# Hybrid modeling in cross-flow ultrafiltration: Predicting the flux and the rejection factor evolution of a binary protein solution

Ignasi Bofarull Manzano, B.Sc.

## Master's Thesis

to achieve the university degree of Master of Science (MSc)

Master's degree programme: Biotechnology

Submitted to

## University of Natural Resources and Life Sciences, Vienna

Head of Department of Biotechnology:    Prof. Reingard Grabherr

Supervisor:                             Priv.Doz. Dipl.-Ing. Dr. Astrid Dürauer

Co-Supervisor:                          Dipl.Ing. Dr. Mark Dürkop

Vienna, March 2021

# Statutory Declaration

"I declare in lieu of an oath that I have written this master thesis by myself and that I have not used any sources or resources other than stated for its preparation. I further declare that I have clearly indicated all direct and indirect quotations."

Vienna, _____          _____

        Date                                       Signature

# Eidesstattliche Erklärung

Ich erkläre eidesstattlich, dass ich die Arbeit selbständig angefertigt, keine anderen als die angegebenen Hilfsmittel benutzt und alle aus ungedruckten Quellen, gedruckter Literatur oder aus dem Internet im Wortlaut oder im we-sentlichen Inhalt übernommenen Formulierungen und Konzepte gemäß den Richtlinien wissenschaftlicher Arbeiten zitiert, durch Fußnoten gekennzeich-net bzw. mit genauer Quellenangabe kenntlich gemacht habe.

Wien, _____          _____

        Datum                                    Unterschrift

# Acknowledgment

# Zusammenfassung

Tangentialflussfiltration wird während der Aufreinigung von biopharmazeutischen Produkten zu verschiedensten Zwecken im Prozess eingesetzt wie z.B. bei der Zellernte, der Virusinaktivierung und Aufkonzentrierung des Produktes. Biotechnologische Prozesse sind immer mit Prozessvariabilität verbunden wie Variationen des Produkttiters und unterschiedliche Verunreinigungsprofile.. Rein mechanistische Modelle sind bislang nicht in der Lage die komplexen Zusammenhänge zwischen Prozessbedingungen und Produktqualität ausreichend zu erfassen.

Im Zuge dieser Arbeit wurden drei Hybridmodelle erstellt, welche die Proteine einer Zwei-Komponenten-Lösung zu jedem Zeitpunkt des Filtrationsprozesses quantifizieren können. Die Modelle wurden für zwei Modellproteinen, BSA und Lysozym, erstellt, die das Produkt und die Verunreinigung nachstellen. Die Modelle unterschieden sich in der Art wie die Membrandurchlässigkeit für die Verunreinigung (Lysozym) berechnet und aktualisiert wurde. In einem Hybridmodell wurde ein konstanter Wert für die Membrandurchlässigkeit der Verunreinigung angenommen, während zwei Weitere einen variablen Wert dafür heranzogen. Die Vorhersage des Permeatflusses mit Hilfe der Hybridmodelle wurde mit der Vorhersage mittels der mechanistischen Filmtheorie verglichen.

Die Hybridmodelle konnten die Permeatflussabnahme und die Konzentration des Produkts und der Verunreinigung gut vorhersagen. Dies war für verschiedenste Kombinationen von Konzentrationen und mechanischen Prozessparametern möglich. Die Vorhersagen basierten nur auf der anfänglichen Konzentration der Proteine, dem Transmembrandruck und dem Querstrom und benötigten nur wenige Trainingsexperimente. Die entwickelten Hybridmodelle bilden die Basis für Softsensoren zur Echtzeitmodellierung und modellbasierter Prozesskontrolle für industrielle Anwendungen. Solche Modelle ermöglichen das Weiterentwickelnd der biopharmazeutischen Industrie nach den Gesichtspunkten der Quality-by-Design Prozessführung.

# Abstract

Cross-flow filtration is a powerful technique used during several purification processes in the biopharmaceutical industry, such as cell harvesting, virus clearance or protein concentration. Due to the intrinsic variability of biological processes, the product titers and impurity profiles in the fermentation broth commonly differ from batch to batch. Solely mechanistic models do not fully describe the complex interactions between these components and the process conditions. In the present work, three hybrid model structures were established that accounted for varying concentrations of a two-component solution, a model system based on BSA and lysozyme, which mimicked product and impurity, respectively. The models differed according to how the impurity rejection factor was calculated, ranging from static values to dynamically updating structures. The flux predictions of the hybrid models were compared to the predictions obtained by the well-established mechanistic stagnant film model and the recently established one-component-hybrid model.

The established two-component-hybrid models accurately described the flux and the concentration evolution of the two proteins over a wide range of process parameters and product-to-impurity ratios. The predictions were solely based on the initial protein concentrations, transmembrane pressure and cross-flow velocity parameters and the models trained on a minimum set of training experiments. On the opposite, the stagnant film theory and the one-component-hybrid models exhibited larger errors for flux and showed poorer prediction for the impurity, since they are based on one component only. The presented hybrid models are the foundation for the implementation of soft-sensors for real-time monitoring and model predictive control of complex multi-component solutions in industry. They thereby pave the way for moving biopharmaceutical manufacturing towards the Quality-by-Design initiative and next generation bioprocessing.

# Contents

# Abbreviations

| | |
|---|---|
| ANN | Artificial neural network |
| BSA | Bovine serum albumin |
| CF | Cross-flow velocity |
| CHO | Chinese Hamster Ovary |
| CP | Concentration polarization |
| CPP | Critical process parameter |
| CQA | Critical quality attribute |
| DF | Diafiltration |
| *E.coli* | Escherichia coli |
| HCP | Host-cell proteins |
| HM | Hybrid model |
| HPLC | High-pressure liquid chromatography |
| Lys | Lysozyme |
| MLR | Multiple linear regression |
| MnLR | Multiple non-linear regression |
| MPC | Model predictive control |
| MWCO | Molecular weight cutoff |
| NRMSE | Normalized root-mean-squared error |
| ocHM | One-component-hybrid model |
| PAT | Process analytical technology |
| PBS | Phosphate-buffer saline |
| QbD | Quality by design |
| SEC | Size exclusion chromatography |
| SFM | Stagnant film model |
| tcHM | Two-component-hybrid model |
| TMP | Transmembrane pressure |
| UF | Ultrafiltration |

# Nomenclature and Symbols

A         membrane area [m$^2$]

A$_n$       n$^{th}$ virial coefficient [m$^{3n}$/kg$^n$]

c$_0$        initial bulk concentration [g/L]

c$_{B,i}$      bulk concentration of component $i$ [g/L]

c$_G$        gel layer concentration [g/L]

c$_m$       concentration at the membrane surface [g/L]

c$_p$        permeate concentration [g/L]

D         diffusion coefficient [m$^2$/s]

dt        time increment [s]

J          permeate flux [LMH] or [m3/m2·s]

k          mass transfer coefficient [LMH]

R$_i$        rejection factor of component $i$

R$_{average}$ average rejection from training/test experiment

$R_{bl}$       hydraulic resistance of the boundary layer [m$^{-1}$]

$r_{bl}$       specific resistance of the boundary layer [m$^{-2}$]

$R_m$       hydraulic membrane resistance [m$^{-1}$]

$r_m$       specific membrane resistance [m$^{-2}$]

V$_0$        initial reservoir volume [mL]

V$_B$        bulk/reservoir volume [mL]

V$_p$        permeate volume [mL]

## Symbols

$\Delta P$      applied pressure [Pa]

$\Delta\Pi_b$      osmotic pressure difference between the feed and the permeate side [Pa]

$\Delta\Pi_{bl}$    osmotic pressure difference along the boundary layer [Pa]

$\Delta\Pi_m$    osmotic pressure difference along the membrane [Pa]

$\delta$        thickness of the boundary layer [m]

$\eta_0$       viscosity of the solvent [Pa·s]

# 1. Theoretical Background

During the production of biopharmaceuticals, the product of interest is usually found together with a complex mixture of components in the fermentation broth or cell culture supernatant. These components are categorized as impurities and contaminants, which have to be removed throughout a number of purification steps, so-called downstream processes, in order to ensure the quality and consistency of the product. Process-related impurities are those that derive from the manufacturing process itself, such as media components, host cellular components - DNA, RNA, proteins, lipids…-, metabolites or endotoxins. Product-related impurities are molecular variants of the product that have the potential to differ from it in terms of activity, efficacy or safety – i.e., product precursors, aggregates, glycovariants, isoforms, oxidized or deamidated forms. On the other hand, contaminants are adventitiously introduced materials, such as viruses, prions, microbial proteases or microbial species. The acceptable amount of each of these compounds in the final product is defined by the regulatory agencies through the so-called product specifications, which must be fulfilled by the manufacturer to ensure the product safety and efficacy[1,2].

In the first part of the introduction, a brief summary of the overall downstream processes used in the production of biopharmaceuticals is provided, with a more detailed information on the role of cross-flow filtration in the following. Later on, the different modeling approaches for filtration processes are described. Finally, the most common impurities found in bioprocessing of bacterial and mammalian production systems and their titer compared to the product are also explained.

## 1.1 Downstream processing

Downstream processes account for the major part of the production costs of biopharmaceuticals, representing up to the 80% of the total costs in some cases[3,4]. It is for this reason, that the development of separation and purification techniques with enhanced efficiency, throughput and selectivity have to deal with increasing product titers and high dose therapies to reduce the cost of goods[5,6]. The purification process of a biotechnological product typically consist of four main sections: primary recovery, capture, purification and polishing[7].

During the first section, primary recovery, whole cells and cell debris are separated from small soluble molecules by mechanical separation methods, mostly centrifugation and depth filtration. In the case of microbial fermentations, where the cell debris densities are usually higher due to the disruption of the cells from the harvest, a primary clarification step by centrifugation usually takes place followed by depth filtration. On the contrary, in mammalian cell culture, where the molecules of interest are secreted out of the cell and the cell debris densities are lower, depth filtration is used more often[8], which has the advantage of lower

4

equipment costs and easier scalability (by the membrane area). Further, some depth filters are charged, which allows for additional chemical separation of the host cell impurities[9]. Finally, a second type of filters, known as sterile filters, of uniform and smaller pore size (0.2-0.4 μm), are usually used after depth filtration in order to ensure that the subsequent downstream processing units are protected from unwanted cell debris or contaminants, thereby preventing them from fouling. After the separation of the product from the biomass, the capture phase is the next step. The main goal of the capture step is to isolate the product and to reduce the volume of the solution. This is essential in order to reduce the material consumption and the size of equipment used in the subsequent steps - which are particularly the most expensive ones, both in fixed and variable costs. Afterwards, during the purification phase, most contaminants are depleted and the concentration of the target component is further increased. The most prominent method for the capture and purification phase is ion exchange chromatography. Finally, during polishing, the removal of trace contaminants and product-related impurities such as product aggregates take place achieving the final desired high purity levels. The type and number of final polishing steps will depend on the product as well as on its required purity. For biopharmaceuticals (purity >90%) the polishing stage usually consists of a combination of chromatographic and filtration steps. The reason for using filtration in this stage is that chromatographic methods are the most expensive steps during downstream processing -representing over one third of the total downstream processing costs[7]-, and therefore, product concentration and volume reduction becomes essential for the economic viability of the process. In addition, UF/DF units allow for buffer exchange while concentrating, thereby ensuring the most optimal conditions for each chromatographic step.

## 1.1.1 Membrane technology

Membrane processes play a critical role in the purification of biotechnological products, being used in several applications such as cell harvest, clarification, sterile filtration, virus removal, protein concentration or buffer exchange. Three main types of membranes are used in downstream process of biopharmaceuticals, depending on the pore size and therefore their removal characteristics. Microfiltration membranes have pore sizes ranging from 0.05 to 10 μm and are used to retain cells and cell debris, whereas ultrafiltration membranes have smaller pores, between 1 and 20 nm, and are designed to retain proteins and other macromolecules. The pore size of ultrafiltration membranes is commonly given as the molecular weight cut-off (hereafter referred to as MWCO), which is defined as the lowest molecular weight of a solute with at least 90% retention by the membrane. Finally, membranes designed for virus filtration fall between micro- and ultrafiltration membranes and have pore sizes from 20 to 70 nm. It is important to mention that depth filters are not typically considered as membranes, since they retain the components throughout their porous structure rather than on their surface[9]. These

filters are made of a matrix of fibers with an appropriate filter aid (e.g., diatomaceous earth, perlite or activated carbon) and a binder, which creates their porous filter media and increases the mass of particles they can retain before start clogging. Some sterile or virus filters may also retain the feed particles throughout their depth, however, in this case this is because they have multiple layers of different pore size or because their pore size distribution is graded in order to increase the permeability or the capacity of the filter.

Regardless of the membrane type, there are two operational modes for pressure-driven filtration systems: normal filtration (also called dead-end filtration), where the feed flow is made pass through (vertical) the membrane filter (Figure 1A), and cross-flow filtration (CFF) (also called tangential flow filtration), where the feed is directed parallel to the membrane (Figure 1B). In both operational modes, the solution that goes through the membrane is known as permeate or filtrate, while the solution that is retained on the feed side is called retentate or filtrate. In cross-flow filtration, the components that are retained on the membrane surface are swept away by the feed flow and the retentate is directed back to the feed reservoir. This tangential flow reduces the buildup of concentration gradients on the membrane (so-called concentration polarization (CP) layers) and fouling effects and allows for direct recovery of the product in solution. Contrary, in dead-end filtration the recovery of the retained material is uncommon. Dead-end filtration is mainly used for removing undesired components from the working solution in applications such as sterile, depth or virus filtration, when the product of interest is recovered on the permeate side.

Hence, when the goal is to concentrate the product, ultrafiltration membranes with MWCOs between 3 and 5 times smaller than the target protein mass are commonly chosen, ensuring that there are no losses of product to the permeate. On the contrary, for applications such as cell and lysate clarification, the filter pore size should be at least 10 times higher than the molecular weight of the target protein, so it can freely pass through the filter and the yields are not compromised.
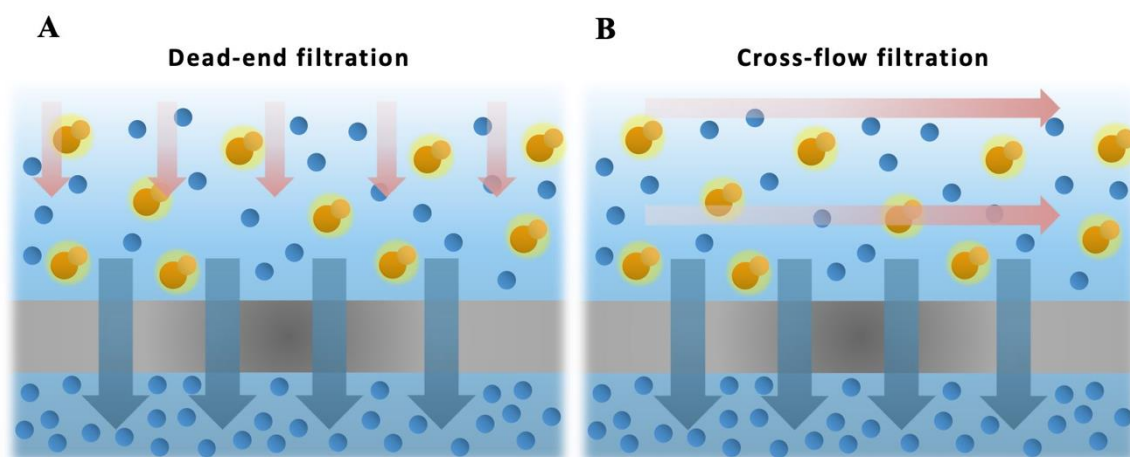


Figure 1: Schematic representation of (A) dead-end and (B) cross-flow filtration modes.

