge difference between native and deamidated forms can be small, resulting in low selectivity.
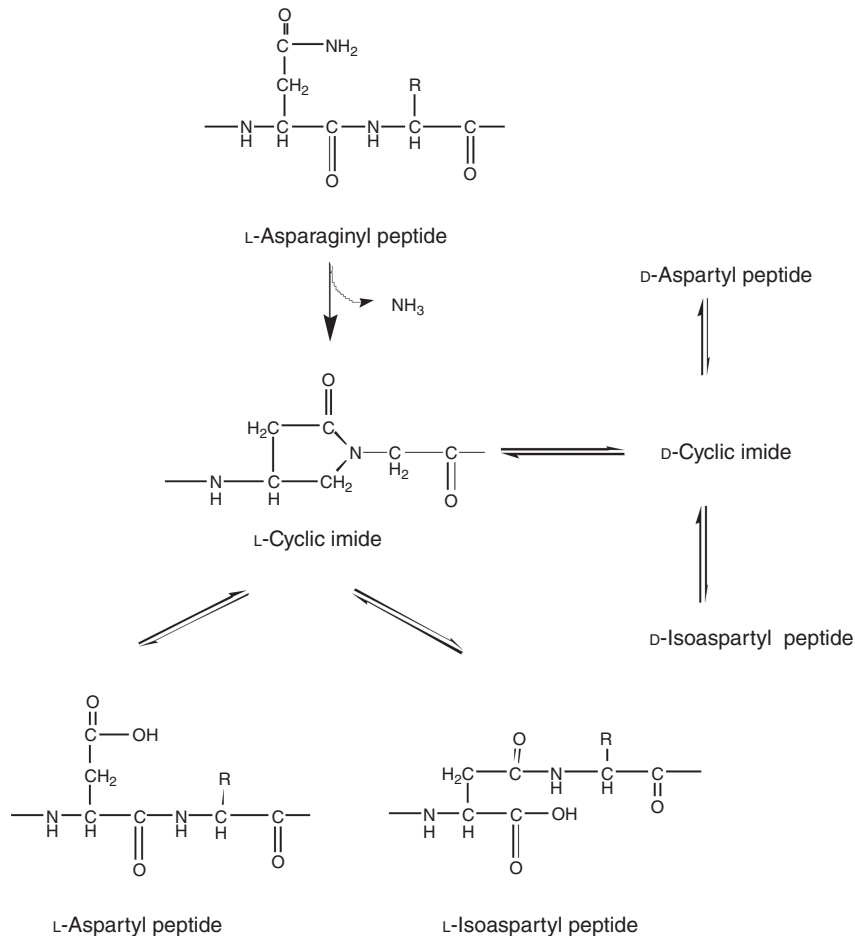
**Figure 1.11** Mechanism of deamidation of proteins associated with asparagine residues.

### 1.2.1.8 Oligonucleotides and Polynucleotides

Oligonucleotides and polynucleotides are either contaminants or may constitute the product. For example, in the production of plasmid DNA for gene therapy applications, genomic DNA is a contaminant [14]. Conversely, in the production of protein pharmaceuticals, both genomic and plasmid DNA are contaminants.

Polynucleotides are present in the cell either as deoxyribonucleic acid (DNA) or as ribonucleic acid (RNA). DNA or RNA encode genetic information. In humans, animals, and plants, DNA is the genetic material, while RNA is transcribed from it. In some other organisms, such as RNA virus, RNA is the genetic material and, in reverse fashion, the DNA is transcribed from it. The building blocks of these molecules are nucleotides, which are composed of a phosphate group, a sugar group, and a nitrogenous nucleoside group. Nucleotides are thus rather hydrophilic and negatively charged because of the acidic phosphate group. In DNA, the nucleotides are arranged in a double-stranded helical structure held
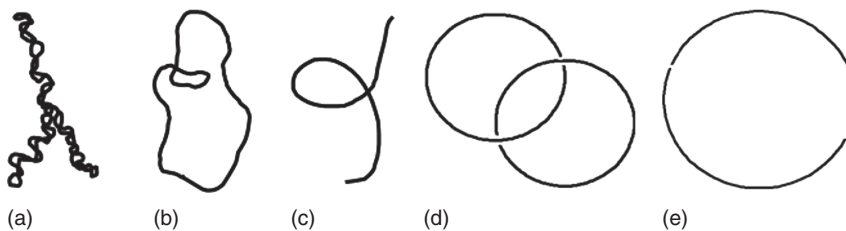
together by weak hydrogen bonds between pairs of nucleotides. The molecule resembles a twisted "ladder," where the sides are formed by the sugar and phosphate moieties and the "rungs" are formed by the nucleoside bases joined in pairs by hydrogen bonds.

There are four nucleotides in DNA, each containing a different nucleoside base: adenine (A), guanine (G), cytosine (C), and thymine (T). Base pairs form naturally only between A and T and between C and G so that the base sequence of each single strand of DNA can be simply deduced from that of its partner strand. RNA is similar to DNA in structure but contains ribose instead of deoxyribose. There are several classes of RNA molecules including messenger RNA, transfer RNA, and ribosomal RNA. These molecules play a crucial role in protein synthesis and other cell activities. miRNAs are global regulators of gene expression. They consist of noncoding double-stranded RNA molecules, comprising 19–22 nucleotides and regulate gene expression at the post-transcriptional level by forming a conserved single-stranded structure and showing antisense complementarity that was identified initially in the nematode *Caenorhabditis elegans*.

DNA and RNA are chemically very stable molecules unless DNAse or RNAse enzymes are present. In the presence of these ubiquitous enzymes, polynucleotide chains are rapidly degraded. Polynucleotides are also very sensitive to mechanical shear. Upon cell lysis, DNA and RNA are released into the culture supernatant and dramatically affect the viscosity of fermentation broths as a result of their size and filamentous structure.

Genomic DNA present in the nucleus of eukaryotic organisms is always associated with very basic proteins called histones. Plasmid DNA, on the other hand, present in the cytoplasm of prokaryotic organisms, is histone free but exists in different physical forms such as supercoiled, circular, linear, and aggregated, as illustrated in Figure 1.12.

These forms differ in size providing a basis for separation by gel electrophoresis or by size exclusion chromatography (SEC). Polynucleotides are negatively charged over a wide range of pH because of the exposed phosphate groups. Thus, they are strongly bound onto positively charged surfaces. As a result, their removal in downstream processing is conveniently and efficiently carried out with anion exchange resins or with positively charged membranes.



(a)          (b)          (c)          (d)                    (e)

**Figure 1.12** Different physical forms of plasmid DNA: supercoiled (a), open circular (b), linear (c), catenane (d), and concatemer (e). The supercoiled form normally has the highest transformation efficiency and is the predominant form found in therapeutic plasmids. The open circular form is generated when one strand is nicked while the linear form is generated with the cleavage of the double strand. Two circular forms generate a catenane or a (e) concatemer.
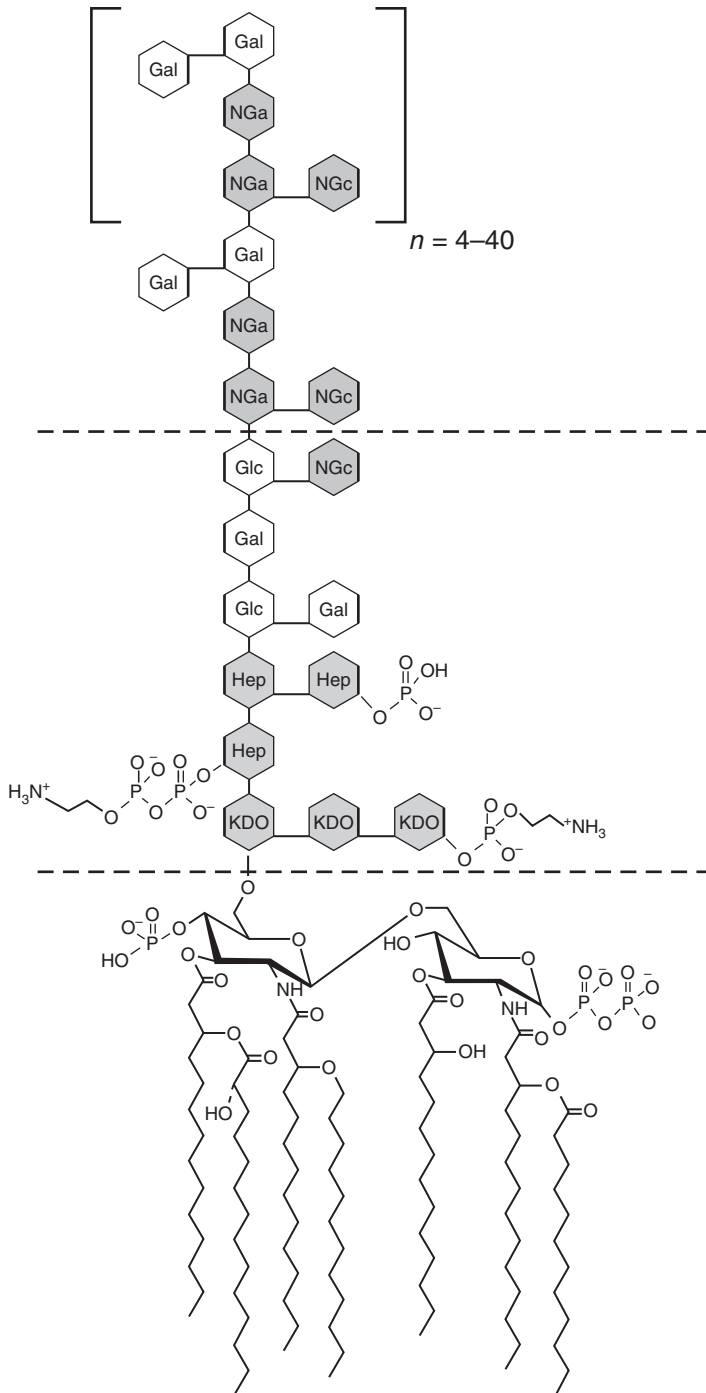
### 1.2.1.9 Endotoxins

*Endotoxins*, also known as *pyrogens*, are components of the cell wall of gram-negative bacteria. They are continuously excreted by bacteria and are ubiquitous in bioprocessing. Endotoxins are extremely toxic when entering the bloodstream and humans are among the most sensitive organisms. Thus, their essentially complete removal from the finished product is required. As shown in Figure 1.13, endotoxins are lipopolysaccharides comprising a *lipid A part*, a *core region*, and a *O-antigen* or a *S-antigen*. The lipid A part is the most conservative and is found in all endotoxins. This is also the toxic part of the molecule. The O- or S-antigens are highly variable and strain specific. The size and structure also depend on the growth conditions.

Endotoxins target immune-responsive cells such as macrophages, monocytes, endothelial cells, neutrophils, and granulocytes and can induce the expression of interleukins, tumor necrosis factor, colony-stimulating factor, leukotrienes, and oxygen radicals in these cells. As a consequence, the recipient of endotoxins will often develop tissue inflammation and fever, drop of blood pressure, shock, increase of palpitation, drop of vessel permeability, respiratory complications, or even death. The same symptoms occur with severe bacterial infections, commonly known as septic shock. Severe hepatic toxicity and hematological disorders have been observed in humans with as little as 8 ng of endotoxins per kg body weight. In contrast, endotoxins are not very toxic for many animals. For example, the $LD_{50}$ in mice is as high as 200–400 μg/animal. For parenteral biopharmaceuticals, the threshold level for intravenous applications is 5 *endotoxin units* (EU) per kg of body weight per hour. EU defines the biological activity of endotoxins with 1 EU corresponding to 100 pg of the EC-5 standard endotoxin or 120 pg of the endotoxin derived from the *Escherichia coli* strain O111:B4. Endotoxin detection is difficult and is performed with bioassays. In the past, rabbits have been used for this purpose. This time-consuming test has been replaced by the so-called limulus amebocyte lysate (LAL) test, which is based on the hemolymph of the horseshoe crab. LAL coagulates in the presence of minute amounts of endotoxins (see Figure 1.14) forming the basis for assays with endotoxin detection limits as low as of 10 pg/ml. General guidelines are described in the United States Pharmacopeia (USP) in Chapter 79 on pharmaceutical compounding and sterile preparations (CSP).

Table 1.4 provides examples of typical endotoxin content of various solutions. Endotoxins are present in large concentration in protein solutions derived from bacterial fermentations but can also be present as adventitious agents in many other systems.

In the industrial production of pharmaceuticals for parenteral use, special care is used to prevent endotoxin contamination. For example, endotoxin-free water, used in the preparation of culture media and chromatography buffers, is recirculated at high temperature in order to avoid bacterial growth and consequent formation of endotoxins. Although endotoxins are heat stable, they are destroyed at alkaline pH. Thus, cleaning processing equipment, tanks, membranes, and chromatography media with a sodium hydroxide solution is generally required to assure complete removal of these contaminants.

**Figure 1.13** Chemical structure of endotoxins. Source: Petsch and Anspach 2000 [15]. Reproduced with permission of Elsevier.

**Figure 1.14** Coagulation test for detection of endotoxins using amebocyte lysate. The lysate forms a gel in the presence of endotoxins. Source: Photograph courtesy of Associates of Cape Cod, Inc.

**Table 1.4** Examples of endotoxin concentrations in various solutions of crude and purified proteins.

| Protein source | Solution | Endotoxin (EU/ml) |
|---|---|---|
| Proteins from high-cell-density culture of *E. coli* TG:p$\lambda$FGFB | Supernatant after homogenization | $\gg$2 000 000 |
| Proteins from shaking flask culture of *E. coli* | Culture filtrate | 70 000–500 000 |
| Murine IgG1 from cell culture | Culture filtrate | $\leq$100 |
| Whey processed from milk of local supermarket | Supernatant after acid milk precipitation | $\sim$10 000 |
| Commercial preparation of BSA | Reconstituted lyophilizate at a concentration of 1 mg/ml | 50 (Supplier I) 0.5 (Supplier II) |

Source: Data from Petsch and Anspach 2000 [15].

### 1.2.2 Biochemical and Biophysical Properties

#### 1.2.2.1 UV Absorbance

The concentration of a protein in solution is often quantified by the UV absorbance, which is primarily due to the aromatic amino acids tyrosine, tryptophan, and phenylalanine and to disulfide bridges. The wavelength absorbance maxima and corresponding extinction coefficients for these components are summarized in Table 1.5.

Because of the strong absorbance of tryptophan, absorption maxima for proteins are typically around 280 nm, and this wavelength is most frequently

**Table 1.5** Absorbance characteristics of aromatic amino acids and disulfide bridges.

| Amino acid | $\lambda_{max}$ (nm) | $\varepsilon_m^{\lambda_{max}}$ (M$^{-1}$ cm$^{-1}$) | $\varepsilon_m^{280}$ (M$^{-1}$ cm$^{-1}$) |
|---|---|---|---|
| Tryptophan | 280 | 5500 | 5600 |
| Tyrosine | 275 | 1490 | 1400 |
| Phenylalanine | 258 | 200 | $\sim$0 |
| Disulfide bridge | | | 134 |

used for quantitative determinations. According to the Lambert–Beer law, the absorbance of a protein solution at a given wavelength, defined as

$$A = -\log \frac{I}{I_0} \tag{1.3}$$

is linearly related to the molar concentration of the protein, C, by the following equation:

$$A = \varepsilon_m l C \tag{1.4}$$

where $I_0$ is the incident light, $I$ is the light transmitted through the solution, $l$ is the length of the light path through the solution, and $\varepsilon_m$ is the specific molar absorbance or extinction coefficient. An analogous expression can be written, of course, in terms of a mass-based extinction coefficient and of the mass concentration of the protein using the protein molecular mass. The validity of Eq. (1.4) is generally limited to relatively dilute solutions and short light paths, for which $A$ is less than 2. At higher values of $A$, the ratio of transmitted and incident light becomes too small to permit a precise determination. Thus, quantitative determinations for concentrated protein solutions require dilution or very short light paths.