All types of pressure-driven membrane separation processes are dominated by two main phenomena: concentration polarization and membrane fouling. Concentration polarization (CP) is the accumulation of solutes or particles in a thin liquid layer (called CP layer or boundary layer) adjacent to the membrane surface as a result of mass transfer limitations during filtration. As a consequence, the retained solutes in this layer will have higher concentrations than in the bulk solution and will create an additional hydraulic resistance to the flow of solvent through the membrane. Moreover, at high concentrations, these retained solutes create higher osmotic pressures towards the permeate side, which will counteract the applied effective pressure driving force. Therefore, concentration polarization is an inherent phenomenon of all membrane separation processes that starts to occur as soon as the solutes are convectively transported to the membrane by pressure-driven feed flows[10]. On the other hand, membrane fouling is a more general term that can be due to either the adsorption on and within the membrane pores and/or the deposition on the membrane surface of particles, colloids and macromolecules. Hence, it is usually used to describe long-term phenomena, where the filtration process has to be stopped and the membrane cleaned[11] due to a significant reduction of the permeate flux compared to the initial values.

It is important to mention that within the fouling definition, the term cake layer is usually used to refer to the accumulation of retained solids on the membrane in microfiltration, while the term gel layer is rather used in ultrafiltration for referring to the precipitation of soluble macromolecules when their concentration close to the membrane surface is too high. However, these terms are sometimes mixed up and used interchangeably by some authors to describe two different mechanisms by which explaining the phenomenon of concentration polarization in ultrafiltration, regardless of the operational mode type (cross-flow or dead-end filtration). In this case, the cake filtration model states that the concentration in the layer on top of the membrane is constant (there is no concentration gradient) and that its thickness increases with increasing the permeate volume, while the opposite being for the film theory (see Figure 2)[10,12,13,14].

### 1.1.1.1 Cross-flow filtration

Cross-flow filtration sweeps out the built-up layer on the filter surface, thereby reducing concentration polarization and fouling, which are the key limiting parameters in this process unit. This leads to higher permeate fluxes and consequently faster filtration processes. Due to this feature, cross-flow filtration systems are commonly used for applications such as cell perfusion and medium exchange during cell cultivation by microfiltration membranes, or protein concentration and buffer exchange by employing ultrafiltration membranes.

The basic configuration of a cross-flow filtration system is depicted in Figure 4. Usually they are operated under one type of process control mode: either constant transmembrane

pressure (TMP) or flux, although they can also be combined by an interplay of several process parameters[15]. In this thesis, the TMP control mode was employed, in which TMP is kept constant by adjusting the retentate valve - and if necessary also the permeate valve. Closing the retentate valve results in increased flux, while closing the permeate value resulted in decreased flux. The TMP control mode can be used at constant feed flow, constant retentate flow or constant pressure drop between feed and retentate (ΔP). On the other hand, in the flux control mode, the flux is maintained at a constant rate either by regulating the feed flow or the retentate and permeate flow, thereby increasing TMP with time. This control mode is mostly used in microfiltration, where the pores of the membrane are relatively large and the main process limitation is the creation of an additional hydraulic resistance to the flow due to the accumulation of cells and cell debris. However, in ultrafiltration, the dominant effect is the reduction in the effective pressure driving force due to osmotic pressure effects, as a result of the high protein retention on one side of the membrane, which increase non-linearly with concentration at high concentrated solutions[16]. Therefore, the flux control mode is not an option for ultrafiltration of proteins, where the TMP necessary to keep the same flux rate would dramatically increase with protein concentration, especially close to high values, known as limiting or critical fluxes. At this point, the flux would no longer depend on the applied transmembrane pressure but only on mass transfer affecting parameters, such as protein concentration in the bulk, buffer composition, cross-flow velocity or temperature. This region is called pressure-independent or mass transfer-limited region (see *1.2.1.1 Concentration polarization*). Additionally, working at high TMPs in turn increases the risk of entering in this region. It is for this reason that the optimal operating TMP range for cross-flow ultrafiltration processes is described as those values that are just below the pressure-independent region[15].
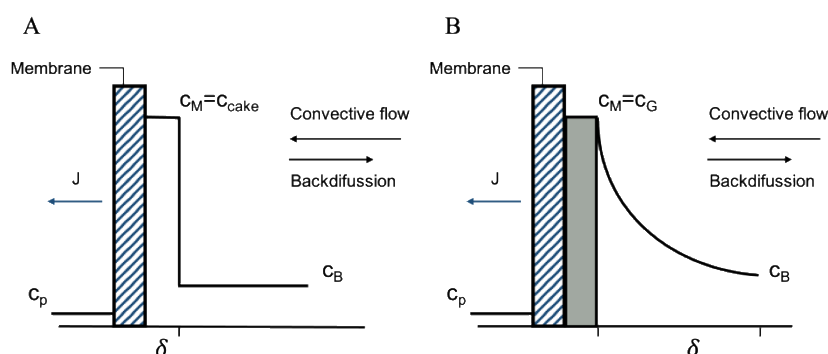


Figure 2: Concentration profile next to the membrane according to (A) cake- and (B) gel-layer concentration polarization models.

## 1.2 Modeling of filtration processes

### 1.2.1 Mechanistic modeling

The flow of a fluid through a porous medium was first described by Darcy's law[17], which established a proportional relationship between the filtrate flux of a fluid, the permeability of the porous medium and its thickness, the dynamic viscosity of the fluid and the applied pressure drop. In modeling of filtration, however, the permeability of the membrane and the membrane thickness are commonly used together under the term membrane resistance ($R_m$), in order to describe a membrane-specific parameter which only depends on the membrane characteristics -material, porosity, pore size distribution, thickness, geometry…- and that can be directly used to relate the flux of any kind of fluid in a certain membrane to the applied pressure drop, in the form of Eq. (1):

$$J = \frac{TMP}{\eta_0 * R_m} \tag{1}$$

Where $J$ is the filtrate flux, $TMP$ the applied transmembrane pressure and $\eta_0$ the viscosity of the solvent. $R_m$ is drawn from the inverse of the slope of Eq. (1), when recording the flux at different TMPs and after dividing it by the viscosity of the fluid.

When filtering a protein solution, the flux is usually lower compared to the pure solvent (Figure 3). The main reasons for this are concentration polarization and membrane fouling, although other phenomena such as protein-protein interactions when using protein mixtures or protein interactions with buffer salts may also lead to further deviations from the expected flux behaviors. Therefore, several factors may explain the flux behavior in a filtration process, being this is the reason why different mechanistic models have been developed to describe them[14]. The extent to which each of these factors occur and influence the flux will depend on many parameters, both of the used solution and membrane, and they will thus not always be the same. Consequently, there is currently -despite all the extensive research conducted in membrane technology- no general model able to describe the flux evolution in filtration processes in a holistic way. The process data must be instead analyzed in every case in order to determine the model that fits the flux best and that therefore better explains the underlying mechanism. However, in some cases it has been shown that different mechanistic models can describe the same flux behaviors -e.g. limiting fluxes can be described by both, gel-layer and osmotic pressure models[13]. In other cases, the main effect driving the flux evolution may change over time -e.g. from concentration polarization to membrane fouling governing flux-, therefore making it difficult to choose a single model to describe the entire filtration process[14]. Finally, phenomena that are not fully understood cannot be represented in mechanistic

models. All this usually leads the modeler to make some assumptions when using mechanistic models, both of the underlying mechanism governing the flux as well as for the determination of its coefficients. However, if the overall behavior of the process changes as a result of variations in the solution composition, which is especially likely in multi-component systems as the treated in bioprocessing, these assumptions might not hold anymore and the predictions lose accuracy.
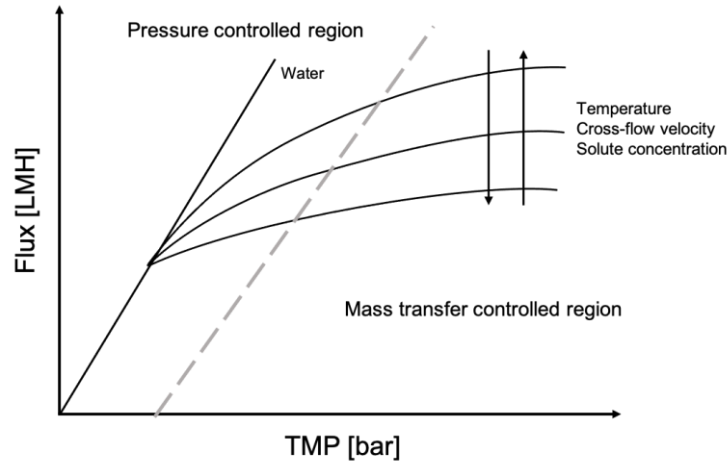


Figure 3: Steady-state fluxes in cross-flow filtration of a macromolecular solute as function of the applied transmembrane pressure. An increase in temperature or cross-flow velocity or a decrease in the solute concentration would lead to higher steady-state fluxes and vice versa.

In this section, a general overview of the classical most used mechanistic models for predicting the flux in cross-flow ultrafiltration of proteins is provided, with their advantages and disadvantages.

### 1.2.1.1 Concentration polarization

Concentration polarization is an inherent and immediately occurring phenomenon of all membrane separation processes. In the case of cross-flow ultrafiltration of proteins, it is the main parameter governing the flux evolution. The effects of concentration polarization on flux might be explained by three different mechanisms: resistance, gel-polarization and osmotic pressure models. They are all based on the universal-known film theory model (Eq. (2))[10], which describes the mass transfer of a solute across a boundary layer of thickness $\delta$ by the convective transport towards the membrane and the back-diffusion caused by the concentration gradient in the layer.

$$\frac{\partial c}{\partial t} + U \frac{\partial c}{\partial x} + J \frac{\partial c}{\partial y} = \frac{\partial}{\partial y}(D \frac{\partial c}{\partial y}) \tag{2}$$

If neglecting axial diffusion ($U$) of the solute, and assuming that the fluid velocity is constant at all positions on the membrane ($dx$), that the diffusion coefficient ($D$) is independent of the solute

concentration and steady-state conditions, Eq. (2) is then transformed to Eq. (3). Steady-state refers to constant thickness and no solute accumulation on the boundary layer under unchanged mass transfer conditions. This is usually reached after a short time in cross-flow ultrafiltration.

$$J = \frac{D}{\delta} \ln \left( \frac{c_m - c_p}{c_B - c_p} \right) = k \ln \left( \frac{c_m - c_p}{c_B - c_p} \right) \tag{3}$$

Where $c_m$, $c_B$ and $c_p$ are the solute concentrations on the membrane surface, bulk solution and permeate, respectively, and $k$ is the solute mass transfer coefficient.

### 1.2.1.1.1 Resistance models

According to resistance models, concentration polarization influences the flux by creating an additional hydraulic resistance to the solvent on top of the membrane, known as boundary layer resistance ($R_{bl}$), due to the higher solute concentrations there ($c_m$). This additional resistance, which depends on the thickness of the boundary layer ($\delta$) as well as on the so-called (solute) specific boundary layer resistance ($r_{bl}$), is then incorporated together with the membrane hydraulic resistance ($R_m$), in Eq. (1) to calculate the flux.

$$J = \frac{TMP}{\eta_0 * (R_m + R_{bl})} \tag{4}$$

Where:

$$R_{bl} = \int_0^\delta r_{bl} \, dy \tag{5}$$

$r_{bl}$ can be calculated by different approaches, depending on the model and whether the particles in the fluid are soluble or not. It strongly depends on the solute concentration at the boundary layer as well as on other solute specific characteristics such as the particle size or the density. In the case of compressible solutes, it depends on the applied pressure drop as well[14].

It is important to highlight that resistance models are mainly used for cross-flow filtration of solid particles[14], but not of soluble solutes, such as proteins. It has been both theoretically and experimentally demonstrated that they are equivalent to the osmotic pressure models in non-gelling, non-adsorption and completely solute rejection conditions[18]. The reason is that the osmotic pressure difference between the bulk solution and the permeate side ($\Delta\Pi_b$) is very small in most ultrafiltration processes, especially if compared to the applied pressure drop (Eq. (6-7)), therefore being the main osmotic pressure difference made on the boundary layer

($\Delta\Pi_{bl}$, Eq. (8)). Thus, under the aforementioned conditions, the derivation into Eq. (9) takes place:

$$\Delta\Pi_m = \Delta\Pi_{bl} + \Delta\Pi_b \tag{6}$$

$$TMP, \Delta\Pi_{bl} \gg \Delta\Pi_b \tag{7}$$

$$\Delta\Pi_m \simeq \Delta\Pi_{bl} \tag{8}$$

$$J = \frac{TMP - \Delta\Pi_m}{\eta_0 * R_m} = \frac{TMP - \Delta\Pi_b}{\eta_0 * (R_m + R_{bl})} \simeq \frac{TMP}{\eta_0 * (R_m + R_{bl})} \tag{9}$$

Where the first equality term of Eq. (9) is the osmotic pressure model (Eq. (11)) and the latter is Eq. (4). Resistance models are more difficult to construct and more prone to errors than osmotic pressure models[18], this being the reason why the latter are preferred.

### 1.2.1.1.2 Gel-layer models

Gel-layer (or gel-polarization) models explain that as a result of concentration polarization, the concentration at the membrane surface $c_m$ increases rapidly with the permeate flux and reaches a maximum value, the gel-layer concentration, $c_G$. Here, the solution on top of the membrane is no longer fluid and a gel-like layer forms due to protein precipitation. $c_G$ will be constant and any further increase in the applied pressure will result in an increased gel-layer thickness, that will prevent the flux from increasing, entering in the pressure-independent region. Under these conditions, the limiting or critical flux, $J^{lim}$, will only depend on parameters affecting the mass transfer rate, as aforementioned. A common assumption of these models is that the solute is completely retained by the membrane -$R_i$=1, and therefore $c_p$=0-, which allows Eq. (3) to be derived into Eq. (10). This assumption simplifies the complexity of the calculations necessary to determine the models parameters, since, under unchanged mass transfer conditions, $J^{lim}$ can then be linearly plotted over the logarithm of the bulk concentration, $c_B$, with the mass transfer coefficient $k$ being the slope and $\ln(c_G)$ the intercept point with the abscissa (Figure 10, Eq. (21)).

$$J = k \ln\left(\frac{c_G}{c_B}\right) \tag{10}$$

The gel-polarization model is one of the most used in modeling of cross-flow ultrafiltration processes, since it is simple, easy to construct, and moreover many experimental fluxes have shown to correlate well with it[14]. However, due to considering $c_m = c_G$ constant all the time, it can only explain the pressure-independent region, therefore failing when predicting fluxes at low TMPs or solute concentrations, such as at the beginning of a process, before the gel layer is formed. Nevertheless, according to this theory, the time to reach the gel-layer formation

should be short in filtration processes of proteins[10]. Furthermore, Zydney et al.[12] demonstrated that Eq. (10) was mathematically true for small $c_G/c_B$ ratios. He additionally extended it to describe the flux under concentration-dependent diffusivity conditions.

The most important factor when using this model is the determination of its mass transfer coefficient *k*. On the one hand, its theoretical calculation from proposed correlations in literature (e.g., Sherwood) is difficult and not very precise[19]. These expressions relate the mass transfer coefficient with physical properties of the solution (protein diffusivity, viscosity and density), geometries of the filtration device and process parameters (flow regime type, velocity…). Additionally, many of these properties depend at the same time on others parameters, for instance, the protein diffusion coefficient might change with protein charge (pI and pH), buffer conductivity, protein concentration, temperature and viscosity - which in turn also depends on temperature. Therefore, their calculation *a priori* is cumbersome and they should be measured in every case. Similarly, the coefficients relating the dimensionless numbers (Sherwood, Reynolds, Schmidt…) in this expression are highly dependent on the flow regime as well as on the velocity and concentration profiles, which in turn are equipment geometries specific. Thus, these coefficients should also be measured for accurate results. Summarizing, even though the theoretical calculation of the mass transfer coefficient *k* is possible, it is very prone to errors. On the other hand, its experimental determination can be very material and time consuming, as well as require good knowledge of the used equipment and extensive data fitting. It is for these reasons that the direct measurement of *k* from the experimental flux data has been suggested as the most precise and easier way[19]. The advantage of gel-layer over other filtration models is that its mass transfer coefficient *k* can be directly drawn geometrically from the slope of Eq. (21), J vs ln($c_B$). However, this requires the selection of just the data that are under the pressure-independent region (linear part, Figure 10), and therefore, it relies on the criterium of the modeler.