As can be seen in Table 1.6, the specific absorbance of typical proteins varies with the relative content of the aromatic amino acids Trp and Try and, to a lesser extent, of the disulfide bridges (Cys). Because the relative content is different for different proteins, an empirical determination is needed for exact quantitative determinations. Alternatively, the molar absorption coefficient can be estimated with relative accuracy as the linear combination of the individual contributions of the Trp and Tyr residues and of the disulfide bridges according to the following equation:

$$\varepsilon_m^{280}(\text{M}^{-1}\,\text{cm}^{-1}) = 5500 \times n_{\text{Trp}} + 1490 \times n_{\text{Tyr}} + 125 \times n_{\text{SS}} \tag{1.5}$$

where $n_{\text{Trp}}$, $n_{\text{Tyr}}$, and $n_{\text{SS}}$ are the numbers of its Trp, Tyr residues, and Cys disulfide bonds, respectively. As noted below, nucleic acids have an absorbance maximum at 260 nm and can interfere substantially with protein determinations at 280 nm. Thus, when nucleic acids are simultaneously present in solution, corrections must be made in order to determine protein concentration from absorbance values at 280 nm

The peptide bonds of proteins also absorb light in the "far-UV" range (180–230 nm) and very high absorbance values are observed in this region even for very dilute conditions. As a result, analytical chromatography detection is often at around 215 nm, where absorbance is about 100 times greater. Proteins with additional chromophores either absorb in the near UV range or in the visible wavelength range. Typical examples are the iron-containing proteins such as cytochrome c, hemoglobin, myoglobin, and transferrin, which have a red color, or Cu–Zn superoxide dismutase, which is green.

*Nucleic acids* show a strong UV absorbance at wavelengths in the region 240–275 nm, which originates from the $\pi$–$\pi^*$ transitions of the pyrimidine and purine nucleoside rings. Polymeric DNA and RNA show absorbance over a broad range with a maximum near 260 nm. The specific mass extinction coefficient

**Table 1.6** Representative values of the mass-based and molar extinction coefficients of proteins at 280 nm.

| Protein | Molecular mass | Number of amino acids Trp-Tyr-Cys | Mass-based extinction coefficient, $\varepsilon_M^{280}$ (ml/(mg cm)) | Molar extinction coefficient, $\varepsilon_m^{280}$ (M$^{-1}$ cm$^{-1}$) |
|---|---|---|---|---|
| Immunoglobulin G[a] | 155 000 | Varies with subclass and individual antibody | ~1.4 | ~217 000 |
| α-Chymotrypsinogen | 50 600 | 8-4-5 | 2.0 | 50 600 |
| Lysozyme (hen egg white) | 14 314 | 6-3-4 | 2.73 | 37 900 |
| β-Lactoglobulin | 18 285 | 2-4-2 | 0.95 | 17 400 |
| Ovalbumin (chicken) | 42 750 | 3-10-1 | 0.74 | 32 000 |
| Bovine serum albumin | 66 269 | 2-20-17 | 0.68 | 45 000 |
| Human serum albumin | 66 470 | 1-18-17 | 0.58 | 39 800 |

a)  May vary with recombinant IgG, when variable domain contains an excess of aromatic amino acids.
Source: Molar extinction coefficients from Mach et al. 1992 [16].

of DNA $\varepsilon_M^{260}$ is 20 (ml/(mg cm)). The purity of DNA is estimated by the ratio of absorbance at 260 and 280 nm. For pure double-stranded DNA and RNA, the ratio $\varepsilon^{260}/\varepsilon^{280}$ is between 1.8 and 2.0, respectively. The measurements are more reliable at alkaline pH. In contrast to proteins, the absorbance of nucleic acids is fairly sensitive to pH and decreases at lower pH values [17].

### 1.2.2.2 Size

Solutions and suspensions found in downstream processing of biotechnology products contain molecules and particles with a broad range of sizes as illustrated in Table 1.7. Globular proteins are in the range of 3–20 nm, while nucleic acids can be much larger. Therapeutic plasmids are in the range of 100 nm. Virus and virus-like particles are in the range from 20 to 500 nm, while cells are in the micrometer range.
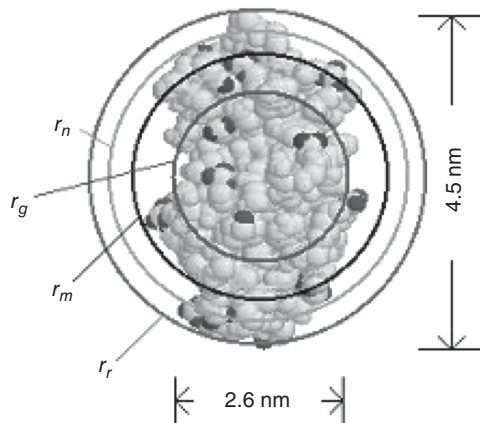
Although cells and cell debris can be separated by centrifugation due to their relatively high sedimentation velocity (Table 1.7), proteins and nucleic acids need different methods such as chromatography and membrane filtration. Separation of proteins by ultracentrifugation is only performed for analytical purposes because extremely high rotation rates (as high as 50 000 rpm) are needed.

The sizes given in Table 1.7 are for folded globular proteins. In this state, native protein structures are quite dense (typical mass density ~ 1.4 g/cm$^3$) and shaped as spheroids or ellipsoids. However, denatured, fibrous, rod, or disk-shaped proteins deviate from these compact shapes. In these cases, the size of proteins and other macromolecules is often described by other parameters, which include

**Table 1.7** Categories of bioproducts and their sizes.

| Category | Example | $M_r$ (Da) | Diameter | Sedimentation velocity (cm/h) |
|---|---|---|---|---|
| Small molecules | Amino acids | 60–200 | 0.5 nm | |
| | Sugars | 200–600 | 0.5 nm | |
| | Antibiotics | 300–1000 | 0.5–1.0 nm | |
| Macro molecules | Proteins | $10^3$–$10^6$ | 3–20 nm | $<10^{-6}$ |
| | Nucleic acids | $10^3$–$10^{10}$ | 2–1000 nm | |
| Particles | Virus | | 20–500 nm | $<10^{-3}$ |
| | Bacteria | | 1 μm | 0.02 |
| | Yeast cells | | 4 μm | 0.4 |
| | Animal cells | | 10 μm | 2 |

**Figure 1.15** Lysozyme with a size of 2.6 nm × 4.5 nm is an ellipsoid-shaped molecule. The molecular mass is 14.7 kDa and its mass density is 1.37 g/cm³. $r_m = 1.6$ nm is the equivalent radius of a sphere with the same mass and particle-specific volume as lysozyme. $r_r = 2.3$ nm is the radius established by rotating the protein about its geometric center, $r_g = 1.4$ nm is the radius of gyration, and $r_h = 2.0$ nm is the hydrodynamic radius.



the radius of gyration, $r_g$, the hydrodynamic radius, $r_h$, the radius established by rotating the protein about its geometric center, $r_r$, and the radius equivalent to a sphere with the same mass and density as the actual molecule, $r_m$. Figure 1.15 illustrates these four different quantities for lysozyme. Note that the first two, $r_h$ and $r_g$, can be obtained from direct biophysical measurements, while the last two, $r_m$ and $r_r$, can only be inferred from a knowledge of the actual protein structure.

The radius of gyration can be measured by static light scattering. This is often done with a light-scattering detector placed in-line with an SEC column allowing measurements for protein mixtures. A general relationship exists between the radius of gyration and the amount of light scattered, which is directly proportional to the product of the weight-average molar mass and the protein concentration. Accordingly,

$$\frac{k \cdot c}{R(\theta)} = \frac{1}{M_w P(\theta)} + 2A_2 c \tag{1.6}$$

where $R(\theta)$ is the excess intensity of scattered light at a certain angle ($\theta$), $c$ is the sample concentration, $M_w$ is the weight-average molar mass, $A_2$ is the second viral coefficient, and $k$ is an optical parameter equal to $4\pi n^2(dn/dc)^2/(\lambda_0^4 N_A)$. $n$ is the solvent refractive index and $dn/dc$ is the refractive index increment, $N_A$ is the Avogadro's number, and $\lambda_0$ is the wavelength of scattered light in vacuum. The function $P(\theta)$ describes the angular dependence of scattered light. The expansion of $1/P(\theta)$ to first order gives:

$$\frac{1}{P(\theta)} = 1 + \left(\frac{16\pi^2}{3\lambda_0^2}\right) r_g^2 \cdot \sin^2\left(\frac{\theta}{2}\right) + f_4 \sin^4\left(\frac{\theta}{2}\right) + \cdots \tag{1.7}$$

At low angles, the angular dependence of light scattering depends only on the mean square radius $r_g^2$ (alternatively called radius of gyration) and is independent of molecular conformation or branching.

The hydrodynamic radius is related to the protein translational diffusion coefficient or diffusivity, $D_0$, through the Stokes–Einstein equation:

$$r_h = \frac{k_b T}{6\pi\eta D_0} \tag{1.8}$$

where $k_b$ is the Boltzmann constant, $T$ is the absolute temperature, and $\eta$ is the solution viscosity. Accordingly, $r_h$ represents the radius of a sphere with the same diffusion coefficient as the actual protein. Tanford [18] has shown that the hydrodynamic radius of a globular protein can be related to its molecular mass, $M_r$, by a simple relationship. For practical calculations, the following equation provides reasonable values:

$$r_h \approx 0.081 \times (M_r)^{\frac{1}{3}} \tag{1.9}$$

where $r_h$ is in nm.

The diffusion coefficient, $D_0$, can be conveniently determined by *dynamic light scattering* (DLS), typically also in conjunction with SEC for the case of mixtures. DLS is based on the fluctuations or Brownian motion of a molecule, which, in turn, cause fluctuations in the intensity of scattered light that are autocorrelated. For a monodispersed sample, the autocorrelation function $G(\tau)$ can be expressed by a single exponential term that allows the determination of $D_0$ from the following equation:
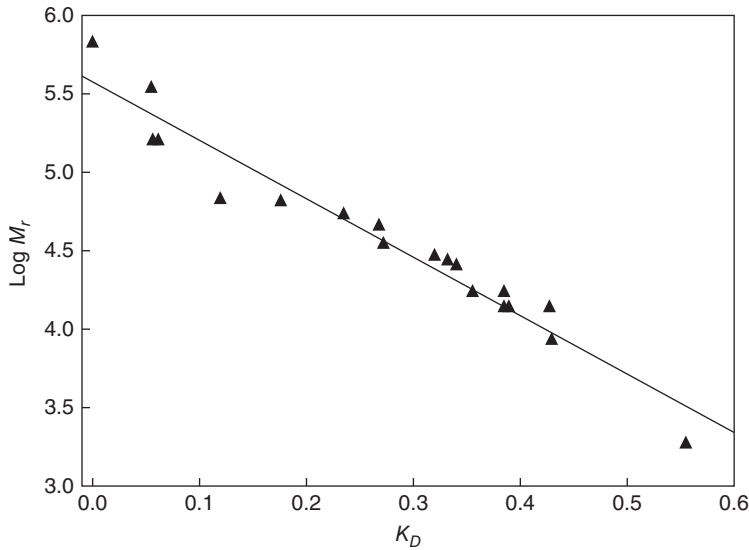
$$G(\tau) = A_1 + A_2\, e^{-2D_0 q^2 \tau} \tag{1.10}$$

with $q = (4\pi n/\lambda_0)\sin(\theta/2)$, where $\theta$ is the angle at which the detector is located relative to the sample cell.

SEC is an alternative, commonly used method for the determination of protein size. In SEC, molecules with different size penetrate differently into the pores of an SEC matrix and thus will be differently retained in the column. The SEC column can be calibrated with protein standards of known molecular mass allowing the size of an unknown protein to be estimated from its retention.

Figure 1.16 shows an example. In this case, the distribution coefficient, $K_D$, is defined as follows:

$$K_D = \frac{V_R - V_0}{V_t - V_0} \tag{1.11}$$

**Figure 1.16** Calibration of a size exclusion column (Superdex 75, GE Healthcare, Uppsala, Sweden) with a set of reference proteins (molecular mass, $M_r$, in parenthesis): thyroglobulin (669 000), fibrinogen (340 000), glucose oxidase (160 000), IgG (160 000), bovine serum albumin (66 430), hemoglobin, (64 500), triosephosphate isomerase (53 200), ovalbumin (45 000), lectin (35 000), carbonic anhydrase (29 000), subtilisin (27 000), chymotrypsinogen (25 000), myoglobin (17 000), calmodulin (16 800), ribonuclease A (13 700), ribonuclease S (13 700), cytochrome c (13 600), ubiquitin (8600), and pep6His (1839).

where $V_R$ is the retention volume, $V_t$ is the total column volume, and $V_0$ is the extraparticle void volume. The latter is determined empirically from the retention of a compound large enough to be completely excluded from the pores of the chromatography matrix. Blue Dextran, a 2000 kDa molecular mass dextran labeled with a blue dye, is often used for this purpose. When the logarithm of molecular mass of standard proteins is plotted vs. $K_D$, an almost linear relationship is obtained (see Figure 1.16). Other methods for the determination of protein molecular size are SDS-PAGE, which provides information on molecular mass, ultracentrifugation, which provides information on hydrodynamic radius, and other scattering methods, such as small-angle X-ray scattering (SAXS).

Mass spectrometry has become a routine method to determine the exact mass of proteins and/or their identity. The main methods for ionization of proteins in mass spectrometry are electrospray ionization (ESI) and matrix-assisted laser desorption/ionization (MALDI). Whole-protein mass analysis is determined primarily using either time-of-flight (TOF) MS or Fourier transform ion cyclotron resonance (FT-ICR). Mass spectrometry can also be combined with SDS gel electrophoresis to identify proteins. In this case, protein bands are excised from the electrophoresis gel and subjected to mass spectrometry. Proteins are either identified via mass search in a database or by the fragment pattern obtained during ionization. It is also possible to enzymatically degrade the protein and subject the fragments to mass spectrometry. From the enzymatic cleavage patterns, the

nature of the protein can be identified via search in databases. Protein mass spectrometry can also be performed as a semiquantitative method.

### 1.2.2.3 Charge

Proteins are amphoteric molecules with both negative and positive charges, which stem from the R-groups of acidic and basic amino acids (see Table 1.1) and from the amino and carboxyl terminus of each polypeptide chain. The latter have $pK_a$ values around 8.0 and 3.1, respectively. Certain post-translational modifications can substantially alter the charge of a protein. Important examples are glycosylation with sialic acid, occurring, for example, at N-glycosylation sites in antibodies or erythropoietin, and deamidation of asparagine and glutamine residues. Both of these post-translational modifications make proteins more acidic and, thus, more highly negatively charged. In many cases, they also affect the in vivo half-life of the protein and/or its biological function, so that their control can be an important goal of downstream processing.

The net charge of a protein depends on the number of ionizable amino acid residues and their $pK_a$ values. The protonation of these residues changes with pH according to the following equations for acidic and basic residues, respectively:

$$K_a = \frac{[R^-][H^+]}{[RH]} \tag{1.12}$$
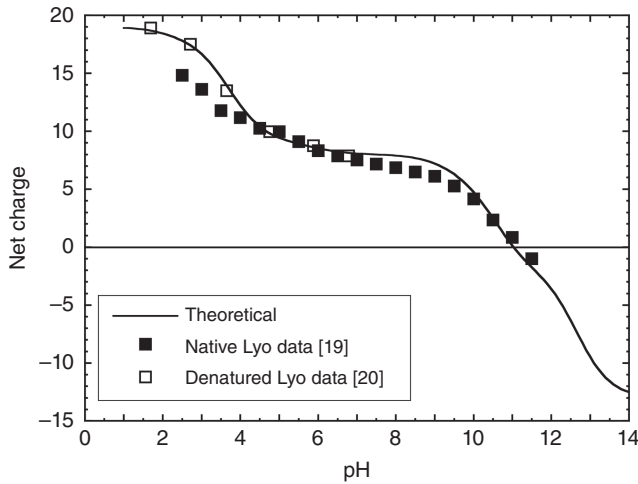
$$K_a = \frac{[R][H^+]}{[RH^+]} \tag{1.13}$$

where the brackets indicate thermodynamic activities. In logarithmic form, we obtain:

$$\log \frac{[R^-]}{[RH]} = pH - pK_a \tag{1.14}$$

$$\log \frac{[R]}{[RH^+]} = pH - pK_a \tag{1.15}$$

where p stands for $-\log_{10}$. From these equations, it is obvious that acidic residues are essentially completely deprotonated and, thus, negatively charged at pH values that are two units higher than their $pK_a$. Conversely, basic amino acids are essentially completely protonated and, thus, positively charged, at pH values that are two units below their $pK_a$. Based on the $pK_a$ values shown in Table 1.1, we can see that, in practice, at the near-neutral pH values typically encountered in bioprocessing, all acidic residues in protein are negatively charged while all basic residues are positively charged. Histidine, however, is an exception to this rule. Its $pK_a$ is near neutral; thus, under typical processing conditions, this residue will be charged to an extent that depends on the exact value of pH.

At a particular pH, called the isoelectric point or pI, the protein net charge becomes zero with an exact balance of positively and negatively charged residues. The net charge and the theoretical isoelectric point of a protein can be calculated from Eqs. (1.12) and (1.13) knowing the $pK_a$ values of the R-groups and the amino acids in the primary sequence. Figure 1.17 shows an example for lysozyme. The calculation shown in Figure 1.17 is only approximate, however, because activity coefficients were neglected and the $pK_a$-values were assumed to be equal to
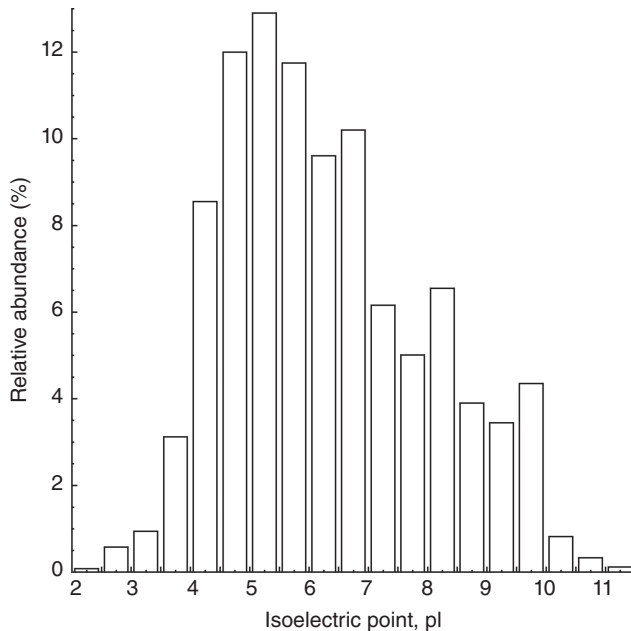
**Figure 1.17** Net charge of lysozyme in denatured and native forms compared to the theoretical calculation based on the amino acid sequence. The pI is around 11. Source: Adapted from Kuehner et al. 1999 [19] and Tanford and Roxby 1972 [20].

those of the free amino acids. This is likely incorrect because the microenvironment where the individual residues are actually found in the protein structure has a significant effect. Nevertheless, the agreement between the theoretical net charge and that determined experimentally as a function of pH is remarkable. The more significant deviations, in this case, appear to occur for native lysozyme at low pH but largely disappear when the protein is denatured, suggesting that the discrepancy may arise because some of the acidic residues may be partially buried in the folded structure. As can be seen in Figure 1.17, the pI of lysozyme is around 11, which is in agreement with IEF measurements. Around neutral pH values, this protein has a high net positive charge with a plateau region where the charge is only slightly affected by pH. Such conditions would be conducive, for example, to a robust adsorption process for the capture of lysozyme with a cation exchanger.

Figure 1.18 shows a histogram illustrating the distribution of the pI values of many common proteins. As can be seen in this graph, a majority of the proteins have a slightly acidic isoelectric point and this is indeed found for the proteins present in many microorganisms such as *E. coli*. As a result, it is often easier to purify alkaline proteins that can be adsorbed onto a cation exchange column because most of the host cell proteins are likely to pass unretained through these columns. Many monoclonal antibodies also have high isoelectric points, allowing the development of effective capture and purification processes with cation exchangers.

A final important consideration with regard to protein charge is the spatial distribution of the charged residues. Figure 1.19 provides an example illustrating the location of positive and negative charges on the surface of lysozyme and on that of human serum albumin at neutral pH. A consequence of this heterogeneous spatial distribution is that, frequently, the net charge of the protein is not sufficient
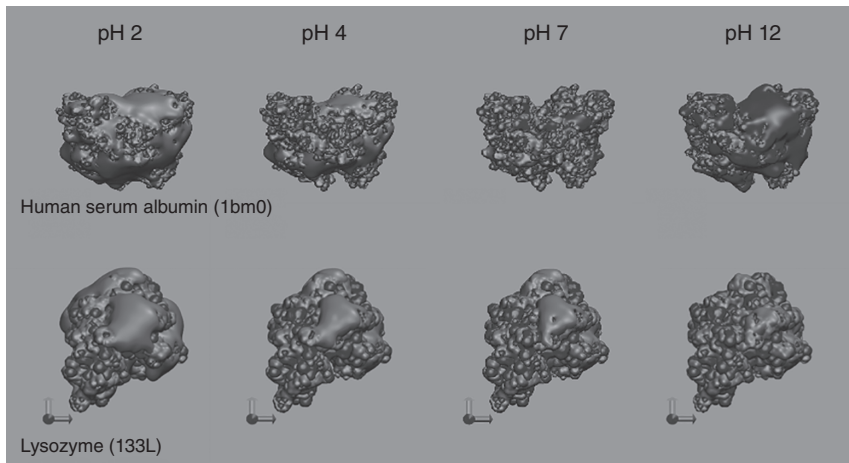
**Figure 1.18** Distribution of isoelectric points of common proteins. Source: Adapted from Righetti and Caravaggio 1976 [21].

to determine whether the protein will bind or not to an oppositely charged surface. For example, as a result of the localized concentration of positively charged residues, it is possible for a protein to bind to a cation exchanger at pH values well above the protein pI, where the net charge is highly negative, or, vice versa, for a protein to bind to an anion exchanger at pH values well below the protein pI, because of localized negatively charged residues.

### 1.2.2.4 Hydrophobicity

The surface hydrophobicity of a protein is determined by the R-groups of its surface-exposed nonpolar amino acids. Although the term hydrophobicity is commonly used, a precise definition is difficult and is extensively debated. The transfer of an apolar compound into a polar liquid, such as water, is associated with heat and quantified as free energy. The hydrophobic effect is strongest when entropic effects are dominant. Hydrophobic effects increase with the surface tension of the solution, which is caused by the attraction between the liquid's molecules by various intermolecular forces. Hydrophobic effects are thus dominated by the strong hydrogen bonds of water, while van der Waals forces generally play a minor role.

A practical approach to measure the hydrophobicity of a protein is by measuring its retention in a chromatography column packed with a hydrophobic matrix. In this case, if unfolding does not occur, retention of the protein is related to the number and hydrophobic amino acid R-groups exposed at the protein surface [22]. The hydrophobicity obtained by this method is relative and depends on the applied methodology, so it is useful only for ranking purposes.
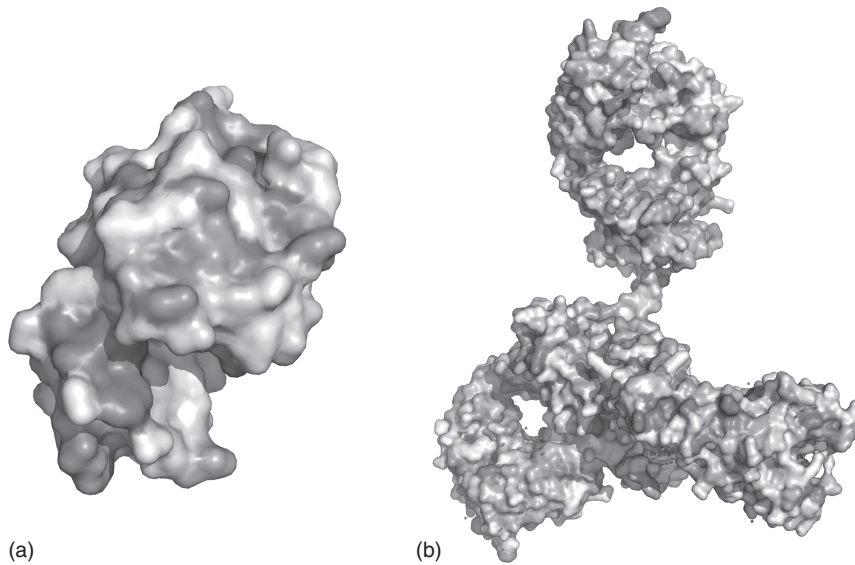
**Figure 1.19** Distribution of positively charged (red) and negatively charged (blue) residues on the surface of lysozyme and of human serum albumin. (*See color plate section for color representation of this figure*).
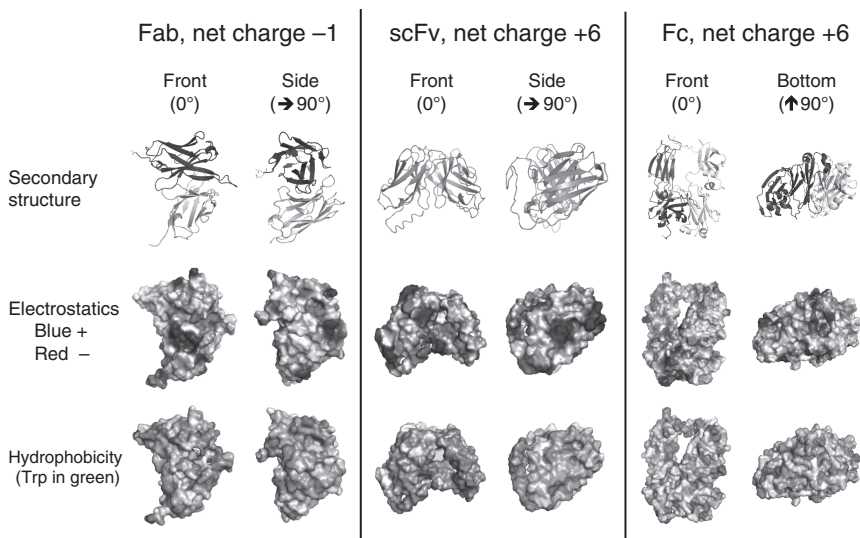
The hydrophobicity of proteins can also be theoretically calculated from the transfer energy of amino acids from an apolar solvent into water as shown by Hopp and Woods [23] and Kyte and Doolittle [24]. However, the free energy of transfer of amino acids does not totally reflect the hydrophobicity of a protein and strongly depends on the hydrophobicity scales that are used for calculating the protein's hydrophobicity.

Analogously to charge, the distribution of surface-exposed hydrophobic residues in proteins is also not homogenous. This is illustrated for lysozyme and for a recombinant antibody in Figure 1.20. The density and distribution of hydrophobic groups at the protein surface provide the basis for hydrophobic interaction chromatography (HIC) where surface hydrophobic residues interact with a mildly hydrophobic matrix. Because incorrect folding may lead to variations in the number of surface-exposed hydrophobic residues, HIC may be used as a tool to separate native proteins from misfolded isoforms.

Figure 1.21 illustrates the distribution of both charged and hydrophobic areas across the three major domains on a bivalent, bispecific antibody from Ref. [25]. In this case, the molecule consists of a framework antibody (IgG) with the same structure as that in Figure 1.20 (b) but with two single-chain fragment variable (scFv) domains linked to the heavy chains of the antibody Fab domains through flexible peptide linkers. As can be seen in Figure 1.21, the different domains (Fab, scFv, and Fc) have very different net charge as well as different heterogeneously distributed charged and hydrophobic surface areas. Tryptophan residues, also present to a different extent in the different domains, exhibit a degree of surface exposure that varies depending on the protein folding state. These differences in charge and hydrophobicity across the different domains of these more complex proteins affect how these molecules interact with each other and with chromatographic surfaces. As shown, for example, in Refs. [25, 26], multiple binding states
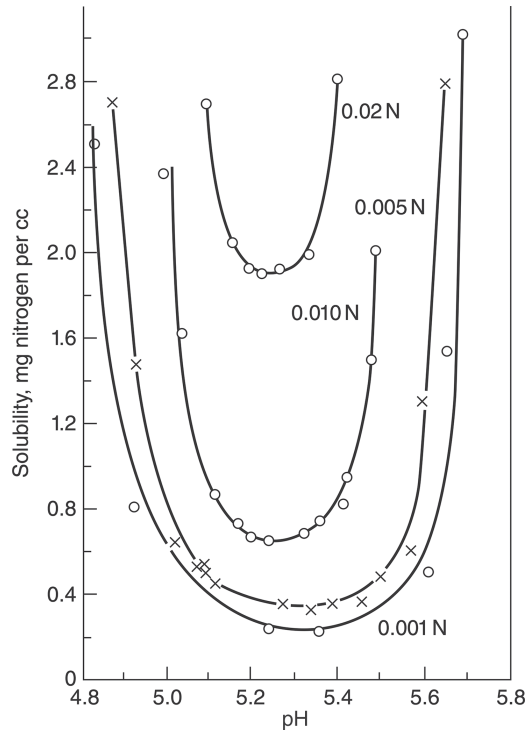
(a)                                        (b)

**Figure 1.20** Distribution of surface hydrophobicity for lysozyme (a) and for a recombinant antibody (b). More hydrophobic areas are colored in deeper shades of red. (*See color plate section for color representation of this figure*).



**Figure 1.21** Secondary structure, surface charge, and surface hydrophobicity of the three major domains (Fab, scFv, and full Fc) of a bivalent, bispecific antibody at pH 7. In the actual molecule, a scFv domain is attached to the heavy chain of each Fab domain via a flexible peptide linker. The net charge is calculated from the amino acid sequence. In the surface hydrophobicity maps, more hydrophobic areas are colored in deeper shades of red. Trp residues are shown in green. Two different views are shown for each domain rotated 90°. Source: Adapted from Kimerer et al. 2019 [25]. (*See color plate section for color representation of this figure*).

**Figure 1.22** Solubility of β-lactoglobulin (measured in terms of amount of protein nitrogen dissolved per cubic centimeter of solution) as a function of pH at four different concentrations of sodium chloride. Source: Fox and Foster 1957 [27]. Reproduced with permission of John Wiley & Sons.

can be observed as a result of interdomain association and interactions between the different domains and the chromatographic surface.
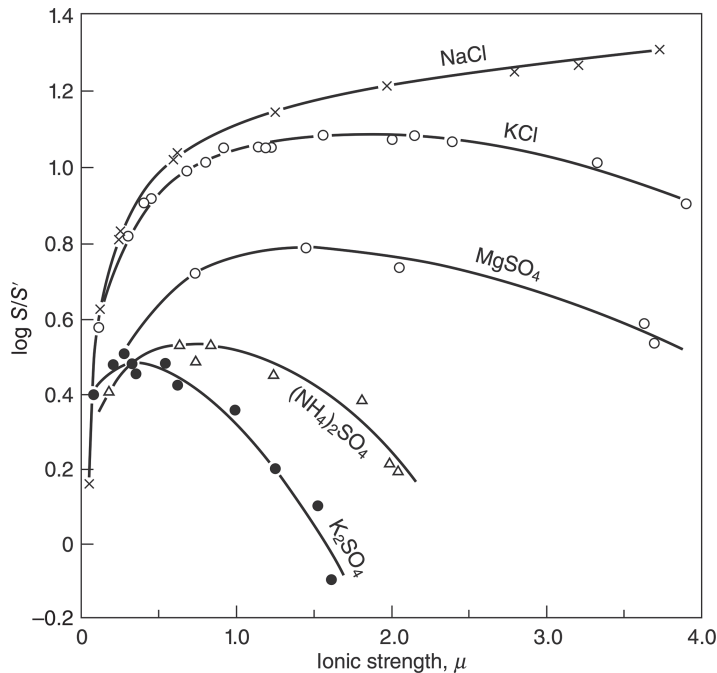
### 1.2.2.5 Solubility

Protein solubility is often a critical consideration in downstream processing because it can vary dramatically with pH, ionic strength, and salt type. Prediction of the solubility in aqueous media from protein structure is difficult and empirical measurements are usually needed. The range of protein solubilities is very broad. Some proteins, e.g. Cu–Zn superoxide dismutase, have solubility as high as 400 mg/ml, while others, e.g. recombinant interferon-$\gamma$, have solubility of less than 10 mg/ml. In general, protein solubility is lowest at the isoelectric point, where the net charge is zero, but varies with ionic strength, which is defined as follows:

$$I = \frac{1}{2} \sum_{j=1}^{n} c_j z_j^2 \tag{1.16}$$

where $c_j$ is the concentration of ion $j$ and $z_j$ is its charge. The solubility of β-lactoglobulin as a function of salt concentration and pH is shown as an example in Figure 1.22.

In general, in the range of very low salt concentrations, adding salt in most cases increases the protein solubility as can be seen in Figure 1.22. On the other hand, in the range of high salt concentrations, e.g. >0.4 M, the addition of salt can either increase or decrease the protein solubility depending on the nature

**Figure 1.23** Solubility of carboxyhemoglobin in aqueous solution with different electrolytes at 25 °C. $S$ and $S'$ are used in lieu of $w$ and $w_0$. Source: Data from Green 1932 [28].

of the salt added. This behavior can be seen, for example, in Figure 1.23 for carboxyhemoglobin. Adding NaCl increases the solubility, which is referred to as "salting in," but adding potassium or ammonium sulfate at concentrations above about 0.4 M depresses the solubility, which is referred to as "salting out." As can be seen in Figure 1.23, for this protein, magnesium sulfate results in salting in at low ionic strengths and in salting out at high ionic strength. In any case, it should be recognized that these effects are quantitatively and sometimes even qualitatively different for different proteins.