### 1.2.1.1.3 Osmotic pressure models

Finally, according to osmotic pressure models, concentration polarization decreases the flux by creating an osmotic pressure difference, which reduces the effective pressure driving force across the membrane[16]. In general, the osmotic pressure of a macromolecular solution is very small compared to a low molecular salt solution, however, the large concentrations of built up material at the membrane surface make the osmotic pressure of the solution increase substantially. Therefore, osmotic pressure models combine the film theory (Eq. (3)) with Darcy´s law (Eq. (1)) to include the osmotic pressure difference in the membrane, $\Delta\Pi_m$:

$$J = \frac{TMP - \Delta\Pi_m}{\eta_0 * R_m} \tag{11}$$

Although the osmotic pressure of an ideal diluted solution can be directly calculated by Van´t Hoff´s equation, for more concentrated non-ideal solutions, virial coefficients are needed (Eq. (12))[16]. These coefficients depend on several properties, and can be either directly measured for the given conditions[20], taken from reported values in literature[16] or calculated as a function of different parameters such as the particle excluded volume, hydration or Donnan effects[21]. Eq. (12) for calculating the membrane osmotic pressure will only be true for completely protein rejection factors ($R_i$=1), since otherwise the osmotic pressure difference between bulk and permeate solution, $\Delta\Pi_b$, should also be considered when calculating $\Delta\Pi_m$ (Eq. (6)).

$$\Delta\Pi_m = A_1 * C_m + A_2 * C_m^2 + A_3 * C_m^3 \tag{12}$$

Wijmans et al.[13] showed that osmotic pressure models indeed shared many characteristics with the gel-layer models: starting from Eq. (11), they derived a ratio between the resistances caused by the osmotic pressures and the ones related with the membrane itself ($\Delta\Pi_m \cdot n$)/($\eta_0 \cdot R_m \cdot k$). They mathematically proofed that the increment of flux with protein concentration in the bulk, i.e. $\partial J/\partial \ln c_B$, was almost equal to $-k$, as it is stated by gel-polarization models, at high values of this ratio. Therefore, even though osmotic pressure models do not initially describe a full limiting flux -since $c_m$ is not constant-, they can also simulate the pressure-independent region at high osmotic pressures, where the osmotic pressure effects would strongly predominate over the membrane resistance. Additionally, they can also explain the flux deviation from the pure solvent at low TMPs or protein concentrations -when $\Delta\Pi_m \approx 0$, $J \sim TMP/(\eta_0 \cdot R_m)$, Eq. (1)-, since they include the membrane resistance $R_m$, which was omitted in the gel-polarization models. Therefore, according to these models, an increase in the applied TMP would lead to a higher $c_m$, which would in turn increase the osmotic pressure difference along the membrane and thereby counteract the expected flux increase. At high protein concentrations, a further increase in TMP would have almost no effect -$\Delta\Pi_m$ increases non-linearly with solute concentration, Eq. (12)- entering in the pressure-independent region. Hence, instead of using a fixed $c_m$ value like gel-polarization models do, $c_m$ is a function of TMP. Wijmans et al.[13] also suggested that depending on the solute and solvent properties, the limiting flux can be better explained either by gel-layer formation or by osmotic pressure effects. According to them -also in agreement with what was published by Vilker et al.[16]- for small and medium molecular weight solutes (≤100 kDa), the osmotic pressure limitations are more likely to occur than the gel-layer, first, because small weight solutes have higher osmotic pressures than high weight molecules, and second, because they would also need higher $c_m$ in order to form a gel-layer, and the contrary being for high weight molecules.

$$c_m = c_B * e^{J/k} \tag{13}$$

$$J = \frac{TMP - \left[ A_1 * C_B * e^{\frac{J}{k}} + \left( A_2 * C_B * e^{\frac{J}{k}} \right)^2 + \cdots \right]}{\eta_0 * R_m} \tag{14}$$

However, although osmotic pressure models provide more information than the gel-layer models, their main limitation is the determination of their mass transfer coefficient $k$. On the one hand, even though its direct measurement from the experimental flux data is possible, it is more complicated and prone to errors than with the gel-layer models due to not considering $c_m$ constant. Contrary, $c_m$ has to be first determined for each $c_B$ by relating the measured fluxes to the membrane osmotic pressure expression (Eq. (11) and (12)). Afterwards, the fluxes are plotted over the logarithmic ratio between $c_m$ and $c_B$ (film theory, Eq. (3)) and the mass transfer coefficient $k$ can be obtained from the slope of the linear regression[19]. By extrapolating the flux to 0, a maximum solute concentration at the membrane surface, $c_{max}$, is also drawn, which corresponds to the maximum solute concentration that could be reached in the bulk, when the osmotic pressure in the membrane would equal the applied TMP. Hence, by determining $k$ in this way, any error when selecting the expression form that relates the osmotic pressure with $c_m$ -Eq. (12), which is only true for completely rejected solutes- or the virial coefficients in it – which are strongly dependent on the pH, the solvent ionic strength as well as on the protein size and protein-protein interactions[16]- will result in $k$ uncertainties. These uncertainties, due to how $k$ is arranged in Eq. (14), will subsequently lead to larger errors in the flux predictions. Summarizing, whereas accurate enough gel-layer models can be constructed with just recording the flux and the solute concentration in the bulk from a set of training experiments, more measurements and calculations are necessary for osmotic pressure models [16],[19]. However, when they are accomplished, more precise and reliable models able to better describe the flux over wider ranges of process conditions are obtained.

### 1.2.1.2 Other flux limiting phenomena

Finally, in addition to concentration polarization, several mechanistic models have also been developed to explain other flux limiting phenomena such as pore-blocking or cake formation in applications such as virus[22], DNA[23], polymers[24] or proteins[25] filtration. However, the main phenomenon governing the flux evolution in most of these processes was also changing over time -which is specially common in fouling processes, where the combination of different component interactions and fouling mechanisms takes place[26]- therefore making difficult to build valid models able to describe the entire process.

Moreover, it is important to highlight that for all the previously explained mechanistic models the mass transfer coefficient $k$ was assumed to be constant (1.2.1.1 Concentration polarization) -and thereby it was also the diffusivity, viscosity and density-, and that therefore,

they were independent of the protein concentration. Although this assumption has shown to be true for many boundary conditions and specific system geometries[12],[13],[14], it is not the case for high concentrated protein solutions[27], especially of high molecular weight, such as monoclonal antibodies[28]. In these solutions, the diffusion coefficient and viscosity strongly change along the boundary layer as a result of significant protein-protein short-range interactions, further complicating the construction of mechanistic models under these conditions. Nevertheless, some approaches have been developed to overcome these limitations, such as introducing correction factors for the viscosity and diffusivity concentration dependencies[12],[29], or modifying the stagnant film theory in more complex models to account for the impact of the intermolecular protein interactions in the thermodynamic -protein diffusivity, related by the chemical potential gradient to the virial coefficients- and hydrodynamic -viscosity- properties[27]. These model were even extended to incorporate the influence of the buffer conditions[28],[30] and the membrane resistance[20]. However, although they showed very good results describing the flux, both in the pressure-dependent and independent regions, these models relied on the experimental measurement of rather physical properties such as the protein diffusivity (by dynamic light scattering), osmotic pressure (by its relation to TMP in a stirred cell device) or the concentration-dependent viscosity (by capillary viscometry), as well as on an extensive data fitting. Finally, they only considered one component in the solution, and therefore the presence of additional components would enormously increase their complexity.

It has been extensively reported in literature that the presence of more than one component can affect both the flux and the membrane selectivity in a filtration system, due to the influence of protein-protein interactions in the mass transport coefficients, from weak[31] and moderate[32],[33],[34],[35], to strong[14],[36], depending on several factors such as the proteins' charge and concentration or the pH and ionic strength of the solution. Similar results have been reported for chromatographic steps[37],[3]. Nevertheless, the most common assumption in modeling of downstream processes is that the overall sample composition is reduced to the target molecule only, so the determination of the mechanistic coefficients and parameters is easier as well as their model structure. For some process units, such as polishing chromatography[38],[39] or ultra/diafiltration before formulation[40], this assumption is realistic, since the impurity content in the treated samples is low. For earlier purification steps, however, such as filtration after capture or between polishing steps or capture affinity chromatography[37], this simplification can lead to erroneous models, since the neglected presence of process-related impurities such as host-cell proteins, DNA or protein aggregates can significantly distort the performance of the purification unit. If these impurities are quantified, such effects can be taken into account. In the case of filtration processes, some mechanistic models have been developed that include the impact of different components during virus[22] or protein dead-end[36]

and cross-flow[41] ultrafiltration. However, the construction of such models is complicated, since in addition to all the labor- and material-intensive experimental work necessary to determine the mechanistic coefficients including all the components, the effects of the impurities are usually complex and not always completely understood, which typically leads to make some assumptions when using these models. If the overall behavior of the process changes due to variations in the solution's composition, these assumptions might not be true anymore, losing model predictability - as mentioned in section *1.2.1 Mechanistic modeling*, different phenomena can influence the flux during a filtration process, and it is therefore necessary to select first which is the dominating phenomenon before determining its coefficients and parameters. However, the selected model might only predominate under certain specific conditions. Finally, different components might have different permeabilities in the membrane, and therefore, in order to include them in the model, it is necessary to first calculate their rejection factor. The rejection factor of a solute is commonly set equal to a constant value, which is empirically calculated from the obtained final concentration in the bulk and permeate solutions. Nevertheless, it has been reported that this value is not constant over time and that it is additionally influenced by the applied process conditions[33],[35]. A mechanistic dynamic description of the protein rejection factor as a function of the concentration of two proteins in cross-flow ultrafiltration was attempted once[36], however, it required the labor measurement of hydrodynamic and thermodynamic parameters from the solution.

## 1.2.2 Data-driven models

An alternative to mechanistic models are data-driven (also known as statistical or non-parametric) models, which capture the data underlying relationships between the variables (inputs and outputs) by mathematical functions, without explicit knowledge of the physical behavior of the system. These models are therefore fast and easy to apply, however, they provide little insight into the functioning of the system and have less extrapolation predictabilities compared to mechanistic models[42]. Extrapolation is defined as the estimation of a variable based on other variables that are outside the original observation range (also known as training data or space, i.e., the data used to construct the model). A model with good extrapolation capabilities will therefore make accurate predictions beyond the training conditions. Due to the structure of these models is inferred from the data, they need larger and more representative amounts of them compared to mechanistic models.

To establish the relationships between variables, data-driven models make use of machine learning, which is an area of computer science that focuses on studying how computer algorithms automatically improve their performance ("learn") through experience[43]. Therefore, based on a set of observations, data-driven models learn the relationship between input and output parameters. After the training, these models can make accurate predictions of new data

(test data)[43]. There are two main types of machine learning algorithms: supervised and unsupervised. Supervised learning algorithms -the ones used in this work- are those where the modeler predefines the causality relationship (input and output variables) in the training data, so the supervised algorithm learns the mapping function between them. Supervised learning algorithms include tasks such as active learning, classification (categorical output variables) and regression (continuous output variables)[44]. On the contrary, unsupervised learning is where no relation between input and output data is given by the user and the algorithm finds the structure between them (by clustering or association tasks).

Within regression, that is the identification of a mathematical expression that relates the data with the least error, different models are used in machine learning, mostly classified as linear, nonlinear and regression trees. Linear regression -simple or multiple, ordinary or partial least squares- models offer good interpretabilities, are simple to construct and require less computational time. However, they might have some limitations such as when defining nonlinear relationships between predictors and response variables, when there are many different predictors, or in cases of small observation points or multicollinearity[45]. Although linear models can be adapted to nonlinear trends by adding model terms (quadratic, logarithmic…), it is necessary to know the nonlinearity of the data. On the opposite, this is not necessary when using nonlinear regression models, such as artificial neural network (ANN), multivariate adaptive regression splines (MARS), support vector machines (SVMs) or K-nearest neighbors (KNNs). Finally, regression trees partition the data into smaller groups based on logic statements, and they are mostly used in other applications such as in decision making.

Regarding the non-linear artificial neural network regression models, due to their increasing popularity as well as their extensive utilization in this work, they are shortly explained in the section below.

### 1.2.2.1 Artificial neural network (ANN)

Artificial neural networks are nonlinear machine learning algorithms that recognize the underlaying relationships in a set of data through a process inspired by the way the human brain works. Particularly, they consist of a collection of interconnected hidden units or nodes, also known as neurons, each of which receives several inputs, calculates the weighted sum of them (linear activation function, Eq. (24)) and transform it by a non-linear activation function (Eq. (25)). Afterwards, these neurons pass the result ("signal") to the next unit, which receives them by also a linear combination of all the previous units or neurons. By combining linear and non-linear activation functions these models are able to establish more complicated input-output relationships. In addition, the attributed weights allow to increase or decrease the strength of a variable or neuron under a certain condition, gaining flexibility. The structure of the ANN model used in this work for flux prediction is shown in Figure 17.

The neurons are organized in hidden layers, which might have different activation functions and each of them have a certain bias value that multiplies the linear activation function of all the neurons that form that layer. The term hidden refers to the fact that these units are unobserved variables, although they modulate the outcome of the model. Hence, in an ANN model, the signal travels from the first layer of input variables to the last layer of outputs, with the different hidden layers in between.

The value of each weight and bias is optimized (trained) to minimize the model error, which is typically the sum of the squared residuals. This process usually starts with randomly choosing the weights and biases and making the first prediction with the inputs provided in the training set. By comparing this prediction with the outputs from the training, an error is calculated and the ANN backward calculated to fit the prediction. The error between real and backward calculated inputs is used each time to update the weights and biases. The optimization can be performed by different algorithms. The learning process is complete whenever the error is not usefully reduced anymore. The established model structure is finally optimized by applying the model to a different data set, called validation, to fine-tune the structure parameters, e.g., the number of nodes Therefore, the number and nature of the parameters in these regression models are not fixed *a priori*, but flexible and determined from the data.

### 1.2.3 Hybrid modeling

Hybrid modeling (also called semi-parametric modeling) combines the benefits of both modeling approaches. They consist of a data-driven non-parametric part whose structure and parameters are inferred from the data, and a mechanistic described parametric part. These parts are hereafter referred to as black box and a white box, respectively. The black box correlates the input with output variables using the weights/coefficients of a nonparametric function - e.g., artificial neural network (ANN), partial least squares regression (PLS), multivariate adaptive regression splines (MARS), etc[45]. On the other hand, the white box represents a mechanistically or phenomenologically well-understood function of fixed structure typically derived from conservation, thermodynamic or kinetic laws. Therefore, the parameters of the white box have a physical meaning. The black and white boxes can be combined in multiple ways in a hybrid model, mostly as serial or parallel. Parallel hybrid models are typically used when full fundamental mechanistic models are available, but their predictions should be improved (corrected) with data-driven models to fit the process data[42]. On the other hand, parallel hybrid models are used when some parts of the process are not well understood by the available mechanistic knowledge, and data-driven models are then used to describe them. In this case, the black box can proceed the white box or vice versa, depending on the degree of available knowledge about the process. In some cases, a combination of parallel and serial hybrid models is applied.