Protein solubility trends can be described by the extended form of the Debye–Hückel theory. Accordingly, we have:

$$\log \frac{w}{w_0} = \frac{0.5 \cdot z_1 \cdot z_2 \sqrt{I}}{1 + A\sqrt{I}} - \kappa_s I \qquad (1.17)$$
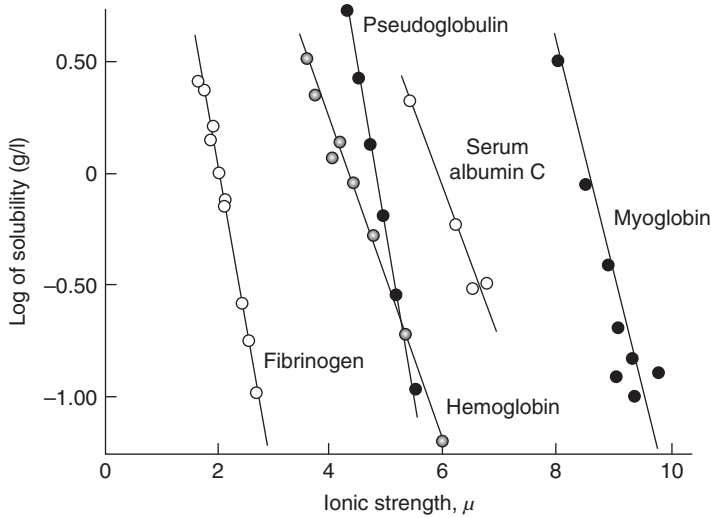
where $w$ is the protein solubility in the actual solution, $w_0$ is the solubility of the protein in water, $z_1$ and $z_2$ are the salt charges, and $\kappa_s$ and $A$ are salt- and protein-specific empirical parameters. At high ionic strengths, Eq. (1.17) reduces to the following log–linear relationship:

$$\log \frac{w}{w_0} = \beta - \kappa_s I \qquad (1.18)$$

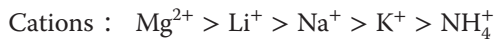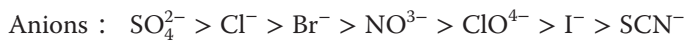which, as seen in Figure 1.24, is observed for many different proteins.

The effect of the salt type on protein solubility has been formally described for the first time by Hofmeister [30] who ranked the anions and cations according

**Figure 1.24** Solubility of various proteins in ammonium sulfate solutions. Source: Cohn and Edsall 1943 [29]. Reproduced with permission of American Chemical Society.

to their ability to precipitate proteins. This ranking is described by the following *Hofmeister series* or *lyotropic series*:

$$\text{Anions}: \quad SO_4^{2-} > Cl^- > Br^- > NO^{3-} > ClO^{4-} > I^- > SCN^-$$

$$\text{Cations}: \quad Mg^{2+} > Li^+ > Na^+ > K^+ > NH_4^+$$

A simple interpretation of this series is that certain ions bind water molecules tightly, thereby decreasing the ability of the protein to stay in solution. Although both the cation and the anion in a given salt are important to this property, the contribution of the anion is usually dominant. Interestingly, the ions in the Hofmeister series also correlate with the so-called Jones–Dole B-coefficient and the entropy of hydration so that both appear to be related to the effects of salts on the structure of water. Finally, it should be noted that, in practice, the selection of salts for use in downstream processing depends not only on the Hofmeister series but also on factors such as price, availability, biocompatibility, and disposal costs.

### 1.2.2.6 Chemical Stability

Two different types of chemical stabilities should be considered for proteins: the *conformational stability* (or *thermodynamic stability*) and the *kinetic stability* (or *colloidal stability*). The conformational stability of a protein is described by the free energy $\Delta G$ of the equilibrium between native and the unfolded states. The transition of the native folded form, $N$, into the unfolded form, $U$, is described by the following quasi-chemical reaction:

$$N \underset{k_{-1}}{\overset{k_1}{\rightleftharpoons}} U \tag{1.19}$$

**Table 1.8** Thermodynamic stabilities of proteins.

| Protein | Conditions | Free energy of the unfolding reaction $\Delta G$ (kcal/mol) | Melting temperature (°C) |
|---|---|---|---|
| Horse cytochrome c at pH 6 at 25 °C | 0 M urea | 31.3 | n.a. |
| | 2 M urea | 22.3 | n.a. |
| | 4 M urea | 14.2 | n.a. |
| | 6 M urea | 3.2 | n.a. |
| Hen egg white lysozyme at pH 3.0 | 24 °C | 41.0 | n.a. |
| | 40 °C | 30.4 | n.a. |
| | 55 °C | 14.7 | n.a. |
| | 75 °C | −5.9 | n.a. |
| Bovine chymotrypsinogen at melting temperature and pH 2.0 | 0% glycerol | 0.015 | 42.9 |
| | 20% v/v glycerol | 0.146 | 44.9 |
| | 40% v/v glycerol | 0.235 | 46.2 |

Source: Data for chymotrypsinogen are from Gekko and Timasheff 1981 [31].

where $k_1$ and $k_{-1}$ are rate constants. The corresponding equilibrium constant $K_{eq} = [U]/[N]$ is usually very low in aqueous solution as protein folding is generally thermodynamically favored as a result of the concentration of the hydrophobic residues in the protein core. The corresponding $\Delta G$ is given by the following equation:

$$\Delta G = -RT \ \ln K_{eq} \tag{1.20}$$

Representative values of $\Delta G$ are given in Table 1.8 along with the corresponding "*melting temperature*," which is defined as the temperature at which half of the protein is in the unfolded state. *Kosmotropic* (or cosmotropic) *salts* and polyols such as sorbitol or sucrose stabilize proteins while *chaotropic salts* or urea at higher concentrations have a destabilizing effect on protein conformation.

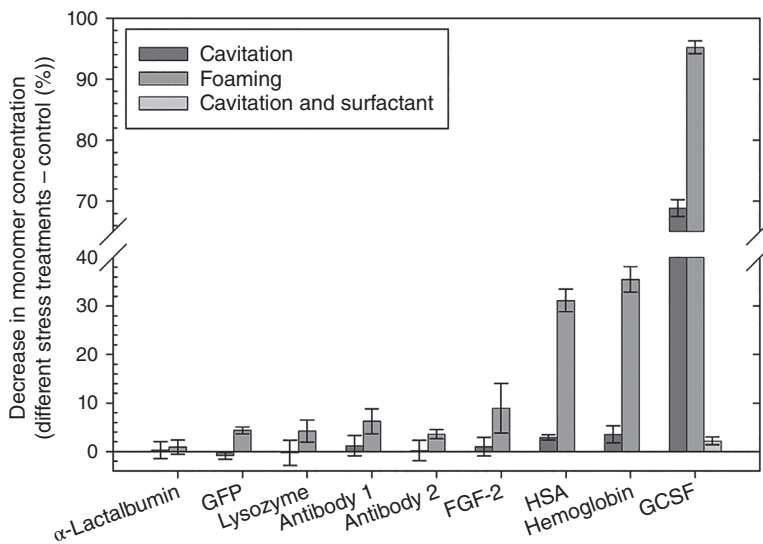Kinetic stability, on the other hand, can be described by the following equation:

$$N \underset{k_{-1}}{\overset{k_1}{\rightleftarrows}} U \xrightarrow{k_2} A \tag{1.21}$$

which shows a further kinetically driven step from the unfolded state to an irreversibly aggregated state $A$. Proteins with a high $k_2$ exhibit low kinetic stability. The overall stability thus depends on both thermodynamic and kinetic effects. It is possible, for example, for an added salt to decrease kinetic stability while enhancing overall stability as a result of thermodynamic effects. This is often difficult to predict, however, so that, in practice, overall stability and shelf life are measured empirically [32].

### 1.2.2.7 Mechanical Stability

It is a frequent misconception that proteins are not mechanically stable because of the contradictory reports in the literature. Proteins are highly mechanically

**Table 1.9** Typical shear rates encountered in bioprocessing operations.

| Operation | Shear rate, $\dot{\gamma}$ (s$^{-1}$) |
|---|---|
| Expanded beds | $<10$ |
| Packed beds | $<10^3$ |
| Stirred tanks | $10^2$–$10^3$ |
| High pressure homogenizers | $10^6$ |



**Figure 1.25** Comparison of protein aggregation from cavitation and foaming for nine different proteins. Data points represent the means of three experiments minus control ± standard deviation. Addition of surfactant decreases the effect of cavitation dramatically. Source: Adapted from Duerkop et al. 2018 [8].

stable. They resist shear rates up to $10^8$ s$^{-1}$. This is 10 times higher than the shear rates obtained in a high-pressure homogenizer (see Table 1.9). Also, high protein concentrations suppress cavitation. Thus, it is almost impossible to unfold an average protein by a mechanical force under ordinary process conditions. Proteins may be sensitive, however, to unfolding at air–liquid interfaces. Therefore, it is highly recommended to avoid foaming when working with protein solutions. In addition, at very high velocity, cavitation may occur, where gas bubbles are generated. This in turn may lead to protein unfolding or aggregation (see Figure 1.25).

#### 1.2.2.8 Viscosity

Many of the solutions and suspensions encountered in bioprocessing are highly viscous. This is especially true for fermentation broths that contain DNA and for highly concentrated protein solutions. In general, viscosity, $\eta$, is related to the

**Table 1.10** Apparent viscosities of various fluids at 20 °C.

| Liquid | Apparent viscosity (mPa s) |
|---|---|
| Water | 1 |
| Glycerol | 1070 |
| Ethanol | 1.20 |
| Acetonitrile | 0.34 |
| Clarified cell culture supernatant | ~5 |
| Blood | 10 |
| *E. coli* homogenate | ~40 |
| *E. coli* broth | ~20 |
| *Penicillium chrysogenum* fermentation broth | 40 000 |
| Heinz Ketchup | 50 000–70 000 |

shear stress, $\tau$, and the shear rate, $\dot{\gamma}$, by the following equation:

$$\tau = \eta \times \dot{\gamma} \tag{1.22}$$

For Newtonian fluids, $\eta$ is a constant and the relationship between shear stress and shear rate is linear. For non-Newtonian fluids, however, $\eta$ varies with shear rate and the relationship becomes nonlinear. For example, the behavior of pseudoplastic fluids is described by the following equation:
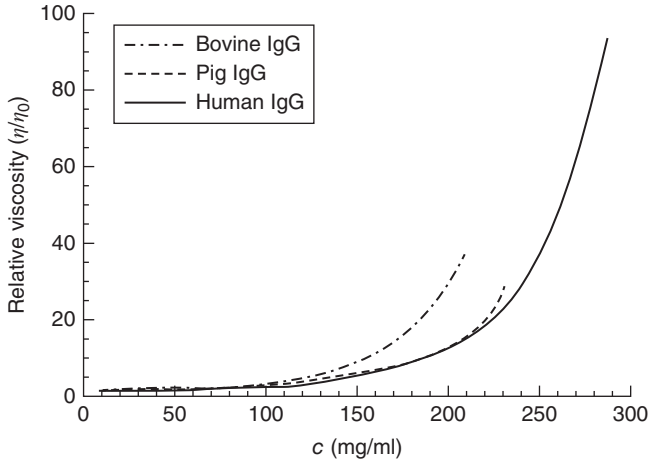
$$\tau = K \times (\dot{\gamma})^n \tag{1.23}$$

where $K$ and $n$ are called the *flow consistency index* and *flow behavior index*, respectively. For highly concentrated protein solutions and for many culture supernatants, $n$ is smaller than unity, indicating that the apparent viscosity, $\eta = \tau/\dot{\gamma}$, decreases with increasing shear rate. The ranges of various shear rates encountered in bioprocessing are shown in Table 1.9.

The viscosities of some typical fluids as well those of some of the solutions encountered in bioprocessing are shown in Table 1.10. In general, cell culture supernatants have viscosities lower than 10 mPa s, while cell homogenates are much more viscous with values of $\eta$ up to 40 mPa s. DNA is, usually, the greatest contributor to the viscosity of raw fermentation broths. Fortunately, however, both genomic and plasmid DNA are very sensitive to shear and are often mechanically degraded early on in the downstream process. DNAse enzymes, naturally occurring or added intentionally, also help degrade these molecules, thereby reducing viscosity. This is especially important for intracellular products that require disruption of the cells with concomitant release of the intracellular components into the product-containing solution.

The *intrinsic viscosity* $[\eta]$ is a measure of a solute's contribution to the solution viscosity and is defined as follows:

$$[\eta] = \lim_{c \to 0} \frac{\eta - \eta_0}{\eta_0 c} \tag{1.24}$$

**Figure 1.26** Relative viscosity of human, bovine, and pig IgG solutions as a function of IgG concentration at room temperature. Source: Adapted from Monkos and Turczynski 1999 [33].

where $\eta_0$ is the solvent viscosity in the absence of the solute and $c$ is the solute concentration. In a semidilute limit, $\eta$ can be described as a function of $c$ by the following series expansion:

$$\frac{\eta - \eta_0}{\eta_0} = [\eta]c + k_1[\eta]^2c^2 + k_2[\eta]^3c^3 + \cdots \tag{1.25}$$

At very high protein concentrations, however, semiempirical models are needed [33]. As an example, Figure 1.26 shows the relative viscosity (defined as $\eta/\eta_0$) of IgG solutions as a function of IgG concentration. For concentrations lower than about 100 mg/ml, the relative viscosity increases approximately linearly with solute concentration conforming to the first term on the right-hand side of Eq. (1.25). However, at higher concentrations, the viscosity increases almost exponentially.

Table 1.11 provides the intrinsic viscosities of representative biomolecules. As can be surmised from these data, the intrinsic viscosity depends on the shape of the molecule. For instance, rod-shaped proteins have a higher intrinsic viscosity than globular ones. An empirical relationship between $[\eta]$ and molecular mass $M_r$ is given by the Mark–Houwink equation:

$$[\eta] = K(M_r)^a \tag{1.26}$$

where $a$ is a parameter related to the "stiffness" of the polymer chains. Theoretically, $a = 2$ for rigid rods, 1 for coils, and 0 for hard spheres. Empirically, however, values of $a = 0.6$, 0.7, and 0.5 have been found for BSA, ovalbumin, and lysozyme, respectively. Literature data (e.g. see Ref. [35]) suggest a general relationship between intrinsic viscosity and the number of amino acid residues, $n_{aa}$, which can be expressed as follows:

$$[\eta] = 0.732(n_{aa})^{0.656} \tag{1.27}$$

with $[\eta]$ in ml/g. Accordingly, larger proteins and protein aggregates have higher intrinsic viscosity than smaller proteins and monomeric forms.

**Table 1.11** Intrinsic viscosities of various biologically important macromolecules in dilute solutions.

| Shape | Solute | Molecular mass | $[\eta]$ (ml/g) |
|---|---|---|---|
| Globular | Ribonuclease | 13 680 | 3.4 |
| | Serum albumin | 67 500 | 3.7 |
| | Ribosomes (*E. coli*) | 900 000 | 8.1 |
| | Bushy stunt virus | 10 700 000 | 3.4 |
| Random coils (unfolded proteins) | Insulin (A-chain) | 2 970 | 6.1 |
| | Ribonuclease | 13 680 | 16 |
| | Serum albumin | 68 000 | 52 |
| | Myosin subunit | 197 000 | 93 |
| Rods | Fibrinogen | 330 000 | 27 |
| | Myosin | 440 000 | 217 |
| | Calf thymus DNA | 15 000 000 | >10 000 |

Source: Sibileva et al. 2001 [34]. Reproduced with permission of Springer.

### 1.2.2.9 Diffusivity

The *molecular diffusion coefficient* or *diffusivity* in solution, $D_0$, is a function of the solute size, the viscosity of the solution, and temperature. As previously noted, the Stokes–Einstein equation describes this relationship as follows:

$$\frac{D_0\eta}{T} = \frac{k_b}{6\pi r_h} \tag{1.28}$$

where $k_b$ is the Boltzmann constant and $r_h$ is the solute hydrodynamic radius. The diffusivities encountered in bioprocessing range widely from $1 \times 10^{-5}$ cm$^2$/s for salts and other small molecules to $1 \times 10^{-9}$ cm$^2$/s for large biomolecules such as DNA. Protein diffusivities in dilute aqueous solution are generally in the range $10^{-6}$–$10^{-7}$ cm$^2$/s. Table 1.12 provides a summary of typical diffusivities in dilute solutions at room temperature.
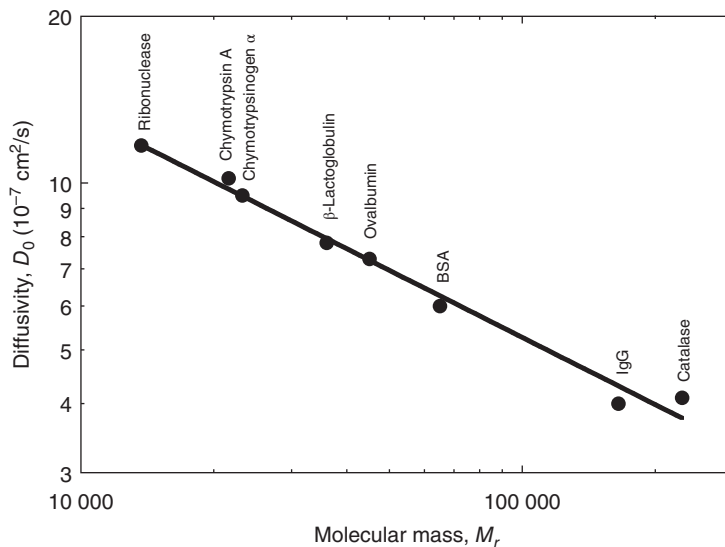
In general, protein diffusivities are 10–100 times lower than those of small molecules. Plasmids have even smaller diffusivity values – as much as 1000 times lower than for small molecules. Figure 1.27 illustrates the effect of molecular mass on the diffusivity of globular proteins in dilute aqueous solution at room temperature. Tyn and Gusek [36] provide the following correlation for such proteins:

$$\frac{D_0\eta}{T} = \frac{9.2 \times 10^{-8}}{(M_r)^{1/3}} \tag{1.29}$$

where $D_0$ is in cm$^2$/s, $\eta$ in cp, $T$ in K, and $M_r$ in Da, which has accuracy better than $\pm 10\%$. Both Eqs. (1.28) and (1.29) indicate that the group $D_0\eta/T$ is a constant for a given protein in aqueous solution. This fact can be utilized conveniently to estimate the effects of solution viscosity and temperature on the protein diffusivity. In general, $D_0$ increases with temperature in part because of the $T$-term in the denominator of this ratio but, more importantly, because the solution viscosity usually decreases as the temperature is increased. For example,

**Table 1.12** Diffusivities in dilute solution at room temperature.

| Solute | Solvent | Solvent viscosity, $\eta$ (mPa s) | Diffusivity, $D_0$ ($10^{-5}$ cm$^2$/s) |
|---|---|---|---|
| Benzoic acid | Water | 1 | 1.00 |
| Valine | Water | 1 | 0.83 |
| Sucrose | Water | 1 | 0.53 |
| Water | Ethanol | 1.1 | 1.24 |
| Water | Glycol | 20 | 0.18 |
| Water | Glycerol | >120 | 0.013 |
| Ribonuclease ($M_r = 14$ kDa) | Water | 1 | 0.120 |
| Albumin ($M_r = 65$ kDa) | Water | 1 | 0.060 |
| IgG ($M_r = 165$ kDa) | Water | 1 | 0.037 |
| pDNA ($M_r = 3234$ kDa) | Water | 1 | 0.004 |



**Figure 1.27** Diffusivity of globular proteins in dilute aqueous solution at room temperature. Source: Adapted from Tyn and Gusek 1990 [36].

going from 25 to 4 °C, the viscosity of a dilute aqueous solution increases from about 1 mPa s to about 1.5 mPa s. As result of this relationship, $D_0$ is expected to become only about 60% of the value at 25 °C. Understanding this effect is important for scale-up as biochromatography processes are often developed in the laboratory at room temperature but then scaled-up for operation at different temperatures, usually lower than room temperature in order to minimize product degradation and inhibit the potential growth of microorganisms.

Predicting the diffusivity of nonglobular proteins is more complex as their shape and not just their radius affect $D_0$. In general, for the same molecular

mass, the diffusivity of an elongated macromolecule, such as a fibrous protein, is substantially lower than that of the same molecule in globular form. Thus, unfolded proteins have smaller diffusivities than the corresponding folded ones. Plasmids, which are rod shaped in their supercoiled conformation, have much smaller diffusion coefficients than proteins with a stronger dependence on molecular mass (2/3 power) than that suggested by Eq. (1.29) (see Ref. [37]). Linear polymers have even lower $D_0$-values that are nearly inversely proportional to molecular mass. Polymer–protein conjugates, such as PEGylated proteins, have smaller diffusivities than the corresponding unconjugated protein. In this case, both the protein molecular mass and the polymer molecular mass contribute to the value of $D_0$, but the former according to the 1/3 power while the latter according to the first power of $M_r$. The following equation has been developed by Fee and van Alstine [38] to predict the hydrodynamic radius $r_h$ (in nm) of PEGylated proteins:

$$r_h = 0.082(M_{r,\text{protein}})^{1/3} + 0.373 + 0.000\,11 M_{r,\text{PEG}} \tag{1.30}$$

where $M_{r,\,\text{protein}}$ is the molecular mass of the protein and $M_{r,\,\text{PEG}}$ is the molecular mass of total conjugated PEG polymers. The value of $D_0$ for the conjugated protein is found replacing $r_h$ in Eq. (1.28) with the value of $r_h$ predicted by this equation. Examples of chromatography of PEGylated proteins including high loading conditions can be found in Ref. [39]. Finally, virus and virus-like particles exhibit diffusion coefficients consistent with Eq. (1.28) if they are spherical but substantially smaller values if they are rod shaped.

In general, if the shape of a protein or of a bioparticle can be approximated as an ellipsoid with radii $r_a$ and $r_b$, using the average radius $(r_a + r_b)/2$ in place of $r_h$ in Eq. (1.28) results in an error of less than 15% in the predicted value of $D_0$ if $r_a/r_b < 2$. Several approaches are available for the experimental determination of diffusivity and are reviewed, for example, by Cussler [40]. Commonly used approaches for proteins include DLS (described above), diffusion cells, Taylor dispersion-based methods, and microinterferometry. Such measurements may be needed for molecules that have complex or unknown shapes such as in the case of aggregates.
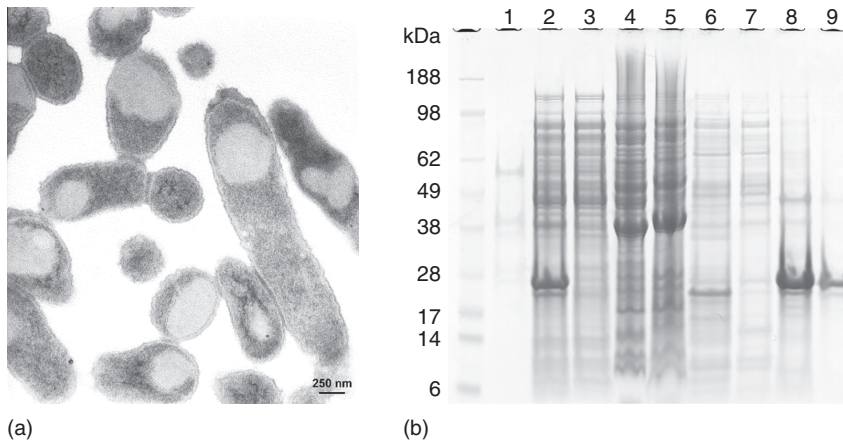
## 1.3 Bioprocesses

This section discusses commonly used expression systems and the general structure of downstream processes needed to achieve the desired product purity. Special emphasis is placed on the production of recombinant proteins by fermentation and cell culture, which play a major role in industrial biotechnology.

### 1.3.1 Expression Systems

Many different expression systems have been developed for recombinant proteins, ranging from very simple bacteria to plants and animals. However, the number of host cells actually used in the industrial production of biopharmaceutical proteins is quite limited. The most popular bacterial strain is *E. coli* BL21, which is used for the production of proteins whose biological activity does

**Figure 1.28** (a) Electron micrograph of *E. coli* cells overexpressing a recombinant protein as inclusion bodies. (b) SDS–PAGE under nonreducing conditions. Lane 1: culture supernatant; lane 2: homogenate; lane 3: soluble fraction; lanes 4–7: supernatants of wash steps; lanes 8 and 9: insoluble fraction containing the product.

not require post-translational modifications. Protein expression in *E. coli* can occur in three different ways. The protein can be secreted into the periplasm, which is the space between cell membrane and cell wall; it can be expressed in the cytoplasm as a soluble protein; or it can accumulate in the cell as inclusion bodies. Each system is effective for different proteins. The cytoplasm of *E. coli* is a strongly reducing environment that hinders the formation of disulfide bridges, whereas the periplasm offers a more oxidizing environment, which allows such a formation facilitating folding. Therefore, antibody fragments and antibody-derived molecules (the so-called "new formats") are often expressed into the periplasm. Some proteins that are toxic to the cells or that have a short half-life in the cytoplasm or periplasm life have been successfully produced as inclusion bodies. In this case, the protein forms aggregates that cannot be attacked by proteases. On the other hand, although often not fully denatured, proteins expressed as inclusion bodies are generally not in their native conformation. Thus, a refolding process is typically required to generate a fully active form. Although refolding can be costly, on balance, this approach is frequently economically viable, as expression levels in inclusion body system are extremely high and simple washing procedures can be used to remove most host cell proteins resulting in relatively high initial purities. Figure 1.28 shows an example of *E. coli* cells overexpressing a protein as inclusion bodies and the corresponding SDS–PAGE analyses at various stages in the process.

The yeast cells *Saccharomyces cerevisiae* and *Pichia pastoris* have been successfully used for over expression of various recombinant proteins including insulin and albumin. *S. cerevisiae* is also used for the production of Hepatitis surface antigen. However, mammalian cells are used for the majority of biopharmaceutical proteins. Although mammalian cell culture is generally more complex, these cells can perform complex post-translational modifications, such as glycosylation, which are often critical to proper biological and

pharmacological activity. Chinese hamster ovary cells (CHO) are the most commonly used mammalian expression system, especially for recombinant antibodies. The human cell line PerC6 has been developed more recently and is also used for production of some recombinant proteins. CHO and PerC6 are able to overexpress antibodies in concentrations as high as 15 mg/ml. Such high product titers are achieved mainly because the expression of proteins in these systems is generally independent of growth. As a result, the cells can be maintained in a perfusion bioreactor for long times, e.g. up to 30 days, in a viable productive state and can be cultivated to very high densities, e.g. up to $2 \times 10^8$ cells/ml.

Other mammalian cell lines used in recombinant protein production include baby hamster kidney (BHK) cells, Vero cells, and Madin–Darby canine kidney (MDCK) cells. Several coagulation factors that require $\gamma$-carboxylation are produced in BHK. Vero cells, derived from monkey kidneys, and MDCK cells are used for the production of vaccines. Insect cells such as cells from pupal ovarian tissue of the fall armyworm *Spodoptera frugiperda* (Sf9) have also been proposed for overexpression of recombinant proteins and, in particular, for the expression of virus-like particles and virus base gene therapy vectors. These cells can be infected with insect viruses such as the Autographa californica nuclear polyhedrosis virus, also known as baculovirus. Although easy to handle, so far, processes including insect cell systems have not been licensed for the industrial production of biopharmaceutical proteins.

In the past, transgenic mammals were considered as excellent production systems, as proteins can be secreted into their milk with titers up to 5 mg/ml. Two decades ago, such titers could not be achieved with mammalian cell culture. Thus, these expression systems were often preferred. Modern advances in cell culture, however, have made it possible to achieve routinely even higher titers in bioreactors while avoiding the enormous complexities of dealing with animals and allowing much more robust operation in a closed, more controllable system. As a result, transgenic animals no longer play a significant role in biopharmaceutical manufacturing.

Plants and plant cells are also potential candidates for expression of proteins. Although different from those that occur in mammalian cells, post-translational modifications are also possible in these organisms. However, expression in leaves, stems, or seeds presents significant recovery and purification challenges because the tissue or seeds must be ground up, pressed, or extracted yielding very complex mixtures. An interesting expression system is rhizo-secretion, where the protein is secreted into a cultivation fluid from hairy roots. Plant cell cultures, on the other hand, present fewer downstream processing difficulties as these cells grow in very simple media. Currently, one licensed therapeutic protein – Elelyso (taliglucerase alfa) – produced in carrot cells is on the market. There is a lot of research in this field with many ongoing trials. Most probably, there will be more therapeutic proteins manufactured by plant cell culture or plant cells in the future. In addition, other manufacturing applications of plant cell culture are progressively tested such as industrial enzymes, polymer degrading enzymes, feed additives, and edible vaccines, which require less purification efforts [41].

## 1.3.2 Host Cell Composition

The composition of the host cells has important effects on downstream processing, especially when the product is expressed intracellularly, as in this case, cellular components are the major impurities. However, host cell components are also found extensively as impurities in secreted products, as cell lysis always occurs at least to some extent during fermentation and cell culture. In fact, in some cases, cultivation procedures that yield high titers, such as those used for antibody production, result in partial lysis of the cells, which, in turn, causes contamination of the product with host cell components. An overview of the composition and physical characteristics of the major host cells is given in Tables 1.13 and 1.14.

As can be seen in Tables 1.13 and 1.14, mammalian cells contain more protein but less nucleic acids compared to bacteria. Bacteria and yeast also contain numerous cell wall components. These components are frequently insoluble and are efficiently removed during early processing steps, but these early steps are affected significantly by the cell density and the broth viscosity. Extremely high

**Table 1.13** Composition of common host cells for expression of recombinant proteins.

| | Composition (% dry weight) | | | | | |
|---|---|---|---|---|---|---|
| Organism | Proteins | Nucleic acids | Lipids | Cell count per ml | Dry mass (mg/ml) | Wet mass (mg/ml) |
| *E. coli* | 50 | 45 | 1 | $10^{11}$ | 20 | 100 |
| Yeast | 50 | 10 | 6 | $10^{10}$ | 80 | 400[a] |
| Filamentous fungi | 50 | 3 | 10 | $10^{9}$ | 130 | 400[b] |
| Mammalian cells | 75 | 12 | Up to 10 | $10^{7}$–$10^{8}$ | 0.17–1.7 | 1–10 |

a) For high density culture of *P. pastoris* grown on glucose medium.
b) For high density culture of *P. pastoris* grown on methanol.

**Table 1.14** Composition of single cells for expression of recombinant proteins.

| Component | Amount per CHO cell[a] | Amount per *E. coli* cell |
|---|---|---|
| Total DNA (pg) | $7 \pm 1.2$ | 0.017[b] |
| Total RNA (pg) | $18 \pm 3.0$ | 0.10 |
| Total Protein (pg) | $146 \pm 22$ ($2.4 \times 10^{7}$ with an approx. $M_r$ of 40 000) | 0.2 ($3 \times 10^{6}$ molecules with an approx. $M_r$ of 40 000) |
| Dry weight (pg) | $263 \pm 31$ | 0.4 |
| Wet weight (pg) | 2500 | 2 |
| Diameter | $18\,\mu m$ | $0.5\,\mu m \times 3\,\mu m$ |
| Volume ($cm^3$) | $1.7 \times 10^{-9} \pm 0.2$ | |

a) Data kindly provided by Nicole Borth from BOKU.
b) A fast-growing *E. coli* cell contains in average the genome in fourfold repetition. Each genomic DNA weight about 0.0044 pg.

cell densities can be obtained for yeast, particularly for *P. pastoris*, for which cell densities up to 400 g of cells per liter have been reported. Such suspensions are extremely difficult to clarify; often, the suspension must be diluted in order to be able to centrifuge the broth. Nucleic acids are present in the form of DNA and all kinds of RNAs. These compounds result in high broth viscosity but are often rapidly degraded by mechanical shear or by endogenous nucleases.

A final consideration is the mechanical stability of the host cells. Bacteria and yeast cells are generally mechanically very stable and shear rates above $10^6 \, \text{s}^{-1}$ are necessary to break these cells. Such high shear rates are attained only using special equipment such as high-pressure homogenizers or French presses. By comparison, mammalian cells are much weaker. The burst force needed to destroy a yeast cell is in the range of 90 μN, whereas the burst force for mammalian cells is in the range of 2–4 μN. The burst force increases with culture length in a batch culture, which is consistent with the observation that older cells are more difficult to disrupt.