Hybrid models have been successfully used for modeling of bioprocesses, both in upstream[46],[47] and downstream[48] applications. In the case of filtration, they have been used to describe the fouling phenomenon[49] and different cleaning strategies[50] in dead-end wastewater ultrafiltration. Recently, a hybrid model for flux prediction of protein solutions in cross-flow ultrafiltration was also developed, highlighting its benefits compared to the mechanistic stagnant film theory[51]. However, this model only accounted for single-protein solutions, which can substantially deviate from reality in some of the purification stages where ultrafiltration is usually used in the biopharmaceutical industry, such as after capture chromatography or between polishing steps, where the still high presence of host-cell impurities and contaminants can strongly affect the filtration performance.

## 1.3 Host cell impurities

### 1.3.1 Typical *Escherichia coli* impurities

Due to its rapid growth, high yield of product and cost-effectiveness, *Escherichia coli* is usually the first-choice microorganism for the production of heterologous proteins including biopharmaceuticals, being about 30% of the current approved therapeutic proteins produced by these host cells[52]. However, due to the presence of both cytoplasmic and outer membrane as well as cell wall, the heterologously expressed proteins are not efficiently secreted to the extracellular space by *E.coli*, making cell lysis usually necessary[53]. This leads to the release of several host cell-derived impurities such as endotoxins, dsDNA or host cell proteins (HCPs) into the supernatant in excess along with the protein of interest. Moreover, in *E. coli* the desired product is, if expressed at high rates, often deposited in a non-soluble form known as inclusion bodies, which have to be re-solubilized afterwards. They are only to minor extend associated with host cell proteins such as elongation factors[54], subunits of RNA polymerase[55], outer membrane proteins[56] or proteins conferring antibiotic resistance[57]. Hence, the starting material for the purification of recombinant proteins from *E. coli* is complex, which makes multistep downstream processes essential to accomplish with the required limits for each of these impurities established by the authorities. An overview of the dominant process-related impurities which have to be considered during overexpression in *E.coli* as well as their imposed regulatorily limits are briefly summarized below.

Endotoxins, also known as lipopolysaccharides (LPS), which are integral components of the outer cell membrane of gram-negative bacteria, are extremely toxic when getting into the human blood stream. Even small concentrations might lead to systemic inflammatory reaction, tissue injury or even the death. For this reason, the current maximum amount per dose for intravenous application products is limited to 5 endotoxin units (EU, related to the biological activity of the endotoxin, usually between ~ 120-200 pg of endotoxin depending on the type)

per kilogram body weight and hour[58]. DNA is also considered a potentially dangerous substance, due to the possibility of cellular transformation by potentially oncogenic DNA. A maximum of 10 ng DNA per dose is limited for intravenous application[59] and 100 µg for orally administrated vaccines[60]. Host cell proteins (HCPs), are a very common impurity found in every host cell line used in bioproduction, and are usually considered critical process parameters due to their potential immunogenicity as well as their adjuvant, proteolytic or direct biological activity[61]. The range of HCPs species in a fermentation broth is very extensive. There is nearly an infinite variety of HCPs and they strongly depend on the cell type and expression modality used[62] as well as on the product[63] and production process itself[64]. It is for this reason that there is not yet neither an absolute control limit established by the regulatory agencies nor an standard quantification method assay regarding HCPs. On the contrary, their acceptance level is usually reviewed on a case-by-case basis by the authorities depending on several aspects (maximum dose of the drug, route of administration, frequency of dosing, pre-clinical and clinical data…). However, it has been reported that most biologic products reviewed by the FDA up to 2004 had HCPs contents below 100 ppm or mg/L[65]. Hence, the imposed limits by the regulatory agencies to ensure the safety of a medicament are quite stringent, especially if considering the typical impurity concentrations of an *E. coli* homogenate – about 1 g/L dsDNA, 20 g/L $HCP_{Total}$ and $1 \cdot 10^6$ EU/mL for a 2 g/L cytoplasmic expressed protein solution. After capture steps, however, these amounts are typically reduced to 0.5-1.5 mg/L dsDNA, 0.5 mg/L $HCP_{Total}$ and 5000-100,000 EU/mL, while the target protein being concentrated to about 10 g/L[66].

### 1.3.2 Chinese hamster ovary (CHO) cell impurities

Mammalian cell lines are also very important host cells for the industrial production of recombinant proteins due to their ability for correct folding, assembly and post-translational modifications of proteins, especially glycosylation, for which bacteria lack the required intracellular machinery[67]. Among them, Chinese hamster (*Cricetulus griseus*) ovary (CHO) cells are the most used cell types, being about 70% of the current clinical recombinant products produced by them[61]. From all biopharmaceuticals, monoclonal antibodies (mAb) are the ones that have been clearly dominating the pharmaceutical market in the last recent years and whose fast growth is expected to keep growing in the future [61],[6]. These molecules are secreted by the cells into the cell culture supernatant, together with many other host proteins, from which they have to be afterwards purified to meet with the required product specifications. Recently, it was reported that CHO cells secrete hundreds of different HCP species into the cell culture supernatant, with a typical range of 300 mg HCP/g mAb[63]. Furthermore, some of these HCPs were described as exceptionally difficult to remove during downstream processes[64], due to either having similar physiochemical properties to the product, strongly interacting with it or

due to being retained on the chromatographic media[63]. The latter, in Protein A chromatography -the most significant step for HCP removal (>90%[37]) in mAbs and Fc-fusion carrying proteins purification processes- was reported to negatively affect its longevity due to the binding of impurities to the Protein A ligand[63],[37]. Additionally, it was described that the primary and secondary clarification steps influenced the Protein A performance by regulating the HCP profile that went through it[3],[37]. Hence, having a control of the content of the impurities throughout the downstream processing units is of high interest in the biopharmaceutical industry, not only for ensuring the safety and quality of the final product and reduce the risk of batch rejection, but also for the efficiency and economic viability of the process[62], since downstream processes represent the major part of the production costs of biologicals.

# 2. Objective

The aim of the present study was to extend the previously developed one-component-hybrid model[51] to account for a second component in the system, lysozyme, which mimicked a process-related impurity, and to evaluate its influence on the flux together with BSA, mimicking the target product. The used ratios between the two model proteins were based on the literature mentioned in *1.3 Host cell impurities* about the typical product and impurity titers in downstream processing solutions. BSA and lysozyme were selected as model proteins because they exhibit different physicochemical properties that facilitate their separation and quantification. The advantage of using hybrid models for two-component UF systems is that the complex and variable effects that each of the components might have on the flux decline do not need to be experimentally quantified, which would be time- and material-consuming, but instead, this is done by machine learning algorithms, by establishing mathematical relationships between inputs and outputs, thereby drawing the most information out of the available data. The classical stagnant film model and the previously established one-component-hybrid model were also reproduced to compare for the flux prediction.

Additionally, due to lysozyme was only partially retained by the selected membrane, the data-driven black box models were also used for predicting its rejection factor ($R_{Lys}$) evolution. For this purpose, three overall different two-component-hybrid model (tcHM) structures were developed, depending on how $R_{Lys}$ was calculated. The structures ranged from static $R_{Lys}$ values (tcHM1), to dynamically updating $R_{Lys}$ in the same or different black box with flux, tcHM2 and tcHM3, respectively. The white box model in turn consisted of a mass balance that calculated the future protein concentrations and remaining reservoir volume based on the predictions from the black box models. Hence, the presented hybrid models aimed, from just the initial concentration of each protein and the transmembrane pressure and cross-flow velocity parameters, to make accurate predictions of both the flux and the $R_{Lys}$ over time. These predictions were subsequently used to calculate the expected final concentration of each component in the bulk and compared to the measured concentrations.

# 3. Materials and Methods

## 3.1 Experimental setup

### 3.1.1 Equipment and chemicals

All ultrafiltration runs were performed on an ÄKTA Crossflow system (Cytiva, Columbia, USA) controlled by UNICORN 5.31 software. The reservoir tank had a maximum volume of 1100 mL and the system featured online pH, temperature, UV-absorbance and conductivity sensors on the permeate side as well as a pressure-based reservoir level sensor. The experiments were performed with a Sartocon Slice hydrophilic, stabilized cellulose-based membrane (Hydrosart) cassette (Sartorius AG, Göttingen, Germany) with a membrane area of 200 cm$^2$. The chosen pore size for the membrane was 30 kDa, so that the mimicked protein of interest -BSA, 66 kDa- was fully and the mimicked impurity -lysozyme, 14 kDa- partially retained. Both model proteins were purchased from Sigma-Aldrich, St. Louis, MO, USA. The filtration buffer used was 50 mM PBS, pH 8.

Before all experiments were performed, the filters were flushed with 1 L deionized water and conditioned with filtration buffer until a stable conductivity signal was reached. After the experiments, the filters were flushed with 1 L deionized water, cleaned for 2 hours with 1 M NaOH clean-in-place (CIP) solution and flushed again with 1 L deionized water. Finally, before and after each experiment, a clean water permeability (CWP) test was made in order to ensure similar membrane permeability properties and thereby the reproducibility of the experiments. All test runs performed within less than a 9% Flux/TMP variation compared to the new membrane.

### 3.1.2 Training and test data generation

The data for training the black box models were collected using the UNICORN built-in process optimization tool. The schematic depiction of the cross-flow filtration unit is shown in Figure 4. During training data generation, for each $c_{B,i}$ the CF was first adjusted to the lowest flow rate by employing the feed pump ($P_F$), followed by stepwise TMP variations from the lowest to the highest value using the valve on the retentate-side ($P_R$). After the TMP variations, the CF was also increased stepwise (Figure 4c). During the scouting, the permeate was redirected into the reservoir in order to keep $c_{B,i}$ constant (Figure 4a). For each $c_{B,i}$, three CFs (100, 200 and 300 mL/min) and five TMPs (0.8, 1.3, 1.8, 2.3 and 2.8 bar) were tested, and the resulting fluxes and UV-absorbances were recorded. The UV-absorbances were afterwards used to calculate the lysozyme rejection factor, $R_{Lys}$, through a calibration curve (see section *3.1.3 UV absorbance-cp,Lys calibration curve*) and fed to the black box together with the fluxes for model training. At the beginning of each scouting round a sample was taken from the reservoir and

analyzed by SEC-HPLC. The obtained $c_{B,i}$s were assumed to be constant during all the scouting round. To increase the $c_{B,i}$ between rounds, the samples were concentrated (Figure 4b) until the desired bulk volumes were reached (Figure 4d). At the beginning of each of these concentration steps, a sample was taken from the permeate side and its concentration, $c_{p,i}$, related to the $c_{B,i}$ from the previous round in order to calculate the $R_i$ factor. With all the calculated $R_i$ values at each concentration step, a training $R_{i,average}$ factor was obtained (*4.1.2.1 RLys average from training set*). The cross-flow and transmembrane pressure during the concentration steps were chosen to be in the middle of the trained process parameters space (Figure 4c and Figure 13), 200 mL/min and 2 bar, respectively.

In total three training experiments were performed. The first one included only BSA, in order to test the reproducibility of the previously developed one-component-hybrid model[51] in the current membrane and its performance for predicting two-component solutions. The second one, with only lysozyme, was used to generate the calibration curve between $c_{p,Lys}$-UV absorbance. Moreover the second training set was used to see how pure lysozyme behaved with the membrane without BSA. Finally, a third training experiment combining both proteins, with the purpose of building the models that simultaneously predicted the flux, rejection factor and concentration of two-protein solutions was performed. A summary of all training experiments is provided in Table 1.

Table 1: Summary of all performed training experiments, with the number of scouted $c_B$s, their training set size and the $c_B$s range for the two proteins.

| Training set | nº scouted $c_B$s | Training set size | $c_{B,BSA}$s range [g/L] | $c_{B,Lys}$s range [g/L] |
|---|---|---|---|---|
| BSA alone | 10 | 150 | 3.7-180 | - |
| Lysozyme alone | 12 | 178 | - | 0.39-43.9 |
| Combined (BSA + Lys) | 6 | 90 | 3.8-77.9 | 0.28-3.81 |

Figure 4: Schematic overview of the setup of the cross-flow filtration system. (a) and (c): TMP and CF scouting rounds. The permeate was redirected into the bulk reservoir in order to keep $c_B$ constant. (b) and (d): concentration steps of the training experiment, the reservoir volume was reduced in order to increase $c_B$. This setup was the same for the test sets. Figure from [51].

As it can be seen in Figure 4c, every time before the TMP started to increase stepwise as well as at the end of each scouting round, the solution was recirculated with the permeate valve closed (TMP=0) during some minutes. The reason for this was to reduce the CP layer on the membrane surface, which would lead to smaller fluxes and therefore to flux underestimations when building the models, since $c_{B,i}$s had to be constant for all the scouting round. However, due to concentration polarization is an inherent phenomenon of all pressure-driven filtration processes, its effects should also be considered. This was accomplished by using the data from the concentration steps between scouting rounds (Figure 4d, *3.1.5 Concentration polarization correction*). Hence, in this way it was possible to both build precise models for flux prediction and correct for the protein accumulation in the membrane due to concentration polarization by just taking samples at the beginning of each scouting round.

After the data collection for training the black box models, UF test runs with varying initial $c_{B,BSA}$, $c_{B,Lys}$, CF and TMP were performed (Table 2). The protein solutions, with a volume of > 1000 mL, were concentrated at constant TMP and CF, and the decrease in flux over the process time was recorded and compared to the predictions from the hybrid models. In addition, several samples were also taken during the filtration process both from the retentate and permeate sides to calculate $R_i$ evolution, which were also compared with the model predictions. The exact initial and final reservoir volume as well as the number of samples taken at each test set are shown in Table A1.

The predictive performance of the models was evaluated by the normalized root-mean-square error (NRMSE), comparing predicted and measured flux and $R_{Lys}$ values as shown in Eq. (15).

$$NRMSE\ [\%] = 100 \cdot \frac{\sqrt{\frac{1}{n}\sum_{i=1}^{n}(y_i - \hat{y}_i)^2}}{y_{max} - y_{min}} \qquad (15)$$

where $n$ is the number of overserved values $y_i$, and $\hat{y}_i$ the corresponding predicted ones. The normalization $y_{max} - y_{min}$ allowed for a fair comparison of fluxes at different conditions. The observed permeate flux was recorded online in regular time intervals. The observed $R_{Lys}$ values were present in irregular time intervals, since they were based on offline measurements.

Table 2: Summary of all performed test sets with varying TMP, CF, $c_{B,BSA}$ and $c_{B,Lys}$ and the ratio between the mimicked protein of interest (BSA) and the impurity (lysozyme), simulating the purity degree of the solution.

| Test set number | TMP | CF | $c_0$BSA [g/L] | $c_0$Lys [g/L] | Ratio (%) |
|---|---|---|---|---|---|
| 1 | 1.8 | 200 | 6.68 | 0.00 | 100.0 |
| 2 | 1.8 | 200 | 4.00 | 0.28 | 93.5 |
| 3 | 2.8 | 300 | 3.82 | 0.32 | 92.3 |
| 4 | 2.1 | 250 | 3.71 | 0.38 | 90.7 |
| 5 | 2.5 | 280 | 4.56 | 0.25 | 94.8 |
| 6 | 1.8 | 200 | 3.79 | 0.50 | 88.3 |
| 7 | 1.6 | 230 | 5.97 | 0.15 | 97.5 |
| 8 | 1.4 | 270 | 8.80 | 0.19 | 97.9 |
| 9 | 1.8 | 260 | 2.38 | 0.57 | 80.5 |
| 10 | 2.0 | 350 | 3.62 | 0.34 | 91.3 |

After training and test runs, the permeate valve was closed and the remaining solution was recirculated at 300 mL/min cross-flow during about five minutes. Afterwards, the reservoir was emptied and the solution weighted and sampled in order to close the mass balance. Later on, fresh buffer (50-100 mL) was added to the system and recirculated during approximately half an hour at very high cross-flow rates (500-550 mL/min). The reservoir was then emptied again and the concentrations measured to evaluate if the missing protein in the mass balance could be recovered.