### 1.3.3 Culture Media

Modern biopharmaceuticals are commonly produced with the so-called defined media whose components are chemically defined. In the past, yeast, meat, and soy extracts, produced by proteolytic degradation and extraction, were commonly used for cultivation of bacteria and yeast cells. The standardization of such raw materials was extremely difficult, resulting in substantial batch-to-batch variations. Similarly, until recently, it was common to supplement cultivation media for mammalian cells with fetal calf serum in concentrations up to 10%. Beside added complexity and cost, such supplements can introduce undesirable adventitious agents, such as prions, which can significantly increase the downstream processing challenge. Although testing for such agents may still be required, the use of defined media greatly simplifies downstream processing.

Media for industrial cultivation of bacteria are usually very simple and provide the essential sources of carbon, nitrogen, and phosphate needed by these simple organisms. Examples are given in Table 1.15. Sometimes, a cocktail of trace elements is added, but frequently trace elements already present in the water are sufficient.

When the fermentation pH is controlled by the addition of NaOH, conductivities as high as 40 mS/cm can be reached at the end of the cultivation period. Such high conductivities can interfere with downstream processing operations such as ion exchange requiring dilution or diafiltration steps. This difficulty may be circumvented using ammonia for pH control, which typically results in lower conductivity of the culture supernatant.

In addition to the salts and sugar that are required for cell metabolism, production of recombinant proteins in bacteria typically requires the addition of an inducer, such as isopropyl-β-D-thiogalactopyranoside (IPTG), which is used to activate protein expression when a certain cell density is reached. The use of natural compounds as inducers is advantageous as such species are readily degraded in the culture and are not significant impurities. On the other hand, detergents or oils added to the culture as antifoaming agents, although present in relatively

**Table 1.15** Composition of defined culture media for cultivation of *E. coli*.

| Compound | Concentration (mg/ml) |
|---|---|
| Glucose | 1.0 |
| $Na_2HPO_4$ | 16.4 |
| $KH_2PO_4$ | 1.5 |
| $(NH_3)_2PO_4$ | 2.0 |
| $MgSO_4·7H_2O$ | 0.2 |
| $CaCl_2$ | 0.01 |
| $FeSO_4·7H_20$ | 0.0005 |

small amounts, can affect downstream processes as they tend to foul membranes and chromatography matrices.

Culture media for yeast are similar to those used for *E. coli*. Methanol is used frequently used as the inducer for systems based on the alcohol oxidase promoter (AOX) expression system. Mammalian cell culture, on the other hand, requires much more complex media including glucose as carbon source, amino acids, vitamins, inorganic salts, fatty acids, nucleotides, pyruvate, and butyrate. This basal medium is supplemented with proteins for oxygen transport, hormones, and growth factors. Oxygen transport proteins such as transferrin have bound iron. In order to create a totally protein-free medium, these proteins are often replaced by iron chelators such as ferric citrate, ferric iminodiacetic acid, ferric ammonium citrate, and tropolone (2-hydroxy-2,4,6-cycloheptatriene-1-one). These compounds can, however, interfere with downstream processing. For example, under slightly acidic conditions, ferric citrate forms a gel, which is difficult to separate from proteins and other biomacromolecules.

Several other additives that may be present in cell culture media also impact downstream processing. pH indicators, such as phenol red, added to laboratory-scale culture media often bind to ion exchange resins and are best avoided for large-scale cultivation. Hydrophilic polymers, such as poly(propylene glycol) or poly(ethylene glycol), are often needed in concentrations up to 0.02% to protect the cells from shear stress.

pH in mammalian cell culture is typically regulated by addition of $CO_2$, although high-density cultures may require addition of NaOH. Final conductivity values of less than 17 mS/cm are typical, making direct capture by ion exchange easier compared to capture from yeast and *E. coli* homogenates.

### 1.3.4 Components of the Culture Broth

In general, before harvest, the culture broth contains the following components: intact cells, debris from lysed cells, intracellular host cell components, unused media components, compounds secreted by the cell, and enzymatically or chemically converted media components. Oxygen is depleted because during primary

recovery, the oxygen supply is shut down and the residual dissolved oxygen is rapidly consumed. The low oxygen content can induce necrosis, and cells may rapidly die and lyse during this phase. Some cell types begin autolysis after just 30 minutes without oxygen. As a result, depending on the cell type, rapid separation of the cells from the broth supernatant is necessary to keep host cell impurities low. The culture broth may contain high concentrations of $CO_2$, which is produced by the residual cells and fragments and shifts pH to acidic region. $CO_2$ has a much higher solubility than oxygen in aqueous solutions, so that substantial amounts can be present. Dissolved $CO_2$ can be rapidly liberated when the pH is lowered to conduct certain downstream processing steps, thereby forming bubbles that may then enter chromatography columns and disrupt the packing.

Intracellular host cell components appearing as impurities in a culture supernatant can be estimated from the following equation:

$$\begin{pmatrix} \text{Impurity} \\ \text{concentration} \end{pmatrix} = \begin{pmatrix} 1 - \dfrac{\text{Fraction of}}{\text{viable cells}} \end{pmatrix} \times \begin{pmatrix} \text{Cell} \\ \text{count} \end{pmatrix} \times \begin{pmatrix} \text{Amount of} \\ \text{impurity per cell} \end{pmatrix}$$

where the amount of intracellular impurity components per cell can be estimated from Tables 1.13 and 1.14. Thus, from a downstream processing perspective, high cell viability is desirable. This is not always possible, however. For example, in high-titer antibody production by cell culture, the cells are often lysed in the final stage of cultivation. As a result, mammalian cell culture supernatants from cultivation with defined media contain host cell proteins in the range of 1–3 mg/ml (1000–3000 ppm). These levels must be reduced to less than 100 ppm in the final product.

### 1.3.5 Product Quality Requirements

Biopharmaceutical product quality and process validation are subject to regulations by the individual governments. In the United States, the regulatory framework is published in the Code of Federal Regulations 21 (21 CFR), Subchapter F Biologics. The US FDA is responsible for its implementation. In the European Union, the regulatory framework is still under the sovereignty of the individual member states, although an EU-wide umbrella organization, the European Medical Agency (EMA), has been founded with the goal of harmonizing the EU regulatory structure. The existence of multiple regulatory frameworks adds complexity to the global biopharmaceutical industry. Thus, the International Council for Harmonisation of Technical Requirements for Pharmaceuticals for Human Use (ICH) (https://www.ich.org/home.html) has been established to develop a common international regulatory framework. ICH brings together the regulatory authorities of Europe, Japan, and the United States, as well as experts from the pharmaceutical industry in the three regions to discuss scientific and technical aspects of product registration. The purpose is to make recommendations on ways to achieve greater consistency in the interpretation and application of technical guidelines and of the requirements for product registration. Although progress is being made in the harmonization process, the regulatory framework of each country remains in effect.

The guidelines ICH Q5B (Analysis of the expression construct in cell lines used for production of rDNA-derived protein products), ICH Q5D (Derivation and characterization of cell substrates used for production of biotechnological/biological products Share), and ICH Q6B (Specifications: test procedures and acceptance criteria for biotechnological/biological products) regulate the medicinal products derived from recombinant DNA.

The following Categories of Therapeutic Biological Products are reviewed and regulated by the Center for Drug Evaluation and Research (CDER), a suborganization of the US FDA:

- Monoclonal antibodies for in vivo use.
- Proteins intended for therapeutic use, including cytokines (e.g. interferon), enzymes (e.g. thrombolytics), and other novel proteins, except for those that are specifically assigned to Center for Biologics Evaluation and Research (CBER, e.g. vaccines and blood products). This category includes therapeutic proteins derived from plants, animals, or microorganisms and recombinant versions of these products.
- Immunomodulators (nonvaccine and nonallergenic products intended to treat disease by inhibiting or modifying a pre-existing immune response).
- Growth factors, cytokines, and monoclonal antibodies intended to mobilize, stimulate, decrease, or otherwise alter the production of hematopoietic cells in vivo

The following categories of therapeutic biological products are defined and regulated by the Center for Biologics Evaluation and Research (CBER), also a suborganization of the US FDA:

- Cellular products, including products composed of human, bacterial, or animal cells
- Gene therapy products
- Vaccines
- Allergenic extracts used for the diagnosis and treatment of allergic diseases and allergen patch tests.
- Antitoxins, antivenins, and venoms
- Blood, blood components, plasma-derived products (for example, albumin, immunoglobulins, clotting factors, fibrin sealants, and proteinase inhibitors), including recombinant and transgenic versions of plasma derivatives (for example clotting factors), blood substitutes, plasma volume expanders, human or animal polyclonal antibody preparations including radiolabeled or conjugated forms, and certain fibrinolytics such as plasma-derived plasmin, and red cell reagents.

It is possible that certain recombinant protein products will fall under the regulatory oversight of multiple regulatory bodies simultaneously resulting in further complexities.

The manufacture of biological products for pharmaceutical applications must follow general guidelines that are established by the regulatory framework. Three keywords summarize the principal product quality requirements: *purity*, *potency*, and *consistency*. Industrially, these requirements must be met with

processes that are economically viable and that can bring products to the market rapidly. Downstream processes must be designed to obtain sufficient purity while maintaining the potency or pharmacological activity in a consistent manner.

### 1.3.5.1 Types of Impurities

The purity requirements of therapeutic proteins are defined by ICH guidelines. Purity requirements for biopharmaceuticals vary depending on the particular application, dose, and cell line used in the manufacturing process. Thus, it is not possible to specify absolute values. In the ICH guideline ICH Q6B, a distinction is made between *process-related impurities* and *product-related impurities*. Process-related impurities are derived from the manufacturing process, from the cell substrates, from the cell culture components, and from the downstream processing. Examples are host cell proteins, host cell DNA, inducers, antibiotics, and other media components, as well as ligands or other chemicals leached from chromatography media. Product-related impurities are defined as molecular variants arising during manufacture and/or storage, which do not have properties comparable to those of the desired product with respect to activity, efficacy, and safety. Examples are precursors, certain degradation products, aberrant glycoforms, and aggregates. An important distinction can be made, however, among the different types of impurities between *critical impurities* or *noncritical impurities*. This is determined during process development by risk analysis studies. A noncritical impurity is an inert compound without biological relevance. This can be, for instance, residual PEG from an extraction process or a harmless host cell component such as a lipid. On the other hand, endotoxins or growth factors secreted into the culture supernatant are examples of critical impurities, as they can exert adverse biological activity. These impurities have to be traced throughout the process and extensive testing and documentation of their removal is generally required.

*Contaminants* in a bioproduct include all adventitiously introduced materials not intended to be part of the manufacturing process, such as microbial proteases and/or microbial species. Contaminants should be strictly avoided and/or suitably controlled with appropriate in-process acceptance criteria or *action limits* for drug substance or drug product specifications (see ICH guideline Q6B). For adventitious virus, mycoplasma, or prion contamination, the concept of the action limit is not applicable. For this case, the strategies proposed in ICH guidance documents Q5A (Quality of Biotechnological/Biological Products: Viral Safety Evaluation of Biotechnology Products Derived from Cell Lines of Human or Animal Origin) and Q5D (Quality of Biotechnological/Biological Products: Derivation and Characterization of Cell Substrates Used for Production of Biotechnological/Biological Products) should be considered. The starting material should be preferably free of the agent. Spiking experiments at small scale must be conducted to demonstrate clearance and their presence in final product must be strictly controlled. Finally, bioburden, originating from microbial contamination from air or personnel or from inadequately cleaned equipment, can also have serious effects and must be carefully monitored and controlled.

Table 1.16 summarizes the measures required to demonstrate the removal of virus, mycoplasma, prions, and other impurities. This demonstration is

**Table 1.16** Measures typically required to demonstrate the removal of adventitious virus, mycoplasma, prions contaminants, and impurities.

| Measure | Virus, mycoplasma, prions | Other impurities |
|---|---|---|
| Spiking experiments to demonstrate clearance | Yes | No |
| Starting material preferably free of agent | Yes | No |
| Clearance measured at each step | No | Yes |
| Control of final product | Yes | Yes |

**Table 1.17** Typical virus clearance values of retroviridae virus for purification of recombinant antibody (LRV).

| Step | LRV | Minimum LRV | Maximum LRV |
|---|---|---|---|
| Protein A affinity chromatography | 5 | 2 | 8 |
| Low pH chemical inactivation | 5 | 2 | 8 |
| Cation exchange chromatography in bind elute mode | 3 | 1 | 7 |
| Anion exchange chromatography in flow-through mode | 5 | 2 | 8 |
| Virus filter | 5 | 2 | 8 |

Minimum LRV and Maximum LRV define the range of LRV values reported in regulatory submissions.
Source: Data from Miesegaes et al. 2010 [42].

usually done experimentally using scale-down models because, obviously, it would counterproductive to intentionally contaminate the production plant with an adventitious agent. For these determinations, also known as spiking experiments, a bolus of an adventitious agent, e.g. a virus, is added to the raw feed stream entering a purification process step. The virus titer before purification $a' = \log_{10}(\text{Feed titer})$ and after purification $a'' = \log_{10}(\text{Harvest titer})$ is determined and the log-virus reduction factor (LVR) is calculated as follows:

$$\text{LVR} = \log_{10}(\text{Feed titer}) - \log_{10}(\text{Harvest titer}) = a' - a'' \tag{1.31}$$

In order to account for the effect of volume changes (e.g. a mere 1 : 10 dilution results in a LVR of 1), the following individual reduction factor, $R_i$, is also calculated:

$$R_i = \text{LVR} - \log_{10}\frac{V'}{V''} \tag{1.32}$$

where $V'$ and $V''$ are the feed and harvest volumes, respectively. Finally, the LVR of the individual process steps are added together to arrive at a cumulative LVR for the entire process. Table 1.17 illustrates a typical virus clearance levels for various steps of an antibody purification process.

Although almost all mammalian cells are infected with virus making viral clearance validation obligatory, efforts are frequently made to omit these procedures by utilizing platform processes that have demonstrated clearance efficiency.

**Table 1.18** Examples of input and output parameters in a chromatographic separation process of proteins.

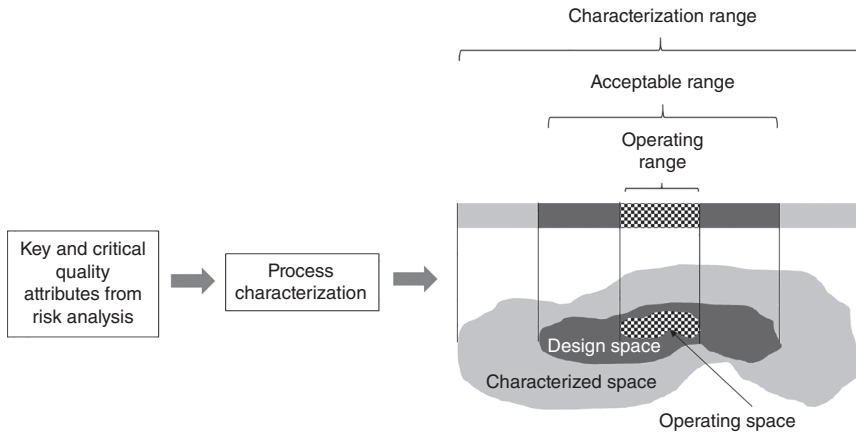|  | **Parameter** | | | | | | | |
|---|---|---|---|---|---|---|---|---|
| Input parameter | pH | Temperature | Ionic strength | Load | Flow rate | Column height (residence time) | | |
| Output parameter (performance) | Purity | Concentration | Stability | Yield | DNA content | Host cell protein content | Endotoxin content | Column back pressure |

### 1.3.5.2 Validation

Validation is a critical aspect of biopharmaceutical process development. According to existing ICH definitions, critical parameters that may influence product quality have to be validated. After validation, a standard operating protocol (SOP) is established, which describes the process and the allowed variations. *Critical operational parameters* are defined as a limited subset of process parameters that significantly affect critical product quality attributes when varied outside a meaningful, narrow (or difficult to control) operational range. Consider, for example, the operation of a chromatographic purification step. As shown in Table 1.18, this operation will require the definition of a number of operating conditions as inputs, which, in turn, will result in certain performance characteristics as outputs.

In order to validate the process, the input parameters must be varied over suitable ranges and the corresponding outputs measured. The critical parameters are then defined based on these experiments, which are usually performed at small scale. Suitably narrow operational ranges are established for these parameters as well as for noncritical parameters. As the latter do not affect critical product quality attributes, their ranges will normally be broader than those for critical parameters (see Figure 1.29).

*Quality by design* (QbD) is a global regulatory initiative having the goal of enhancing pharmaceutical development through the proactive design of pharmaceutical manufacturing process and controls that consistently deliver the intended product characteristics. The "*design space*," a critically important concept in QbD, is defined as the range of conditions under which the process can be operated while maintaining the desired product quality. Although there is no regulatory requirement to have a design space, it is very commonly used to set operating limits. A design space can be described in terms of ranges of material attributes and process parameters. The design space can also be described through a complex mathematical model. The design space can be determined by the following:

- First-principles approach, which is a combination of experimental data and mechanistic knowledge of chemistry, physics, biology, and engineering to model and predict performance.
- Statistically designed experiment such as design of experiments (DOE).

**Figure 1.29** Definition of operation ranges for critical process parameters.

- Scale-up correlations, which are semiempirical approaches to translate operating conditions between different scales or pieces of equipment.

A design space allows more operational flexibility and is highly encouraged by health authorities.

*Process analytical technologies* (PAT) concepts are required to ensure consistent quality of the biopharmaceutical. PAT is a system for designing, analyzing, and controlling processing through timely measurements of critical quality and performance attributes of raw materials, in-process streams, and process parameters with the goals of ensuring final product quality.

### 1.3.5.3 Purity Requirements

It is difficult to describe absolute purity requirements as they depend on the intended use of the biopharmaceutical, dose, and risk–benefit ratio. Table 1.19 provides only approximate values meant to serve as general guidelines.

Aggregates are an important concern for many biopharmaceutical proteins. It has been shown that aggregates can induce immune reactions or cause other side effects. Moreover, aggregates may constitute seeds for precipitation, thereby reducing the shelf life of a product. As a result, controls on the percentage of dimers, oligomers, and higher aggregate forms have been tightened and the maximum allowed values reduced. The leakage of ligands or other leachable chemicals from chromatographic media and membranes is also an important concern because these materials can be immunogenic or toxic.

Viral contamination is obviously a critical issue as in the past it has been responsible for many iatrogenic diseases, such as those occurring because of contaminated blood products. As absolutely complete removal of these adventitious agents is not possible, limits are often established on the basis of a risk–benefit analysis. For example, the World Health Organization (WHO) accepts for a vaccine one adverse case in $10^9$ applications – hence, the probability value of $10^{-9}$ value suggested by Table 1.19. More recently, limits on the allowed amount of host cell DNA per dose have been relaxed from 10 pg per dose to

**Table 1.19** General guidelines for purity, consistency, and potency of protein biopharmaceuticals.
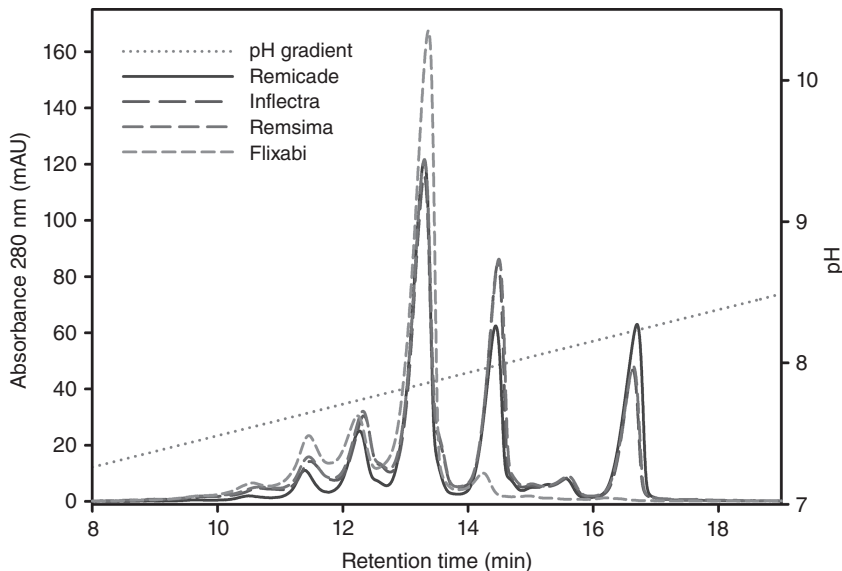
| Criteria | Requirement |
| --- | --- |
| *Purity* | |
| Specific protein content | >99.9% |
| Dimer/oligomer content | <1.0% |
| Ligand leakage | Usually <1 ppm |
| Virus content | Absence with a probability of $<10^{-9}$ |
| DNA content | <10 ng/dose |
| Endotoxin content | <5 EU/(kg h) |
| Prion content | Absence with a probability of $<10^{-9}$ |
| *Consistency* | |
| Microheterogeneity | Permitted, but consistent |
| Impurities | Permitted, but consistent |
| *Potency* | |
| Folding | Correctly folded |
| Mutations | Correctly expressed, no mutations |
| Processing | Correctly processed |

10 ng of DNA per dose. Clinical practice and postmarket studies have shown that, in general, host cell DNA does not pose a high risk, except for some potent compounds such as growth factors or hormonally active compounds.

Many biopharmaceuticals do not consist of individual molecular entities – rather, they consist of a large number of similar isoforms or variants (some recombinant antibody products contain as many as 2000 identifiable variants). Because the biological and pharmacological activity can vary dramatically among different isoforms, it is important to maintain the distribution of these variants within established acceptable ranges. Because of the complexity of bioproduction systems, similar consistency must also be maintained for the impurity profiles, as determined from analytical assays, in order to assure product safety. Finally, test systems must be established to control the potency in vitro and, where necessary, in vivo.

## 1.4 Biosimilars

A *biosimilar* is a product that is similar in terms of quality, safety, and efficacy to an already licensed reference biopharmaceutical defined as the "originator product." Most protein biopharmaceuticals comprise a high number of molecular entities or variants. These molecular entities are dependent on the primary sequence of the protein, on the nature of the expression system, and on the production process. Even minute variations on the production parameters and

**Figure 1.30** Example of differences in variant composition between originator product and various biosimilars. The variants in each product were separated by analytical cation exchange chromatography with a linear pH gradient. The originator antibody Remicade is compared to the three anti-TNF-α biosimilar antibodies Remsima, Inflectra, and Flixabi. Source: Adapted from Beyer et al. 2019 [43].

materials can affect the composition of product variants. Therefore, it is usually not possible to develop a perfect copy of a biopharmaceutical (see Figure 1.30 as an example). The most important goal is thus to achieve similarity with respect to clinical properties.

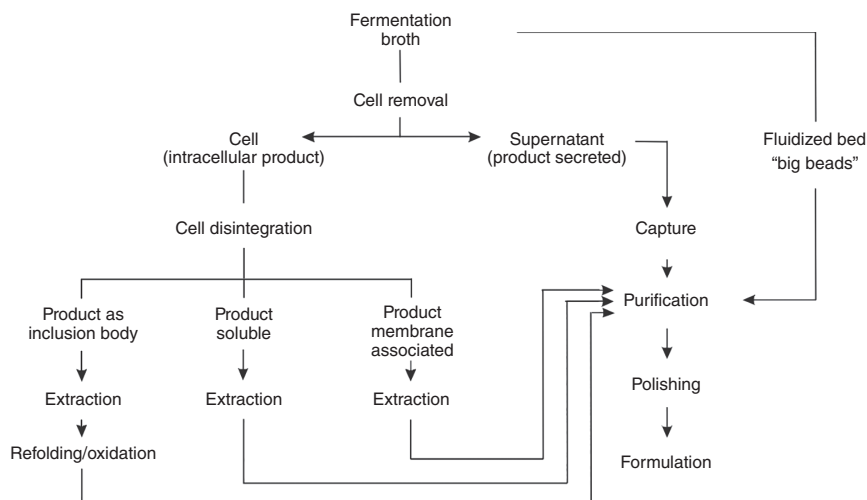## 1.5 Role of Chromatography in Downstream Processing

Chromatography is the principal tool for the purification of biopharmaceuticals. This can be explained by certain advantages of chromatography over other unit operations. Firstly, chromatography provides very high separation efficiencies, which allow the resolution of complex mixtures of components having very similar molecular properties. Properly designed chromatography columns can have the separation efficiency of hundreds or even thousands of theoretical plates. By comparison, extraction and membrane filtration are usually limited to only a few theoretical plates. Secondly, chromatography columns packed with high capacity adsorbents are ideal for capture from the dilute solutions found in bioprocessing. In such systems, a large volume of solution can be contacted efficiently with a small amount of a high-capacity adsorbent packed in a column, resulting either in the rapid concentration of the product. When operated in a flow-through mode, such columns can result in the nearly complete removal of contaminants present in small concentrations without affecting the unbound product. By comparison, liquid–liquid extraction systems typically require similar volumes of the two

phases in order to function properly, so that concentration is not very feasible. A further advantage is that chromatography can be performed in an almost closed system and the stationary phase can be easily regenerated. Finally, chromatography is well established in many practical biopharmaceutical manufacturing processes, and suitable equipment and packing materials are readily available. A perceived disadvantage of chromatography is the difficulty of scale-up within the constraints of the biopharmaceutical industry. However, as will be shown in the remaining chapters of this book, proper application of engineering tools in combination with adequate measurements allows the design of optimum columns for large-scale applications. Indeed, as shown by Kelley [3], chromatographic purification processes can be considered and can be technically and economically viable for protein purification at scales as high as 20 tons/yr. Although no current product is currently made at such a large scale, the popularity of biopharmaceuticals is increasing rapidly so that one could envision such scales in the future.

Figure 1.31 illustrates the structure of a generic process for the recovery and purification of a biological product produced by microbial fermentation or animal cell culture. The initial steps where cells are separated are often referred to as *primary recovery*. These steps require different strategies depending on whether the product is secreted into the culture medium or expressed in the cell, either as inclusion body, in soluble form in the cytosol or periplasm, or anchored in the membrane. Generally, chromatography plays a minor role in these initial steps, which are focused on the removal of suspended solids such as cells or cell debris. Sedimentation, centrifugation, deep bed filtration, and microfiltration or combinations thereof are normally used for these early steps. However, chromatography, implemented through the use of fluidized or expanded beds, can also be used for the direct capture of secreted proteins from cell culture supernatants.

In these systems, the liquid flows upward through an initially settled bed of dense adsorbent particles. Above a certain flow velocity, the bed expands and the particles become fluidized allowing free passage of cells and other suspended matter while the product is directly captured by the adsorbent. The approach can be effective for dilute suspensions. However, as bed expansion is directly influenced by the feed density and viscosity, the operation tends to be critically affected by variations in the composition of the broth. In practice, the high viscosity and cell density encountered in modern fermentation technology (up to 400 mg/ml wet cell mass for *P. pastoris* or 200 mg/ml for *E. coli*) make it difficult to implement this approach reliably at the industrial scale. An alternative possibility for early capture without clarification is to use adsorption beds packed with large particles, sometimes referred to as "big beads." If the particles are larger than about 400 μm in diameter, the interparticle spaces are sufficiently large to allow passage of small cells and cell debris. Although the efficiency of capture is reduced by the diffusional limitations that accompany the larger particle diameter, the ensuing reduction in the number of processing steps can provide overall economic and operational advantages. Unlike expanded beds, packed bed processes are not very sensitive to feed viscosity so that reliable operation with large diameter beads can be achieved even with viscous feedstocks.

As can be seen in Figure 1.31, following primary recovery, the general downstream processing scheme consists of successive *capture*, *purification*,
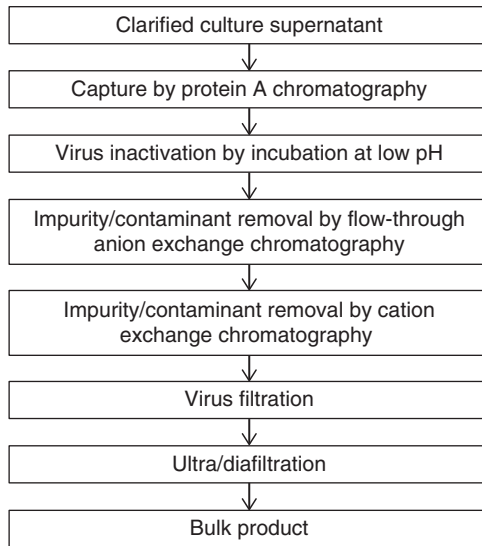
**Figure 1.31** Generalized downstream processing flow sheet for purification of proteins starting with the unclarified fermentation broth. Note: "big beads" are large-size (400–500 μm diameter) chromatography media used to directly capture small proteins from viscous solutions that could contain some particulate matter.

and *polishing* steps, each comprising one or more unit operations. With only a handful of exceptions, current industrial processes for biopharmaceuticals almost exclusively employ chromatography for these three critical steps.
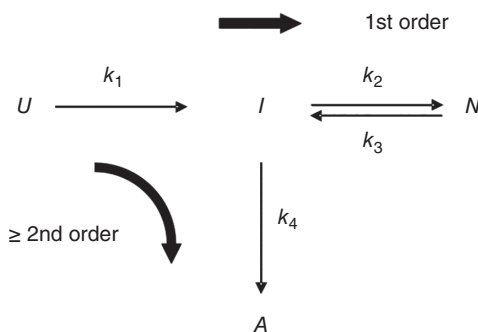
For purification of recombinant antibodies, capture is almost always realized using a selective adsorbent comprising Staphylococcal Protein A immobilized in porous beads (see Figure 1.32). This highly selective ligand allows direct loading of the clarified culture broth on the capture column, which binds selectively the antibody. In the subsequent steps, purification and polishing are conducted with ion exchange and hydrophobic interaction columns to remove host cell proteins and aberrant protein variants. Note that intermediate, nonchromatographic steps are also included. Firstly, incubation at low pH for virus inactivation and then "virus filtration" are implemented for viral clearance. Secondly, an ultrafiltration/diafiltration step is included for buffer exchange and final formulation.

Besides purification, chromatography also finds other uses in bioprocessing. An important example is the use of chromatography to facilitate refolding of solubilized protein, which is sometimes a bottleneck in industrial processes. Without simultaneous separation, misfolding and, especially, aggregation compete with the correct folding pathway. Aggregation may originate both from nonspecific (hydrophobic) interactions of predominantly unfolded polypeptide chains as well as from incorrect interactions of partially structured folding intermediates. As can be seen in Figure 1.33, aggregation reactions are second- (or higher) order processes, whereas correct folding is generally determined by first-order reactions [44].