### 3.1.3 UV absorbance-$c_{p,Lys}$ calibration curve

The calibration curve between the UV absorbance at 280 nm and $c_{p,Lys}$ was built during the training experiment with only lysozyme. The linear regression curve showed an $R^2$ of 0.998 and was afterwards used to correlate all the recorded absorbance values at each TMP and CF equilibrium combination during the BSA with lysozyme training experiment, to their $c_{p,Lys}$ and therefore $R_{Lys}$ – in combination with the $c_{B,Lys}$s from the previous scouting round. The regression curve is shown in Figure 11.

Due to during the training set with only BSA very small amounts of protein were detected on the permeate side -≤ 0.05 g/L, below the lowest point of the standard calibration curve in SEC-

HPLC-, the BSA rejection factor, $R_{BSA}$, was assumed to be 1. Consequently, all the UV absorbance values recorded during the BSA with lysozyme training experiment were assumed to be only from lysozyme.

### 3.1.4 Protein analysis

BSA and lysozyme concentrations were determined by an analytical high-performance size exclusion chromatography (SEC-HPLC) using a TSKgel G3000SWXL column (5 μm, 7.8 × 300 mm; TOSOH, Shiba, Tokyo, Japan). The separation was performed under isocratic conditions with 50 mM sodium phosphate and 200 mM NaCl at pH 6.5 as running buffer, at a flow rate of 0.4 mL/min. The samples were previously diluted with the same buffer to a final concentration of 0.08 to 1.5 g/L for BSA and 0.08 to 1 g/L for lysozyme, and filtered through a 0.22 μm Millex-GV filter (Merck Millipore, USA). The injection volume was 10 μL per sample. Due to their difference in size, BSA and lysozyme peaks could be completely separated and quantified independently using the standard calibrations for BSA and lysozyme made from the stock solutions. The mean elution time for BSA was 23.5 min, while for lysozyme it was 30 min.

### 3.1.5 Concentration polarization correction

During the concentration steps of the training experiment (Figure 4d, TMP 2 bar and CF 200 mL/min), the measured $c_{B,BSA}$ at the beginning of each scouting round was lower than expected if considering the filtered volumes. This difference -that increased with $c_{B,BSA}$- between expected and measured concentrations was due to the formation of a CP layer on the membrane. This phenomenon was accounted by the models by introducing a quadratic polynomial function to the calculated $c_{B,BSA}$s from the training data. In the case of lysozyme, no correction function was introduced due to the accumulated amount of protein observed at the end of the training set was very small.

## 3.2 Hybrid modeling

### 3.2.1 Black box models

The black box part of the hybrid models aimed to predict the flux (output parameter) for a combination of input parameters (TMP, CF, $c_{B,BSA}$ and $c_{B,Lys}$). For that, artificial neural network (ANN) was used as the black box model, which was set up using the *feedforwardnet* function and trained with *trainbr*, a function that uses Bayesian regularization backpropagation to avoid overfitting by minimizing the combination of squared errors and weights, thereby penalizing the large weights (weight decay). All computations were performed using MATLAB 2018b software. The ANN was optimized by varying the number of neurons from 1 to 6 in one hidden layer. The inputs and outputs were normalized between 1 and 2 for processing in the black

box. The black box models were constrained by setting the flux to 0 when the TMP was below 0.1 bar - e.g., at the beginning of the process. In this TMP range, the fluxes are neglectable. In addition, the predicted fluxes in this range are only based on extrapolations, rending them prone to errors.

Moreover, since lysozyme was only partially retained by the membrane, the black box models were also used for predicting $R_{Lys}$ evolution in tcHM2 and tcHM3 - in tcHM1, $R_{Lys}$ was set equal to a certain value from the training set. On the one hand, in tcHM2 $R_{Lys}$ was introduced in the same black box with flux, thereby giving rise to a hybrid model with one black box of four inputs (TMP, CF, $c_{B,BSA}$ and $c_{B,Lys}$) and two outputs (flux and $R_{Lys}$). On the other hand, in tcHM3 an additional black box was introduced entirely for $R_{Lys}$ prediction. Different black box models were utilized for this purpose: ANN, MLR and MnLR. The ANN model for $R_{Lys}$ prediction was optimized by varying the number of hidden neurons between 1 and 4. MLR, stating for multiple linear regression, was optimized by the *stepwiselm* function, which uses stepwise regression for adding or removing predictors from the model based on a certain statistical criterion. In this case, the model was both forwards and backwards stepwise optimized both by the p- and AIC values for the four input parameters, allowing linear and interaction terms. The MnLR black box model, which is a MLR model with including a logarithmic term for $c_{B,Lys}$, was forward stepwise optimized by the p-value criterion. In MLR and MnLR, the terms which exhibited a p-value higher than 0.05 were excluded from the model one after the other.

### 3.2.2 White box model

The white box model is the mechanistic part of a hybrid model. In this case, it used the predicted fluxes and $R_{Lys}$ from the black boxes to calculate the future proteins' concentrations, both in the bulk and permeate, $c_{B,i}$ and $c_{p,i}$, as well as the remaining reservoir volume, $V_B$, after a certain time interval *dt* by using the mass balances, where A is the membrane area.

$$\frac{dV_p}{dt} = -\frac{dV_B}{dt} = J \cdot A \tag{16}$$

$$\frac{d(c_{B,i} \cdot V_B)}{dt} = \left(A \cdot J \cdot c_{B,i}\right)^{R_i} \tag{17}$$

$$c_{P,i} = (1 - R_i) \cdot c_{B,i} \tag{18}$$

### 3.2.3 Multi-step ahead hybrid models

In the present work, different hybrid model structures were investigated in order to describe the flux and the rejection factor evolution of two-protein solutions. All models consisted of a serial hybrid model structure, where the black box predicted first the change in flux for a combination of input parameters (CF, TMP, $c_{B,BSA}$ and $c_{B,Lys}$), which was fed to the white box together with $R_{Lys}$ for calculating $c_{B,i}$ and $V_B$ after a given time interval *dt*. The predicted $c_{B,BSA}$ and $c_{B,Lys}$ from the white box were fed back to the black box for future flux and $c_{B,i}$ predictions

in a multi-step ahead structure (Figure 7B). Multiple iterations were performed until a desired stop criterion was reached - in this case the final retentate volume from the test sets, $V_{B\_final}$ (Eq. (19), Table A1).

$$V_B \geq V_{B\_final} \tag{19}$$

Depending on the way $R_{Lys}$ was calculated, the models were classified in three overall structure types: tcHM1, tcHM2 and tcHM3, which are all summarized in Table 3. These models were further subdivided in different submodels depending on which data they were trained on and the black box model type used for $R_{Lys}$ prediction. The well-known mechanistic stagnant film theory model (SFM) and the previously established one-component-hybrid model (ocHM) were also reproduced for comparison of flux prediction. Before feeding $c_{B,BSA}$ back to all models for the next iteration step, it was corrected by using Eq. (23) for protein accumulation on the membrane as a result of concentration polarization.

Hence, in addition to the multi-step ahead flux and $R_{Lys}$ predictions, the presented hybrid models also yielded time-resolved $c_{B,BSA}$ and $c_{B,Lys}$ forecasts. The predictions of final $c_{B,BSA}$ and $c_{B,Lys}$ were compared to the measured concentrations by SEC-HPLC to assess the performance of the models.

Table 3: Summary of all submodels derived from the one- and two-component-hybrid models, ocHM and tcHM, respectively, and the stagnant film model (SFM) structures, depending on the used training data and the black box model for $R_{Lys}$ prediction.

| Model structure | Submodels | Explanation |
|---|---|---|
| ocHM | ocHM$_{BSA}$ | One-component-hybrid model based on BSA alone training data |
| | ocHM$_{comb}$ | One-component-hybrid model based on combined training data (neglecting lysozyme) |
| tcHM1 | tcHM1$_{R1}$ | Two-component-hybrid model, $R_{Lys}$ set to 1 |
| | tcHM1$_{Raverage}$ | $R_{Lys}=R_{average}$ from training data |
| tcHM2 | tcHM2 | One multi-output black box for J and R prediction |
| tcHM3 | tcHM3$_{ANN}$ | Neural network black box for $R_{Lys}$ |
| | tcHM3$_{MLR}$ | Multiple linear regression black box for $R_{Lys}$ |
| | tcHM3$_{MnLR}$ | Multiple non-linear regression black box for $R_{Lys}$ |
| SFM | SFM$_{BSA}$ | Stagnant film model based on BSA alone training |
| | SFM$_{comb}$ | Stagnant film model based on combined training but neglecting lysozyme concentration |

### 3.2.3.1 One-component-hybrid model

In order to analyze the impact that the mimicked impurity, lysozyme, had on the flux evolution compared to only BSA, the one-component-hybrid model developed by Krippl et al.[51] was first reproduced and trained on the BSA alone training data (ocHM$_{BSA}$). Moreover, due to a common assumption in modeling of downstream processes is to neglect the process-related impurities, either because they do not significantly influence the performance of the purification unit, they are not CQAs or just for simplifying the models, the ocHM was also trained on the BSA with lysozyme training data but neglecting the lysozyme concentration (ocHM$_{comb}$). The aim of this model was therefore to investigate how precise the flux predictions could be without quantifying the impurity in the solution. Furthermore, this model also allowed for a fairer comparison strictly of flux prediction with the stagnant film model (SFM), since the latter is only suited for one-component systems. The structure of the reproduced one-component-hybrid model, consisting of one black box with three inputs (TMP, CF and c$_B$) and one output (flux), can be consulted in [51]. In these models, the rejection factor of the sole component was set to 1 all the time.

### 3.2.3.2 Two-component-hybrid model 1 – static lysozyme rejection factor

In tcHM1, the ocHM1 structure was extended to introduce lysozyme as second component and its influence on the flux prediction. This was done by adding c$_{B,Lys}$ as additional input of its black box, as it is shown in Figure 5. The white box was also expanded with the mass balance for this second component.

In tcHM1, the R$_{Lys}$ was set equal to a certain fixed value. In the case of tcHM1$_{Raverage}$, it was used the average rejection factor from the training set, 0.77, which was calculated by doing trapezoid integration of all obtained R$_{Lys}$ and the accumulated permeate volume, V$_p$, during the concentration steps of the training experiment – Eq. (22) (see *4.1.2.1 RLys average from training set*). Furthermore, in order to evaluate the impact of calculating R$_{Lys}$ on flux and c$_{B,Lys}$, R$_{Lys}$ was set equal to 1 in tcHM1$_{R1}$.



Figure 5: Schematic representation of the two-component-hybrid model 1 (tcHM1) structure. The lysozyme rejection factor, R$_{Lys}$, was provided to the white box as a constant value.

### 3.2.3.3 Two-component-hybrid model 2 – dynamic lysozyme rejection factor with a multi-output black box

In tcHM2, $R_{Lys}$ was introduced as output in the black box together with the flux, with the objective of having a more precise description of the lysozyme rejection factor, which could be updated over time and depend on different process parameters. The tcHM2 structure is shown in Figure 6. This model was only trained on the BSA with lysozyme training data. $R_{Lys}$ predictions were constrained to be ≤1.



Figure 6: Schematic representation of the two-component-hybrid model 2 (tcHM2) structure with a multi-output black box for simultaneous prediction of flux and lysozyme rejection factor from TMP, CF, $c_{B,BSA}$ and $c_{B,Lys}$ input parameters.

### 3.2.3.4 Two-component-hybrid model 3 – dynamic lysozyme rejection factor with two independent black boxes

In tcHM3, a second black box was introduced entirely for $R_{Lys}$ calculation (Figure 7A). This was done in order to fully separate the training procedure for J and $R_{Lys}$ prediction and to optimize the optimal number of hidden neurons independently. In addition, it also allowed to test different black box functions for $R_{Lys}$. The $R_{Lys}$ ANN black box showed the lowest error for $R_{Lys}$ at one hidden neuron (Figure 22), suggesting a rather simple correlation. As an alternative to the ANN black box, a multiple linear regression (MLR) model was tested, which could provide simpler models with less computation times and easier interpretabilities.

Finally, due to the observed strong influence of $c_{B,Lys}$ over $R_{Lys}$ compared to the other input process parameters (Figure 12) and its clearly logarithmic-like trend (Figure 24), a logarithmic term for $c_{B,Lys}$ was introduced in the MLR model, yielding tcHM3$_{MnLR}$.

Figure 7: Schematic representation of (A) the two-component-hybrid model 3 (tcHM3) structure with two independent black boxes for flux and $R_{Lys}$ prediction from the same input parameters: TMP, CF, $c_{B,BSA}$ and $c_{B,Lys}$; and (B) the multi-step ahead hybrid model structure. This model calculated $c_{B,BSA}$ and $c_{B,Lys}$ iteratively for each time increment until Eq. (19) was fullfiled.

### 3.2.4 Stagnant film theory

The presented hybrid model structures were compared to the classical well-known mechanistic stagnant film model, SFM, which derives from the mass-transfer equations that describe the boundary layer next to the membrane for flux prediction (*1.2.1.1 Concentration polarization*). This model, which only accounts for one component in the system, predicted the flux depending on the BSA concentration in the bulk solution -which was between 4 and 46 times in excess compared to lysozyme-, by using Eq. (20) and assuming that BSA was completely retained by the membrane:

$$J = k \ln \left( \frac{c_G}{c_B} \right) \tag{20}$$

where $c_G$ is the gel-layer concentration at the membrane surface and $k$ the mass transfer coefficient, which mostly depends on the diffusion coefficient and the thickness of the boundary layer.

Therefore, for a given mass transfer conditions (TMP, cross-flow velocity, temperature…), the plot J vs ln($c_B$) will be a linear regression of slope $k$ and interception point with the abscissa $c_G$ (Figure 10 and Eq. (21)).

$$J = -k \cdot \ln(c_B) + k \cdot \ln(c_G) \tag{21}$$

Thus, different values of $k$ and $c_G$ were obtained for every TMP and CF combination from the training experiment data (Table A2). Hence, in SFM the black box from the hybrid models was replaced by Eq. (20) in order to predict the flux. In this case, however, the provided parameters

were $k$ and $c_G$ (instead of TMP and CF) and $c_{B,BSA}$ , as it is shown in Figure 8. For those test experiments with TMP and CF not covered by the training set, linear interpolation was made from the closest points. After flux prediction, the same white box as in the hybrid models was applied for $V_B$ and $c_{B,BSA}$ calculation for a certain interval $dt$. These models also used the multistep-ahead structure for several iteration predictions. The SFM was trained both on the BSA alone and the BSA with lysozyme training data, yielding $SFM_{BSA}$ and $SFM_{comb}$, respectively (Table 3).



Figure 8: (A) Schematic representation of the mechanistic stagnant film model (SFM) structure. The flux was calculated using Eq. (20) for the combination of $k$, $c_G$ and $c_{B,BSA}$ parameters. The flux was then applied to the same white box as in the hybrid models, but just for one component, which calculated $c_{B,BSA}$ and $V_B$ for a certain time interval $dt$. (B) Multi-step ahead capability.

# 4. Results and discussion

Many of the results obtained in this thesis have been published as a research article in *Processes* by Krippl et al.[68]. Here, some parts will be discussed in more detail than in the in the manuscript and the graphs that have been published are explicitly referred to the publication. The manuscript can be found in *9. Appendix* as an attachment.

## 4.1 Training and test data

### 4.1.1 Flux

A total of three training sets containing BSA, lysozyme, and a combination of both were generated, covering three CF (100, 200, and 300 mL/min), five TMPs (0.8, 1.3, 1.8, 2.3, and 2.8 bar) and different protein concentrations in the bulk, $c_{B,i}$ (Table 1). The recorded fluxes from each training experiment for the combination of $c_{B,i}$ and TMP at 200 mL/min CF are shown in Figure 9. The x-axis of Figure 9C and D were reduced for a better comparison between training sets. The entire graph with the full $c_{B,i}$ range is shown in Figure A1.



Figure 9: Training data sets including different protein concentrations and TMPs at CF 200 mL/min: two-component training set containing (A) BSA and (B) lysozyme in the same solution (blue). One-component solution of (C) BSA (red) and (D) lysozyme (green). Figure from [68].

The recorded fluxes versus the logarithmic concentration of BSA in the bulk for the BSA alone and BSA with lysozyme training experiments are also shown in Figure 10A and Figure 10B, respectively, at different TMPs and 200 mL/min CF. From the linear part of these plots (pressure-independent region), the mass transfer coefficient $k$ (negative slope) and the gel concentration $c_G$ (intercept with the abscissa) are drawn for the stagnant film models (SFM, Eq. (21)).



Figure 10: Permeate flux vs logarithmic $c_{B,BSA}$ used to estimate the mass transfer coefficient $k$ and the gel concentration $c_G$ of the stagnant film model (SFM), when based on the (A) BSA alone and (B) BSA with lysozyme training data. For the former, only the first concentration range is shown in order to allow for a better comparison of the x-axis. Fluxes recorded at 200 mL/min CF.

In general, increasing $c_{B,i}$ led to lower fluxes, while increasing TMP and CF led to higher fluxes in all training sets. This phenomenon is explained differently according to the theoretical model (*1.2.1.1 Concentration polarization*). On the one hand, according to resistance and osmotic pressure models, the increase in $c_{B,i}$ led to higher $c_m$, which increased the hydraulic and osmotic pressure resistances, respectively. These resistances decreased the effective pressure driving force along the membrane and thereby the flux. The increase rate of $c_m$ decreased with $c_{B,i}$, explaining the negative logarithmic shape seen in Figure 9. On the other hand, in line with the stagnant film model, which considers $c_m=c_G$, the decrease in flux with increased $c_{B,i}$ was due to a reduction in the logarithmic concentration driving force expressed in Eq. (10), as a result of a more prominent protein back diffusion effect along the concentration gradient. All models agree that an increase in TMP leads to higher fluxes by increasing the convective flow towards the membrane. Finally, higher CF decreased the thickness of the CP layer next to the membrane by rectangular displacement thereby increasing the permeate flux. When comparing the data between training sets, it was observed that the training set with only BSA (Figure 9C) exhibited higher fluxes compared to the training experiment with BSA and lysozyme (Figure 9A) - especially at TMPs above 1.3 bar, where the system reached the

pressure-independent region even at low $c_{B,BSA}$. This indicated that lysozyme influenced the flux of solvent through the membrane, either by interacting with BSA and/or with the membrane. The flux also decreased faster if comparing the two-component solution (Figure 9B) to the filtration with only lysozyme (Figure 9D). The reason was that being smaller than the pore size, lysozyme was only partially retained by the membrane, and therefore less accumulated at the boundary layer, giving rise to smaller resistances and subsequently higher fluxes. Contrary, when BSA was present, the $R_{Lys}$ increased notably -by almost a factor of 3, see *4.1.2.1 RLys average from training set*- and so did the lysozyme concentration in the bulk and the amount of protein accumulated on the membrane. This observation led to the assumption that lysozyme influenced the flux most probably by interacting with BSA rather than by membrane fouling – which would also have been shown in the training set with only lysozyme. Due to the isoelectric points (pI) of BSA and lysozyme were about 5 and 11, respectively, and the used buffer had a pH of 8, it was assumed that both proteins were interacting electrostatically, in accordance with some previously reported results[34]. However, it is important to mention that at higher lysozyme concentrations, a more pronounced flux decrease in the training set with only lysozyme was observed at high TMPs (see Figure A1), suggesting that fouling of the membrane by adsorption/pore blockage mechanisms also occurred at high lysozyme concentrations. Hence, both phenomena should be considered when building the models. The interactions between BSA and lysozyme most probably influenced the flux by affecting the mass transfer coefficients of the solution, as already shown for these two model proteins[21],[27],[36]. In addition, these interactions increased the lysozyme rejection factor, and thereby the probability of fouling to occur, since higher amounts of protein were retained on the membrane.

## 4.1.2 Lysozyme rejection factor

In addition to the flux, the UV absorbances on the permeate side were also recorded at each TMP and CF combination during the BSA with lysozyme training set. These values were afterwards related to the lysozyme concentration in the permeate, $c_{p,Lys}$, by using the calibration curve from the lysozyme alone training experiment, which is shown in Figure 11.

Figure 11: Calibration of UV absorbance at 280 nm on the permeate side versus $c_{p,Lys}$, measured by SEC-HPLC, in the lysozyme alone training set,.

The obtained $c_{p,Lys}$ were subsequently used to calculate $R_{Lys}$ together with the $c_{B,Lys}$, which were considered to be constant during each scouting round - and therefore independent of TMP and CF parameters. All the calculated $R_{Lys}$ from the training set, for each combination of $c_{B,BSA}$, $c_{B,Lys}$, TMP and CF, were fed to the black box models of hybrid models 2 and 3 (tcHM2 and tcHM3). The dependence of $R_{Lys}$ over these parameters is depicted in Figure 12.



Figure 12: Lysozyme rejection factor, $R_{Lys}$, values recorded during the BSA with lysozyme training set, for a combination of TMP and CF (y- and x-axis) and different $c_{B,BSA}$ and $c_{B,Lys}$ (different scouting rounds, colored legend).

As shown in Figure 12, the influence of TMP and CF on $R_{Lys}$ was smaller compared to the impact of $c_{B,BSA}$ and $c_{B,Lys}$, and it was mostly only noticeable at low solute concentrations. The main reason for this was that the influence of these two parameters was only recorded for $c_{p,Lys}$, but not for $c_{B,Lys}$. Consequently, the higher the protein concentration was, the smaller

was the effect of a variation in $c_{p,Lys}$ on $R_{Lys}$, since $c_{B,Lys} \gg c_{p,Lys}$. Further investigations on the influence of TMP and CF on $R_{Lys}$ would have required a more granular covering of the training space and additional measurements of $c_{B,Lys}$ at each TMP-CF combination. However, in order to not drastically exceed the number of offline measurements, it was assumed that TMP and CF did not influence $c_{B,Lys}$ and consequently only their influence on $c_{p,Lys}$ was measured. It is important to choose an equilibrium between model precision and complexity/laboriously.

### 4.1.2.1 $R_{Lys}$ average from training set

During the concentration steps of the training experiment -those between scouting rounds to increase $c_{B,i}$, Figure 4b and d-, samples were taken both from the retentate and permeate and analyzed by SEC-HPLC to calculate $R_{Lys}$ over the filtered volume. From these values, a weighted average lysozyme rejection factor from the training experiment, $R_{Lys,average}$, was obtained by doing trapezoid integration (Figure 13, Eq. (22)):



Figure 13: Evolution of the lysozyme rejection factor, $R_{Lys}$, over the permeate volume in the concentration steps of training experiment with (A) only lysozyme and (B) BSA with lysozyme.

$$R_{Lys,average} = \int_0^{V_p} R_{Lys} \, dV_p \tag{22}$$

The calculated $R_{Lys,average}$ for the two training sets containing lysozyme is shown in Table 4 and Figure 13. The presence of BSA strongly influenced $R_{Lys}$ evolution, increasing it by almost three times (Figure 13B).

### 4.1.3 Concentration polarization correction

A correction function was built in order to account for the accumulation of BSA on the membrane surface as a result of concentration polarization (Figure 14A). This function was constructed with data from the concentration steps of the training experiment (2 bar TMP and 200 mL/min CF). Since these parameters were in the mid-point of the training space (Figure 18A), the same correction function was used in all models to predict all test sets, regardless of

their TMP and CF parameters. Therefore, no sampling steps other than for the construction of the black box model for flux prediction were necessary, since the correction curve was built with just the $c_{B,BSA}$ measured at the beginning of each scouting round. Similarly to $R_{Lys}$ prediction, including the influence of TMP and CF on the amount of protein accumulated on the membrane would have required more sampling steps as well as a separate black box for its calculation, which would rather have complicated the usability of the models. This is one of the main advantages of the proposed hybrid models, that offer more precise and robust predictions not only of the flux, but also of the product and the impurity concentration evolution, without requiring any additional sampling step compared to the previously developed one-component-hybrid model, thereby making their implementation in industry easier.



Figure 14: Difference between measured and calculated concentration in the bulk of (A) BSA and (B) lysozyme, during the combined training experiment as a result of concentration polarization. Figure 14A from [68].

$$c_{B,measured\ BSA} = -0.0016 \cdot c_{B,calculated\ BSA}^{2} + 1.039 \cdot c_{B,calculated\ BSA} \tag{23}$$

Regarding lysozyme, the deviations between expected and measured final $c_{B,Lys}$ in the training set were very small (Figure 14B). As a result, no correction function for this protein was introduced in the models. The reasons for the smaller accumulation of lysozyme were, first of all, its lower concentrations compared to BSA, and additionally, that lysozyme was only partially retained by the membrane.

The measured amount of each model protein accumulated on the membrane at the end of each training set, calculated by the mass balance, is shown in Table 4.

In the case of the training set with only BSA, only the final missing amount of the first concentration range is shown. For the second training concentration range, the accumulated amount of BSA was 35.3%, since the final $c_{B,BSA}$ was 277.3 g/L - confirming again that the accumulated amount of protein is dependent on $c_{B,i}$.

Table 4: Summary of accumulated amount of BSA and lysozyme on the membrane at the end of each training set, in percentage, as well as the $R_{Lys,average}$ for the training sets containing lysozyme.

| Training set | BSA accumulated [%] | Lys accumulated [%] | $R_{Lys,average}$ |
|---|---|---|---|
| BSA alone | 12.7 | - | - |
| Lysozyme alone | - | 0 | 0.27 |
| Lysozyme + BSA | 10.3 | 1.68 | 0.77 |

The missing amount of each protein at the end of each test run as well as the amount that was recovered after recirculation with fresh buffer, in percentage, is shown in Table 5. As it can be seen, there were some deviations depending on the applied process conditions ($c_{B,BSA}$, $c_{B,Lys}$, TMP and CF, see Table 2). This will be further discussed in *4.4 Final protein concentration prediction* section. The $R_{Lys,average}$ for each test run was also calculated by using Eq. (22).

Table 5: Summary of the amount of each model protein missing in the mass balance at the end of each test run and the total amount that was recovered after recirculation with fresh buffer, in percentage. The $R_{Lys}$ average for each test set is calculated according to Eq. (22).

| Test set number | BSA missing [%] | Total BSA recovered after recirculation [%] | Lys missing [%] | Total Lys recovered after recirculation [%] | $R_{Lys,verage}$ |
|---|---|---|---|---|---|
| 1 | 26.1 | - | - | - | - |
| 2 | 22.0 | 80.3 | 4.3 | 100.0 | 0.78 |
| 5 | 11.8 | 100.0 | 2.7 | 100.0 | 0.77 |
| 6 | 8.8 | 92.0 | 6.0 | 100.0 | 0.81 |
| 7 | 9.2 | 100.0 | 6.9 | 100.0 | 0.74 |
| 8 | 25.3 | 96.3 | 8.1 | 100.0 | 0.80 |
| 9 | 26.8 | 93.7 | 7.7 | 100.0 | 0.75 |
| 10 | 21.5 | 97.0 | 9.5 | 100.0 | 0.79 |

As it can be seen in Table 5, after recirculating the system with fresh buffer (*3.1.2 Training and test data generation*), most of the protein that was missing in the mass balance was recovered. For test set 1 no recirculation was made.

## 4.2 Flux prediction

### 4.2.1 Comparison between mechanistic and hybrid models

#### 4.2.1.1 One-component models trained on one-component protein solution (ocHM$_{BSA}$ and SFM$_{BSA}$)

As a first step, the previously developed one-component-hybrid model for a membrane with 10 kDa MWCO[51] was reproduced with the current membrane cut-off of 30 kDa. A training run based on solely BSA was first performed, and the hybrid model built on that data (ocHM$_{BSA}$) was evaluated for the flux prediction of test set 1 (only BSA, Table 2 and Figure 15A). The classical one-component stagnant film model (SFM) was also built on the same data (SFM$_{BSA}$) for comparison of flux prediction. As it can be seen in Figure 15A, both models were able to predict the ultrafiltration process with only BSA - with the hybrid model having an NRMSE of 1.8% compared to the 4.9% of the SFM. However, they both failed when predicting the flux of solutions where lysozyme was additionally present - Figure 15B gives an example for test run 4. This demonstrated that the presence of even low amounts of an additional component can significantly influence the flux evolution in UF processes. This is also shown in Table A3 for the rest of the test sets, where the NRMSE of ocHM$_{BSA}$ went up to 12.1% for some of the test runs with both proteins. Hence, in order to build accurate models, the generation of training data under conditions where all the components were present in the system was necessary. In this direction, a training experiment consisting of a solution of BSA with lysozyme was carried out (Table 1).



Figure 15: Comparison of flux prediction of (A) test set 1, with only BSA (TMP 1.8 bar, CF 200 mL/min, initial $c_{B,BSA}$ 6.68 g/L), and (B) test set 4, with BSA and lysozyme (TMP 2.1 bar, CF 250 mL/min, initial $c_{B,BSA}$ 3.71 g/L and initial $c_{B,Lys}$ 0.38 g/L) between the one-component-hybrid model 1, trained on only BSA (ocHM$_{BSA}$, short-dashed light red line) and BSA with lysozyme (ocHM$_{comb}$, long-dashed dark red line), and the SFM model based on k and $c_G$ values

from only BSA (SFM$_{BSA}$, dot-dot-dashed light blue line) and BSA with lysozyme (SFM$_{comb}$, dot-dashed grey line) training sets.

## 4.2.1.2 One-component models trained on two-component protein solution (ocHM$_{comb}$ and SFM$_{comb}$)

A common approach in modeling of downstream processes for simplification is to neglect the process-related impurities. The previous one-component models were trained on the BSA with lysozyme training data but just including the BSA concentration as model input. The aim of these models (ocHM$_{comb}$ and SFM$_{comb}$), was to investigate the flux predictions without having to quantify the impurity concentration from the solutions. Additionally, ocHM$_{comb}$ allowed for a fairer comparison strictly of flux prediction between the hybrid and the stagnant film model, since the latter is only suited for one-component systems.

As shown in Figure 15B and Table A3, both ocHM$_{comb}$ and SFM$_{comb}$ could predict the flux of two-component test runs when built on the BSA with lysozyme training data, with the hybrid model providing again superior predictions in all the test sets compared to the SFM (Table A3). This confirmed that the hybrid model was a better candidate for flux prediction, both in one and two-component solutions. The main reason for this was that, while the SFM just considered the filtration process to take place under the pressure-independent region -i.e., that a gel-layer was already formed from the beginning on the membrane surface and the flux was thereby only influenced by mass transfer affecting parameters-, the hybrid models could explain both the pressure dependent and independent regions. This was possible due to the good interpolation properties of the ANN functions (*1.2.2.1 Artificial neural network (ANN)*) in their black box models. These functions can represent both linear and non-linear relationships, thereby making the models more flexible and precise. This phenomenon also explained why the SFM typically tends to overestimate the fluxes at the beginning of the process[14] (see Figure 20), when the solute concentrations are small, or at low TMPs and high CFs conditions – since the assumed linear relationship between flux and ln($c_G$) is not true under the pressure-dependent region, see Figure 10 and Eq. (21).