In practice, refolding conditions (e.g. denaturant concentration) are adjusted so that the equilibrium distribution favors the formation of native protein (i.e. $k_2 \gg k_3$). The formation of intermediates is generally very fast so that $k_1$ can be

| Clarified culture supernatant |
| Capture by protein A chromatography |
| Virus inactivation by incubation at low pH |
| Impurity/contaminant removal by flow-through anion exchange chromatography |
| Impurity/contaminant removal by cation exchange chromatography |
| Virus filtration |
| Ultra/diafiltration |
| Bulk product |

**Figure 1.32** Example of a process flow diagram for the purification of recombinant antibodies.



**Figure 1.33** Simplified reaction scheme for protein refolding with aggregation of intermediates.

neglected. For the case where $k_3 \to 0$, we effectively have competing first- and second (or higher)-order reaction. For these conditions, refolding in a batch system is described by the following equations:

$$\frac{d[U]}{dt} = -(k_2[U] + k_4[U]^n) \tag{1.33}$$

$$\frac{d[N]}{dt} = k_2[U] \tag{1.34}$$

where the brackets denote concentrations, $k_2$ is the net rate constant of folding, $k_4$ is the net rate constant of aggregation, and $n$ is the reaction order. An analytical solution of these equations is available for $n = 2$ and is given by the following equation [45]:

$$Y(t) = \frac{k_2}{[U]_0 k_4} \ln \left\{ 1 + \frac{[U]_0 k_4}{k_2}(1 - e^{-k_2 t}) \right\} \tag{1.35}$$

where $Y$ is the yield of the refolding reaction and $[U]_0$ is the initial concentration of unfolded protein. As time approaches infinity, the final yield of native protein

**Table 1.20** Refolding and aggregation rate constants of a protein in refolding by batch dilution and matrix-assisted refolding using size exclusion chromatography.

| Process | Folding $k_2$ (min$^{-1}$) | Aggregation $k_4$ (ml/(mg min)) |
|---|---|---|
| Batch dilution | 0.0012 | 0.3 |
| Matrix-assisted refolding by SEC | 0.0012 | 0.01 |

Source: Data from Schlegl et al. 2005 [46].
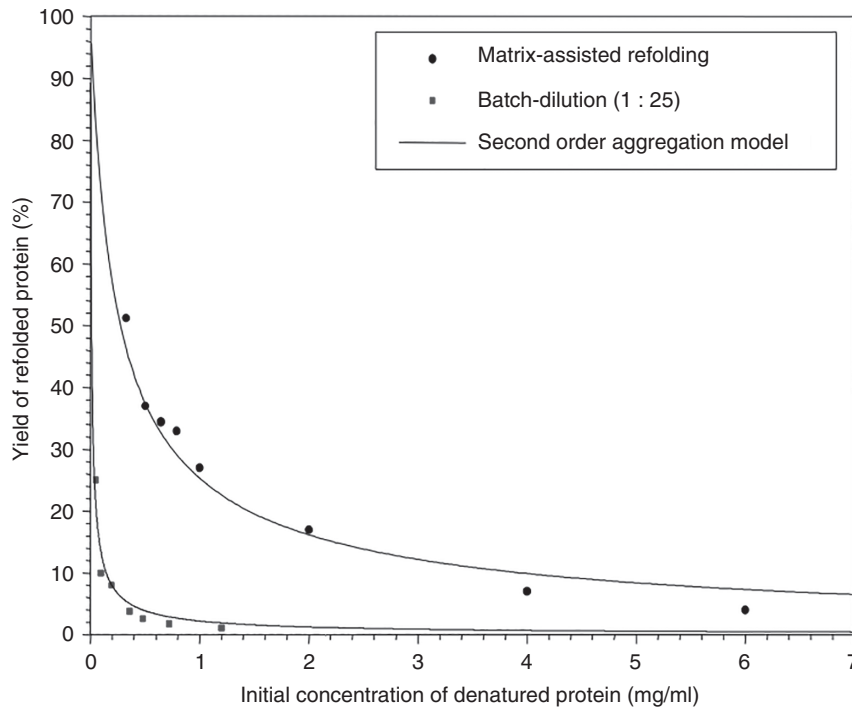
is then given by the following equation:

$$Y(t \rightarrow \infty) = \frac{k_2}{[U]_0 k_4} \ln \left( 1 + \frac{[U]_0 k_4}{k_2} \right) \tag{1.36}$$

This result suggests that dilution (i.e. low $[U]_0$) is a simple and effective way of ensuring high refolding yields. Although this is effective and widely used in practice, the ensuing large solution volumes complicate further downstream processing and increase cost. Refolding in chromatographic columns, also known as matrix-assisted refolding, can be a valuable refolding alternative to reduce the need for extensive dilution. The underlying mechanism leading to improved folding in chromatographic columns is not completely understood and may depend on the specific nature of the protein and the selected conditions. However, the effects can be dramatic as shown, for example, in Figure 1.33. In this case, refolding was conducted by separating the denaturing agent (urea) from the unfolded protein by SEC, thereby allowing refolding to occur within the column. This resulted in a greater yield of folded protein compared to a simple dilution process. The apparent aggregation rate constant in this case was about 30 times smaller compared to that for the dilution process (Table 1.20). A possible explanation of this result is that aggregation may be inhibited within the matrix pores by steric hindrance allowing a greater portion of the protein can follow the path toward correct folding.

In the example given in Figure 1.34, the unfolded protein was passed over a size exclusion column and the denaturant was slowly removed. Comparison of kinetic constants between conventional refolding by dilution into a refolding buffer and matrix-assisted refolding confirms that aggregation is suppressed in the column (Table 1.20).

SEC-promoted refolding is also possible in continuous processes, which can also include a recycling system for aggregated protein. Yield and productivity of a continuous refolding system using pressurized continuous annular chromatography (P-CAC) considering initial protein concentration, residence time, and recycling rate were extensively studied for α-lactalbumin as a model protein [47]. Also, countercurrent chromatography systems such as the simulated moving bed (SMB) can be used for continuous matrix-assisted refolding.
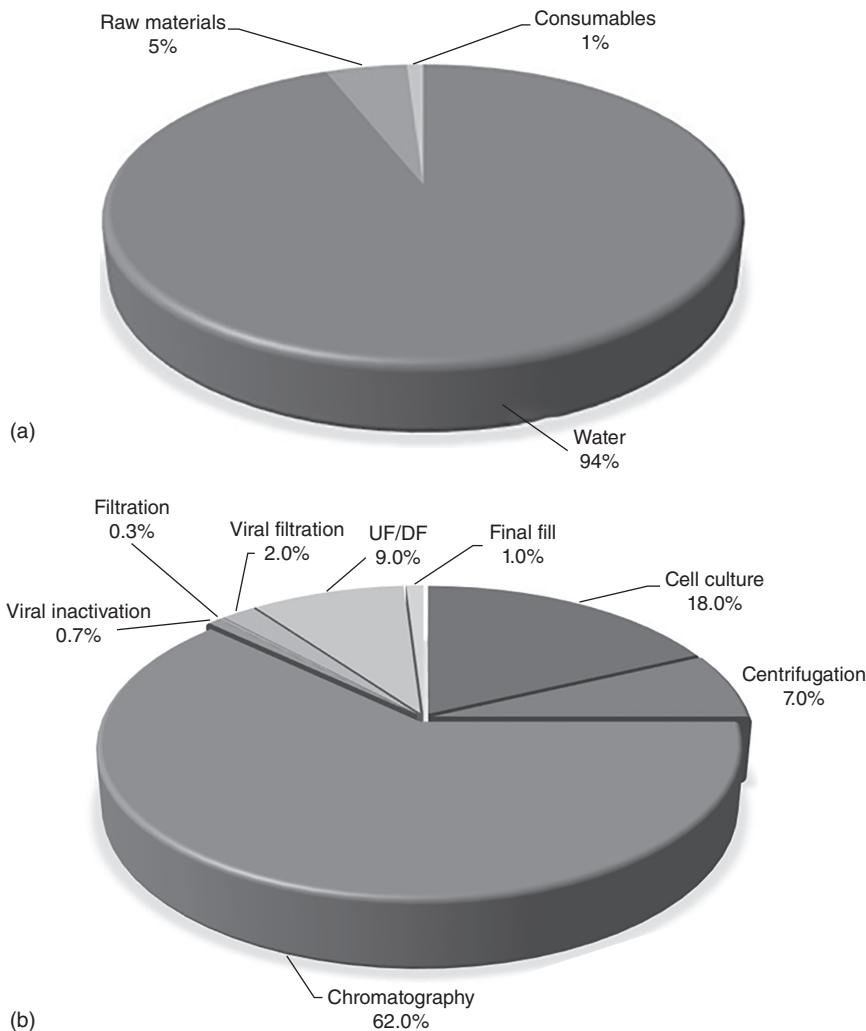
Ion exchange, affinity adsorption, and hydrophobic interaction have also been used to facilitate refolding. A method based on the adsorption of the unfolded protein on an ion exchange resin was introduced by Creighton [48].

**Figure 1.34** Refolding yield of a protein by batch dilution and with matrix-assisted refolding using size exclusion chromatography. Source: Adapted from Schlegl et al. 2005 [46].

Further improvements of this method rely on using more sophisticated buffers during loading and elution. The methods can also be executed in a continuous manner. The surface contact can initiate refolding. In many instances, it has been observed that final refolding takes place after the protein has been eluted from the column. Immobilization of denatured proteins on a solid support often limits its flexibility and therefore its ability to regain the native configuration because of multipoint interaction with the matrix. Introduction of an N- or C-terminal poly-histidine-tag allowed the reversible one-point immobilization of the denatured protein on a solid support based on immobilized metal affinity chromatography (IMAC). Refolding can be achieved by a simple buffer exchange in a stepwise or gradient manner. The use of HIC was described for the refolding of lysozyme, BSA, α-amylase, and recombinant γ-interferon [49].

Immobilized folding catalysts and artificial chaperones have also been suggested as refolding aids. Mimicking in vivo folding systems was a further step of improving in vitro refolding yield. The chaperones or compounds mimicking chaperones are immobilized on a chromatography matrix. The protein solution is passed through such columns. The folded proteins are slightly retarded, and the denaturant is exchanged. The immobilized chaperones prevent aggregation. Thus, a refolding can be achieved at higher concentrations or yield. One has to keep in mind that a chaperone acts in a stoichiometric manner. Thus, large amounts of chaperone protein are necessary to avoid aggregation.

Figure 1.35 (a) Main contributions to the process mass intensity for monoclonal antibody manufacturing. The PMI metric was determined from 14 process datasets from biopharmaceutical firms for both large-scale (≥12 000 l) and small-scale (≤5000 l) antibody manufacturing operations. (b) Percentage of total water use by each unit operation. Source: Adapted from Budzinski et al. 2019 [50].

## 1.6  Environmental Impact of Biopharmaceutical Manufacturing

The *environmental footprint of biopharmaceutical manufacturing* has been largely neglected in the past. However, water scarcity is increasingly being recognized as a major future threat to business. Especially in arid areas, biopharmaceutical manufacturing may become very expensive because of the large amounts of water used by both upstream and downstream processes.

Additionally, pharmaceuticals have become pervasive in the environment, which poses a risk for people and animals. In response to these pressures, the pharmaceutical industry is focusing increasingly on the environmental impact of pharmaceutical and biopharmaceutical manufacturing aiming to improve the efficiency of its operations. In addition to reducing the environmental impact, processes with a reduced environmental footprint also tend to be more economic. The process mass intensity (PMI), already in use in other industrial sectors, has been proposed as a standardized metric to assess environmental impact in pharmaceutical manufacturing. Such a metric is independent of fluctuations of raw material, water, and energy costs. Total PMI is defined based on the amount of active pharmaceutical ingredient (API) produced as follows:

$$\text{Total PMI} = \frac{\text{Total water, raw materials, and consumables used in the process (kg)}}{\text{Amount of API produced (kg)}}$$

As shown in Figure 1.35 for a typical biopharmaceutical manufacturing process producing a recombinant antibody, chromatography is the greatest contributor to the PMI, especially in terms of the amount of water consumed. There are ample opportunities to optimize chromatography with respect to water consumption including maximizing the utilization of binding capacity in processes, and thus reducing the consumption of elution buffers, implementing process intensification by using smaller, more efficient separation columns that can be washed, eluted, and cleaned with smaller amounts of aqueous buffers, and by implementing process integration, thereby reducing or eliminating intermediate steps, such as diafiltration or dilution, which can consume large amounts of buffers. The engineering fundamentals covered in this book provide, we believe, not only opportunities to maximize productivity and reducing downstream processing costs but also have the potential to help reduce the environmental impact of chromatography processes contributing to our quest for sustainability and to creating environmentally friendly production processes.

## References

1 Walch, G. (2018). *Nat. Biotechnol.* 36: 1136.
2 Buchacher, A. and Iberer, G. (2006). *Biotechnol. J.* 1: 148.
3 Kelley, B. (2007). *Biotechnol. Progr.* 23: 995.
4 Rinaldi, A. (2008). *EMBO Rep.* 9: 1073.
5 Berg, J., Tymoczko, J., and Stryer, L. (2006). *Biochemistry*. Palgrave Macmillan.
6 Voet, D.J. and Voet, J.G. (2006). *Biochemistry, Textbook and Student Solutions Manual*. New York: Wiley.
7 Branden, C. and Tooze, J. (1991). *Introduction to Protein Structure.* New York: Garland Publishing, Inc.
8 Duerkop, M., Berger, E., Dürauer, A., and Jungbauer, A. (2018). *Biotechnol. J.* 13: 1800062.
9 Ueberbacher, R., Haimer, E., Hahn, R., and Jungbauer, A. (2008). *J. Chromatogr. A* 1198: 154.

**10** Schmid, F.X. (1997). Optical spectroscopy to characterize protein conformation. In: *Protein Structure: A Practical Approach*, 2e (ed. T.E. Creighton), p. 261. Oxford: IRL Press.

**11** Dayhoff, M.O. (1974). *Fed. Proc.* 33: 2314.

**12** Nötling, B. (1999). *Protein Folding Kinetics.* Berlin: Springer.

**13** Gokana, A., Winchenne, J.J., Ben-Ghanem, A. et al. (1997). *J. Chromatogr. A* 791: 109.

**14** Kneuer, C. (2005). DNA Vaccines – An overview 1. *DNA Pharmaceuticals* (ed. M. Schleef). Weinheim: Wiley-VCH.

**15** Petsch, D. and Anspach, F.B. (2000). *J. Biotechnol.* 76: 97.

**16** Mach, H., Middaugh, C.R., and Lewis, R.V. (1992). *Anal. Biochem.* 200: 74.

**17** Wilfinger, W.W., Mackey, K., and Chomczynski, P. (1997). *Biotechniques* 22: 474.

**18** Tanford, C. (1976). *Biochemistry* 15: 3884.

**19** Kuehner, D.E., Engmann, J., Fergg, F. et al. (1999). *J. Phys. Chem. B* 103: 1368.

**20** Tanford, C. and Roxby, R. (1972). *Biochemistry* 11: 2192.

**21** Righetti, P.G. and Caravaggio, T. (1976). *J. Chromatogr. A* 127: 1.

**22** Mahn, A., Lienqueo, M.E., and Salgado, J.C. (2009). *J. Chromatogr. A* 1216: 1838.

**23** Hopp, T.P. and Woods, K.R. (1981). *Proc. Natl. Acad. Sci. U.S.A.* 78: 3824.

**24** Kyte, J. and Doolittle, R.F. (1982). *J. Mol. Biol.* 157: 105.

**25** Kimerer, L.K., Pabst, T.M., Hunter, A.K., and Carta, G. (2019). *J. Chromatogr. A* 1601: 121.

**26** Kimerer, L.K., Pabst, T.M., Hunter, A.K., and Carta, G. (2019). *J. Chromatogr. A* 1601: 133.

**27** Fox, S. and Foster, J.S. (1957). *Introduction to Protein Chemistry.* New York: Wiley.

**28** Green, A. (1932). *J. Biol. Chem.* 95: 47.

**29** Cohn, E. and Edsall, J.T. (1943). *Proteins, Amino Acids, and Peptides.* New York: Academic Press.

**30** Hofmeister, F. (1888). *Arch. Exp. Pathol. Pharmakol.* 24: 247.

**31** Gekko, K. and Timasheff, S.N. (1981). *Biochemistry* 20: 4677.

**32** Ahrer, K., Buchacher, A., Iberer, G., and Jungbauer, A. (2006). *J. Biochem. Biophys. Methods* 66: 73.

**33** Monkos, K. and Turczynski, B. (1999). *Int. J. Biol. Macromol.* 26: 155.

**34** Sibileva, M.A., Zatyaeva, A.A., and Matveeva, N.I. (2001). *Mol. Biol.* 35: 73.

**35** Reisner, A.H. and Rowe, J. (1969). *Nature* 222: 558.

**36** Tyn, M.T. and Gusek, T.W. (1990). *Biotechnol. Bioeng.* 35: 327.

**37** Prazeres, D.M.F. (2008). *Biotechnol. Bioeng.* 99: 1040.

**38** Fee, C.J. and Van Alstine, J.M. (2004). *Bioconjugate Chem.* 15: 1304.

**39** Zhu, M. and Carta, G. (2014). *J. Chromatogr. A* 1326: 29.

**40** Cussler, E.L. (1997). *Diffusion – Mass Transfer in Fluid Systems*, 2e. Cambridge: Cambridge University Press.

**41** Park, K.Y. and Wi, S.J. (2016). *J. Plant Biol.* 59: 559.

**42** Miesegaes, G., Lute, S., and Brorson, K. (2010). *Biotechnol. Bioeng.* 106: 238.

**43** Beyer, B., Walch, N., Jungbauer, A., and Lingg, N. (2019). *Biotechnol. J.* 14: 4, 1800340.

**44** Clark, E.D. and Hevehan, L.D. (1997). *Biotechnol. Bioeng.* 54: 221.

**45** Kiefhaber, T., Rudolph, R., Kohler, H.H., and Buchner, J. (1991). *Biotechnology* 9: 825.

**46** Schlegl, R., Necina, R., and Jungbauer, A. (2005). *Chem. Eng. Technol.* 28: 1375.

**47** Schlegl, R., Iberer, G., Machold, C. et al. (2003). *J. Chromatogr. A* 1009: 119.

**48** Creighton, T.E. (1990). Process for the Production of a Protein. US Patent 4,977,248 [filed on 29 March 1985 and issued on 11 December 1990].

**49** Geng, X. and Chang, X. (1992). *J. Chromatogr. A* 599: 185.

**50** Budzinski, K., Blewis, M., Dahlin, P. et al. (2019). *New Biotechnol.* 49: 37.