Nevertheless, both ocHM$_{comb}$ and SFM$_{comb}$ failed in predicting those test runs where the initial $c_{B,Lys}$ was lower than in the training experiment – test sets 1, 7 and 8 (Table A3 and Figure 15A). Being based on the training run containing BSA with lysozyme, these models incorporated the fouling effects related to the training lysozyme concentration. Therefore, regardless the test sets had lower lysozyme concentrations, the models kept making predictions assuming the same degree of fouling as in the training experiment and therefore underestimating the flux. In the SFM, for example, the $k$ values from the two-component training set were generally smaller than the ones calculated with BSA only (Table A2), since membrane fouling by lysozyme hindered the solute mass transfer in the boundary layer - reaching the pressure-independent region earlier, Figure 10B.

43

Hence, even though ocHM$_{comb}$ and SFM$_{comb}$ gave acceptable results where test and training concentrations were comparable, variations in the amount of the not included impurity component would lead to wrong flux predictions.

## 4.2.2 Comparison between two-component-hybrid models (tcHMs)

To further improve the model predictions and incorporate $c_{B,Lys}$ and $R_{Lys}$, three types of two-component-hybrid models (tcHMs) were built. These models differed in the way $R_{Lys}$ was calculated. The performance of the models was compared for flux, $R_{Lys}$ and final $c_{B,i}$ prediction. In the first hybrid model, tcHM1$_{Raverage}$, $R_{Lys}$ was assumed to be equal to the $R_{Lys,average}$ from the BSA with lysozyme training set, namely 0.77 (*4.1.2.1 RLys average from training set*). Furthermore, in order to evaluate the impact of calculating $R_{Lys}$ on the flux, a hybrid model with no lysozyme rejection factor ($R_{Lys}=1$) was built as control (tcHM1$_{R1}$). Finally, $R_{Lys}$ was predicted dynamically by a black box model, either by the existing one for flux (tcHM2) or by introducing a completely new black box solely for $R_{Lys}$ (tcHM3).

The optimal ANN structure for flux prediction in all the hybrid models was determined by varying the number of hidden nodes from one to six in one hidden layer and testing the obtained models on a validation data set, which was not used for training. The flux NRMSE out of 20 repetitions was recorded and averaged. The ANN with four hidden nodes yielded the lowest NRMSE in all models. Further increase of the hidden nodes led to higher errors due to training set overfitting as well as to less model consistencies (more prediction variabilities, higher standard deviations). As an example, the results from the ANN structure optimization process for tcHM1$_{R1}$, tcHM2 and tcHM3 (with an ANN black box model with one hidden node for $R_{Lys}$) are shown in Figure 16.



Figure 16: ANN node size optimization. The NRMSE is plotted over the number of neurons in the hidden layer of the ANN flux predicting black box model of tcHM1$_{R1}$, tcHM2 and tcHM3$_{ANN}$.

44

The structure of the chosen optimal ANN black box model for flux prediction is shown in Figure 17. This model consisted of four input parameters (TMP, CF, $c_{B,BSA}$, $c_{B,Lys}$), one output parameter (J) and one hidden layer with four hidden nodes. The result of each hidden node in the hidden layer (see Eq. (24) as an example for the first neuron $x_{12}$) is the sum of each multiplication of an input parameter and the corresponding weight ($w^1_{11}$, $w^1_{21}$, $w^1_{31}$, $w^1_{41}$), multiplied with the bias ($b_1$) of the entire hidden layer. The sigmoid activation function then transforms the output of each hidden node by using Eq. (25). The inputs were normalized between 1 and 2, rendering them comparable.



Figure 17: Structure of the ANN black box model used for flux prediction, indicating the input, output and hidden layers and the number of hidden nodes (neurons) and activation functions. Figure adapted from [68].

$$x_{1,2} = b_1\left(w^1_{11}TMP_{scaled} + w^1_{21}CF_{scaled} + w^1_{31}c_{B,BSA,scaled} + w^1_{41}c_{B,Lys,scaled}\right) \tag{24}$$

$$h_1(x) = \frac{1}{1+e^{-x_{1,2}}} \tag{25}$$

In order to evaluate both the inter- and extrapolation capabilities of these models, different test runs were conducted in a variety of conditions, some of them being only partially covered by the training sets. Ten test experiments were performed in total at different TMP, CF, initial $c_{B,BSA}$ and $c_{B,Lys}$, which are all depicted in Figure 18 and listed in Table 2. Test runs 2-6 were performed within the training space, meaning that TMP and CF were within the trained parameters (Figure 18A) and that the initial $c_{B,BSA}$ and $c_{B,Lys}$ were higher than the initial training

concentrations (Figure 18B – blue area). From them, test runs 2 and 6 were performed in the center of the TMP and CF training space, while test run 3 was performed at the outer limit in order to investigate how the models behaved at the border. Additionally, test runs 4 and 5 were performed under TMP and CF conditions that were within the training space but not under the exact conditions used in the training set, in order to investigate the interpolation capabilities of the models. On the contrary, the test runs 1 and 7-10 were performed under conditions that were partially outside the training space such as initial $c_{B,Lys}$ (1, 7, 8), initial $c_{B,BSA}$ (9) and CF (10), to test the models' extrapolation capabilities. The observed and predicted fluxes are depicted in Figure 20 for the test sets inside the training space, and in Figure 21 for the test sets partially outside. The NRMSE of the flux prediction of each test run by each developed model is shown in Table A3.



Figure 18: Schematic depiction of (A) training space for TMP and CF parameters of the training (white dots) and test (grey dots) runs and (B) initial to final $c_{B,BSA}$ and $c_{B,Lys}$ of the three training (white dots with black solid lines) and the test (grey dots with grey solid lines) runs. The covered concentration area of the BSA with lysozyme training set (training space) is shown in blue. Figure adapted from[68].

In Figure 19, the average NRMSE of the different models for flux prediction of all test runs is shown. As it can be observed, the two-component-hybrid models (tcHM) performed very well, both inside (Figure 19B, 3.2% NRMSE) and outside (Figure 19C, 4.5% NRMSE) the training space - tcHM2 was an exception, which is explained in the following. On the contrary, the one-component models particularly failed when predicting the test runs with initial $c_{B,Lys}$ lower than in the training set. The average NRMSE for the flux prediction outside the training space was 6.4% and 7.6% for ocHM$_{comb}$ and SFM$_{comb}$ models, respectively. This demonstrated that the established two-component-hybrid models were able to differentiate between the influence that each of the components had on the flux decline. This is in line with the fact that the best

simulated test set by all tcHMs was test set 1, without lysozyme, with just a flux NRMSE of 1.7%. Furthermore, no significant difference in flux prediction was observed between tcHMs, indicating that the way of calculating $R_{Lys}$ did not have a notable impact. The reason for this was that the concentration of lysozyme was 4 to 46 times lower than BSA and therefore, a difference in its rejection factor calculation only led to small differences in $c_{B,Lys}$ predictions, which did not change the flux predictions notably. If the processes would take longer (e.g. larger bulk volume or small membrane area) or if the initial $c_{B,Lys}$ had been higher, the way of calculating $R_{Lys}$ would have had a higher impact on flux, since it would lead to bigger differences in $c_{B,Lys}$ prediction and thereby in flux. Finally, the only two-component-hybrid model where the calculation of $R_{Lys}$ influenced the flux predictions was tcHM2. In this model, the fact of combining $R_{Lys}$ and flux under the same black box showed to negatively affect the calculation of both parameters, especially for the test sets that were partially outside the training space (Figure 19C). This was most probably due to a different optimal node size structure as well as optimal weight values for both predicted parameters. This will be further discussed in section *4.3 Rejection factor prediction*.



Figure 19: Average NRMSE for flux prediction of test sets (A) all, (B) inside and (C) outside the training space, out of 20 repetitions, by the different constructed models.

Due to the very similar predicted fluxes of the tcHMs, only the predictions from tcHM1$_{Raverage}$, tcHM2 and tcHM3$_{ANN}$ are shown in Figure 20 and Figure 21, together with the predictions from ocHM$_{comb}$ and SFM$_{comb}$ and the experimentally measured values. Despite the overall good flux predictions by these models, small initial flux underestimations were observed in all test sets (Figure 20 and Figure 21). This phenomenon was, first of all, due to a difference in the experimental setup between training and test experiments. While the recorded fluxes in the training set were taken at equilibrium for each combination of $c_{B,BSA}$, $c_{B,Lys}$, CF and TMP, in the test sets the processes started with a completely clean membrane and therefore the

equilibrium had to be reached during the test process (*3.1.2 Training and test data generation*). Nevertheless, the time to reach such concentration polarization steady-state in cross-flow ultrafiltration of proteins is described to be very short[14]. This was the reason why this phenomenon was quickly corrected and actually only noticeable for the test sets that exhibited high TMPs, where the concentration polarization effects were more pronounced, such as test sets 3 and 5 (Figure 20B and D, see the initial recorded fast decline in flux (black solid line)). Therefore, the further observed deviations between predicted and measured fluxes were due to other phenomena that are explained in the following.

The initial offset between measured and predicted fluxes became more pronounced for those test sets with initial $c_{B,Lys}$ higher than 0.3 g/L (test runs 3, 4, 6, 9 and 10), suggesting that there was a stronger membrane fouling behavior at these concentrations. In the test runs it took some time until the increasing membrane resistance due to fouling of lysozyme reached the equilibrium. However, the hybrid models were falsely simulating equilibrium membrane fouling conditions already from the beginning. After this time, however, the fluxes were predicted correctly. Related to this phenomenon, the test set with highest obtained initial offset -and also highest flux prediction error- was test run 9 (Figure 21C) with an NRMSE of 8%  for flux prediction. This test set had the biggest initial $c_{B,Lys}$ and therefore the strongest fouling. Additionally, its initial $c_{B,BSA}$ was lower than in the training set (Table 2), thereby reducing its influence on the flux and further contributing to the error. The second worst predicted test run, although with a much lower NRMSE, of 4.9%, was test set 7 (Figure 21A). In this case, the initial $c_{B,Lys}$ was lower than in the training set and therefore required the models to extrapolate for these initial input values. Moreover, it is thought that no fouling at all occurred during this test set, further increasing the error and explaining why the observed initial offset in this test run was not corrected until late in the process (Figure 21A).This is in accordance with the results shown in section *4.3 Rejection factor prediction*, where the highest obtained error for $R_{Lys}$ prediction was in this test set (Table A4). Finally, the third largest obtained deviation in flux prediction was in test set 10 (Figure 21D), which was performed at CF 350 mL/min, which was outside of the training space. However, its small error, of only 4.5% NRMSE, demonstrated that the models had good extrapolation capabilities for this parameter and that therefore were not limited by the training space. The rest of test runs were predicted very well, with an average NRMSE of 3% in spite of being tested in a variety of different process conditions. For example, for TMP and CF interpolations, test runs 4 and 5 were predicted with an average of 3.3% NRMSE, while test set 3, at the outer edge of the training space, was predicted with just a 3.2% NRMSE.

The predictions were also very good for the test sets with $c_{B,BSA}$s more than two times higher than in the training space -test set 8, 3.3% NRMSE-, or with only BSA and no lysozyme - test set 1, 1.7% NRMSE. The reason why ocHM_comb yielded such good flux predictions within the

training space (Figure 19B, with 3.1% average NRMSE), was because it did not include lysozyme, and therefore did not underestimated the initial flux by falsely assuming lysozyme related fouling equilibrium conditions. Similarly, this was the reason why the best test set predicted by tcHMs was test set 1, with only BSA.

Summarizing, good flux predictions of two-component test sets could only be achieved if both components were present during the training experiment. Furthermore, the incorporation of each of the components in the models was also necessary for accurate results, especially for input parameters combinations that were outside of the training space. This was particularly crucial when extrapolating for the mimicked impurity concentration, which was not included in the one-component models. Therefore, both one-component models (the mechanistic SFM and the ocHM), were outperformed by the two-component models. Regarding one-component models, however, the ocHM outperformed SFM due to its good interpolation properties and its capability to shift between pressure-dependent and independent region. Moreover, the hybrid models could include both components by simply adding an input parameter that represented the new component. The relationship between variables was established by machine learning algorithms, without requiring in-depth knowledge of the underlaying mechanism -fouling, CP layer formation rate, protein interactions…- nor measuring any physical property, which are the main limitations for implementing mechanistic multi-component models. The established two-component-hybrid models were very precise and consistent, over a broad range of TMP, CF, $c_{B,BSA}$ and $c_{B,Lys}$ input parameters, exhibiting excellent interpolation capabilities as well as very good extrapolations for CF and higher $c_{B,BSA}$. Extrapolations for smaller $c_{B,Lys}$ were also good, however, these models faced some problems when predicting test runs with $c_{B,Lys}$ higher than the training set due to not incorporating the time-dependent fouling of lysozyme. Nevertheless, in most of the downstream processing units where ultrafiltration is commonly used, the impurities content is rarely that high - e.g., after capture or between polishing steps.

Figure 20: Comparison between observed and predicted fluxes over time of the test sets within the training space: (A) test set 2, (B) test set 3, (C) test set 4, (D) test set 5 and (E) test set 6.

Figure 21: Comparison between observed and predicted fluxes over time of the test sets partially outside the training space: (A) test set 7, (B) test set 8, (C) test set 9, (D) test set 10 and (E) test set 1.

## 4.3 Rejection factor prediction

The ANN black box model for $R_{Lys}$ prediction was optimized by changing the number of hidden nodes from one to four and recording the average $R_{Lys}$ NRMSE out of 20 repetitions for the validation data set. In tcHM2, where the same black box was used for flux and $R_{Lys}$ prediction, independent hidden node optimization was not possible. Therefore, the optimal number for flux prediction was chosen. Contrary, in tcHM3, the number of nodes for flux was kept at four during the optimization, since it was the optimum. The results for the ANN model optimization are shown in Figure 22.



Figure 22: ANN structure optimization of tcHM2 (black) and tcHM3 (grey) by changing the number of hidden neurons from one to four in one hidden layer. The average NRMSE for $R_{Lys}$ (circles) and flux (triangles) prediction out of 20 repetitions is shown. In the case of tcHM3, only the black box for $R_{Lys}$ was optimized, since the black box for flux was kept at four nodes.

As it can be seen in Figure 22, one hidden node exhibited the lowest error for $R_{Lys}$ prediction in both models. When increasing the number of nodes, the error and standard deviation increased due to model overfitting. Since in tcHM2 the flux and $R_{Lys}$ were predicted by the same black box, increasing the number of hidden nodes led to a decreasing error for flux while increasing the $R_{Lys}$ error. Due to this discrepancy, the flux error of tcHM2 was higher than in tcHM3. Contrary, in tcHM3 changing the $R_{Lys}$ node size did not affect the flux prediction, due to being both outputs in two independent black boxes.

Since the correlation between $R_{Lys}$ and input parameters (TMP, CF, $c_{B,BSA}$ and $c_{B,Lys}$) was shown to be rather simple, with only one hidden node as optimum, a multiple linear regression (MLR) model was set up as alternative to the black box of tcHM3. MLR required less computation times and was easier to interpret than ANN due to its fewer coefficients.

The MLR model was both forwards and backwards optimized both for the p- and the AIC statistical values. The resulting regression equation out of the four different optimization

strategies was the same, which is shown in Eq. (26) together with its statistical coefficients (Figure 23).

$$R_{Lys} = -0.4 \cdot TMP + 0.3 \cdot CF + 1.2 \cdot c_{B,BSA} + 3.2 \cdot c_{B,Lys} + 0.2 \cdot TMP \cdot c_{B,Lys} - 0.1 \cdot CF \cdot c_{B,Lys}$$
$$-1.3 \cdot c_{B,BSA} \cdot c_{B,Lys} \tag{26}$$

```
Estimated Coefficients:
                    Estimate        SE         tStat       pValue

    (Intercept)     -1.7726      0.28524      -6.2144     2.0367e-08
    x1              -0.40896     0.10275       -3.98      0.0001481
    x2               0.25564     0.088535      2.8875     0.0049643
    x3               1.1908      0.35446       3.3595     0.001187
    x4               3.2476      0.34854       9.3177     1.6804e-14
    x1:x4            0.24044     0.073576      3.2679     0.0015839
    x2:x4           -0.13092     0.063383     -2.0655     0.042034
    x3:x4           -1.3088      0.094472    -13.854      3.476e-23

Number of observations: 90, Error degrees of freedom: 82
Root Mean Squared Error: 0.0852
R-squared: 0.921,  Adjusted R-Squared 0.914
F-statistic vs. constant model: 137, p-value = 1.87e-42
```

Figure 23: Statistical coefficients of the multiple linear regression (MLR) model with the p-value for each of the predictors.

As it can be seen in Eq. (26) and Figure 23, all four input parameters were included in the final regression model both as linear and as interactions terms with $c_{B,Lys}$. Quadratic terms were not included in the model because they provided much worse results during model validation.

Finally, due to the clear predominance of $c_{B,Lys}$ over the other three input parameters on determining $R_{Lys}$ evolution (Figure 12), and its clearly logarithmic trend, shown in Figure 24, the MLR model was modified by introducing a logarithmic term for $c_{B,Lys}$, resulting in tcHM3$_{MnLR}$.



Figure 24: Lysozyme rejection factor values from the combined training set over $c_{B,Lys}$ for different TMPs at 200 mL/min cross-flow velocity.

The MnLR model was forward stepwise optimized by the p-value criterion. The final regression equation and its statistical coefficients are shown in Eq. (27) and Figure 25, respectively.

$$R_{Lys} = 4.5 - 0.4 \cdot TMP + 0.3 \cdot CF - 3.2 \cdot c_{B,BSA} + 5.4 * log(c_{B,Lys}) + 0.2 \cdot TMP \cdot c_{B,BSA} + 1.5 \cdot CF \cdot c_{B,BSA} - 1.67 \cdot CF \cdot c_{B,Lys} \tag{27}$$

|     | Estimate | SE       | tStat   | pValue      |
| --- | -------- | -------- | ------- | ----------- |
| b1  | 4.5081   | 0.29366  | 15.352  | 7.8007e-26  |
| b10 | -1.6652  | 0.26165  | -6.3641 | 1.0633e-08  |
| b2  | -0.39782 | 0.093953 | -4.2342 | 5.9526e-05  |
| b3  | 0.28966  | 0.0812   | 3.5673  | 0.00060568  |
| b4  | -3.1906  | 0.27273  | -11.699 | 3.6472e-19  |
| b5  | 5.4279   | 0.34738  | 15.625  | 2.6359e-26  |
| b7  | 0.23285  | 0.067408 | 3.4543  | 0.0008759   |
| b9  | 1.5135   | 0.26417  | 5.7292  | 1.6166e-07  |

Figure 25: Statistical coefficients of the multiple non-linear regression (MnLR) model.

The average NRMSE for $R_{Lys}$ prediction out of 20 repetitions of all the test sets and all the developed hybrid models is shown in Figure 26A. The errors were differentiated between the prediction of the test sets inside (Figure 26B) and partially outside (Figure 26C) the training space.



Figure 26: Average NRMSE for $R_{Lys}$ prediction of the test sets (A) all, (B) inside and (C) outside the training space, out of 20 repetitions, by the different constructed hybrid models.

As clearly seen in Figure 26A, tcHM3$_{ANN}$ was the model that best predicted $R_{Lys}$ over time, with an overall average NRMSE of 14.3% and keeping good performances both inside -12.1% NRMSE, Figure 26B- and outside -16.6% NRMSE, Figure 26C- the training space. tcHM2

performed second best. However, it was not a good model candidate due to its aforementioned worse flux predictions, especially outside the training space.

MLR and MnLR models were able to make acceptable $R_{Lys}$ predictions of the test sets that were inside the training space (Figure 26B), but they failed for those outside (Figure 26C). This was due to the estimated coefficients in their regression equations. In particular, in MLR model a negative interaction term between $c_{B,BSA}$ and $c_{B,Lys}$ (Eq. (26)) made that the predicted $R_{Lys}$ values started to decrease close to the end of the process, when the concentration of both proteins started to be high enough to make this term of the equation governate over the others. This was the reason why this model especially failed for the test sets with the highest initial $c_{B,BSA}$ and $c_{B,Lys}$: test sets 8 and 9, respectively (Table A4 and Figure 28B and C). In a similar way, MnLR model gave very bad predictions in test set 9 -Figure 28C, with a 89.9% NRMSE- and to a smaller extent in test 10 -Figure 28D, 34.5% NRMSE-, due to the incorporation of a negative interaction term between $c_{B,Lys}$ and CF (Eq. (27)), which counteracted the positive logarithmic term for $c_{B,Lys}$ at high lysozyme concentrations. This phenomenon highlighted one of the main limitations of data-driven models: their extrapolation capabilities. Since the structure and coefficients of these models were inferred from the data, they were adjusted in order to fit them in the best possible way. Consequently, they faced some problems when predicting data at completely different conditions. This is why larger and more representative amounts of data are necessary to build these models compared to mechanistic models - where the parameters have a physical meaning and thereby have better extrapolation capacities. In this case, it was assumed that the relationship between $R_{Lys}$ and input variables was not linear, and consequently, MLR and MnLR models over-fitted to explain the training data as a linear combination of parameters. This yielded acceptable results for the test sets inside the training space but strongly failed outside.

Finally, the tcHM1$_{Raverage}$ predictions of the test sets inside and outside the training space were very similar (Figure 26B and C), despite exhibiting higher errors due to its static $R_{Lys}$ value, with an average NRMSE of 37.7% for all test sets. This indicated that the $R_{Lys}$ average obtained from the training set fitted all independently generated test data quite well. The $R_{Lys}$ values were overestimated at the beginning and underestimated at the end of the process, regardless of the applied process parameters. This was true as long as the total concentration factor - ratio between initial and final reservoir volume- was kept similar. This feature made tcHM1$_{Raverage}$ a good model candidate for those scenarios where model simplicity and interpretability are preferred over accuracy and complexity. This will be further discussed in section, *4.4.2 Lysozyme* prediction.

The measured and predicted $R_{Lys}$ evolution over $V_p$ by the different models is shown in Figure 27 and Figure 28 for the test sets inside and outside the training space, respectively. In spite of the overall good $R_{Lys}$ predictions of tcHM3$_{ANN}$, test sets 7 and 8 showed a great discrepancy

between the predictions inside and outside the training space (Figure 28A and B). These test sets had the highest $c_{B,BSA}$s and lowest TMPs and $c_{B,Lys}$s values of all tested runs (Table 2). Due to not incorporating the time-dependent protein accumulation on the membrane (see *4.2 Flux prediction*), the models assumed that the BSA CP layer was at equilibrium from the beginning of the process. This led to a higher $R_{Lys}$ compared to the test runs, since the BSA CP layer was still building up. Moreover, the lower TMP and $c_{B,Lys}$ of these test sets additionally prolonged the time to reach the CP layer -and also the lysozyme-related membrane fouling-equilibrium, further increasing the $R_{Lys}$ overestimations and thereby the errors. Furthermore, it is important to highlight that the $R_{Lys}$ values that were fed to the black box for training the models already accumulated some small errors from their calculation (Eq. (18)). First of all, these values were calculated assuming that the $c_{B,Lys}$ was constant during each scouting round, and therefore independent of TMP and CF. This also explained why the effect of these two parameters in $R_{Lys}$ was basically only seen in the first scouting rounds (Figure 12), where $c_{B,Lys}$ was lower and therefore variations in $c_{p,Lys}$ had a higher impact on $R_{Lys}$. Moreover, in some parameters -in particular at small TMPs and high CFs- of the first scouting rounds of the training experiment, the recorded UV values were below the recorded absorbances for constructing the calibration curve (Figure 11), and therefore outside the linear regression range. This resulted in underestimated $c_{p,Lys}$ values under these conditions and in turn in overestimated $R_{Lys}$ (Eq. (18)). This would further explain the initial $R_{Lys}$ overestimations -and therefore the higher errors- seen for the predictions of the test sets with lower initial $c_{B,Lys}$ than the training set, since in addition to extrapolating for $c_{B,Lys}$, the models used data with some errors in this vicinity. It is for this reason that the quality of the provided data and in this particular case the determination of the regression curve, is of utmost importance for the predictions of all hybrid models that iteratively update $R_{Lys}$.

Hence, taking into account all the aforementioned results and the simple performed experimental setup, it can be concluded that the $R_{Lys}$ predictions made by tcHM3$_{ANN}$ were very good and robust over a wide range of process conditions, with excellent interpolation capabilities as well as extrapolations for CF and smaller $c_{B,BSA}$. However, this model faced some problems when predicting test sets with higher $c_{B,BSA}$ and lower initial $c_{B,Lys}$ than the training set, due to not including the time-dependent BSA CP and lysozyme fouling formation, which led to initial $R_{Lys}$ overestimations. These overestimations were additionally accentuated by the fact that the provided $R_{Lys}$ values for training the models had some errors accumulated from their calculation.

Figure 27: Comparison between measured and predicted $R_{Lys}$ values over permeate volume, $V_p$, for (A) test set 2, (B) test set 3*, (C) test set 5 and (D) test set 6, by the different developed hybrid models.

*The obtained $R_{Lys}$ curve in test set 3 (Figure 27B), with such a sharp shape, is thought to be due to a problem with the HPLC column when analyzing the samples. This column was dispensed after this run due to it was difficult to build a good calibration curve with the standards. It was not possible to analyze these samples again with a new column due to Covid-19 lockdown issues. However, in order to have more test sets for comparing the models performance, these results are shown in the present work, despite increasing the overall average NRMSE. Without this test set, the tcHM3$_{ANN}$ average $R_{Lys}$ NRMSE for all the test sets would be 13.5%, and 9.3% inside the training space.

Figure 28: Comparison between measured and predicted $R_{Lys}$ over permeate volume, $V_p$, for (A) test set 7, (B) test set 8, (C) test set 9 and (D) test set 10, by the different developed hybrid models.

## 4.4 Final protein concentration prediction

### 4.4.1 BSA

The final goal of modeling the rejection factor of any protein during a filtration process is to be able to accurately describe its concentration evolution over time, both in the bulk and permeate solutions. In the case of the mimicked product in this work, BSA, its rejection factor with the 30 kDa MWCO membrane was 1, and consequently, all models predicted the same $c_{B,BSA}$ at the end of the process. The predicted $c_{B,BSA}$ was calculated solely by using the ratio between filtered and reservoir volumes (Eq. (16-18) and Table A1) and the correction function accounting for the accumulated protein on the membrane as a result of concentration polarization (*4.1.3 Concentration polarization correction*, Eq. (23)). By using these equations, the average difference between predicted and measured final $c_{B,BSA}$ for all the test sets was 9.8%, which would have been 24.3% if no correction function would have been used at all, as it is shown in Figure 29.



Figure 29: Difference between observed and predicted final $c_{B,BSA}$ by all the hybrid models using the correction function (dark red) or not (light red) for protein accumulation on the membrane, in all the test sets, and in the test sets inside and outside the training space.

As it can be seen in Table 6, in most of the test sets the predicted final $c_{B,BSA}$ was higher than the measured one, giving place to the errors displayed in Figure 29. The reason for these final $c_{B,BSA}$ overestimations could be derived from the fact that the predictions were made using the $c_{B,BSA}$ correction function from the concentration steps of the training experiment (Eq. (23)). These data were generated using the training set process conditions (TMP of 2 bar and CF of 200 mL/min), without considering their potential effects nor the effect of $c_{B,Lys}$ on $c_{B,BSA}$. Nevertheless, as shown in Table 6, these parameters indeed had an impact on the final amount of protein accumulated on the membrane – one of the clearest examples was test set

7, where despite its high final $c_{B,BSA}$, only a 9.2% of the total initial amount was missing at the end of the process, while in test set 1, with similar final $c_{B,BSA}$ but no lysozyme, this amount increased to 26.1%. Furthermore, the fact that during the training experiment the reservoir solution was recirculated for about five minutes at the end of each scouting round with zero TMP, it is thought to have contributed to the overestimations in final $c_{B,BSA}$. In the test runs, the reservoir solution was completely filtered at once, without any recirculation step in between, thereby leading to higher amounts of accumulated protein on the membrane. Finally, in some test sets the $c_{B,BSA}$ range was partially outside the training space -from 77.9 to 162.5 g/L, Figure 18B-, thus requiring the correction function to extrapolate and leading to final $c_{B,BSA}$ underestimations. Actually, except for test run 1, all the test sets with underestimated final $c_{B,BSA}$ were those with higher $c_{B,BSA}$ range than the training set: test sets 5, 7 and 8, (Figure 18B and Table 6).

Table 6: Measured and predicted final BSA concentration in the bulk, $c_{B,BSA}$, for each test set by the models, using or not the correction function for BSA accumulation on the membrane.

| Test set # | Measured $c_{B,BSA}$ [g/L] | Predicted $c_{B,BSA}$ [g/L] | | Difference [%] | |
|---|---|---|---|---|---|
| | | c_corr | no corr | c_corr | no corr |
| 1 | 132.7 | 135.0 | 179.6 | 1.7 | 35.4 |
| 2 | 78.1 | 87.9 | 100.1 | 12.6 | 28.2 |
| 5 | 97.6 | 95.4 | 110.7 | 2.2 | 13.4 |
| 6 | 62.5 | 63.6 | 68.5 | 1.8 | 9.6 |
| 7 | 132.8 | 117.7 | 146.3 | 11.3 | 10.2 |
| 8 | 162.5 | 150.3 | 217.5 | 7.5 | 33.9 |
| 9 | 45.4 | 58.3 | 62.0 | 28.3 | 36.5 |
| 10 | 73.6 | 83.3 | 93.8 | 13.2 | 27.4 |
| | | | **Average** | **9.8** | **24.3** |

Summarizing, all models predicted the final $c_{B,BSA}$ rather well, especially in the test sets inside the training space, with only an average error of 5.5% compared to the experimental measured values. These differences were most probably due to using a correction function from the training experiment, which was generated with a different experimental setup compared to the test sets and without considering the influence of TMP, CF and $c_{B,Lys}$ parameters as inputs of the equation. Consequently, higher errors were observed for the predictions of the test sets with some parameters outside the training space. Nevertheless, the errors in $c_{B,BSA}$ prediction were too small to notably influence the flux predictions (*4.2 Flux prediction*). The errors increased with $c_{B,BSA}$, that is, exponentially with time, and therefore they started to gain importance only close to the end of the process. However, without using the correction function, the errors for final $c_{B,BSA}$ prediction increased to almost 3.5 times, with an average

error of 24.3% for all the test sets, and 17.1% and 28.7% for the test sets inside and outside the training space, respectively.

It is important to highlight, that the correction function used in this work was developed from just the $c_{B,BSA}$ values used for flux prediction, without requiring any additional sampling step. A more precise correction function could be established if incorporating the other parameters as inputs to the regression model. This, however, would require much more experimental efforts, since one sample should then be taken at each TMP-CF combination -i.e., 15 measurements compared to currently only 1 per scouting round (Figure 4A)-, as well as more complex models, due to a separated black box just for the $c_{B,BSA}$ correction function calculation would be necessary.

### 4.4.2 Lysozyme

Regarding the prediction of the final concentration in the bulk of the mimicked impurity, lysozyme, the obtained results were different than the expected if considering the $R_{Lys}$ NRMSE (Figure 26). It is for this reason that the discussion of these results was divided into the test runs inside and outside the training space, since the hybrid models performed differently in these two scenarios.



Figure 30: Average difference between observed and predicted final $c_{B,Lys}$ values by the different hybrid models in (A) all test sets and (B) test sets inside and (C) outside the training space.

When reviewing Figure 26, one would expect that the model with smallest $R_{Lys}$ errors also yielded the smallest errors in final $c_{B,Lys}$. However, this was only true for the prediction of the test sets that were inside the training space (Figure 30B), where tcHM3$_{ANN}$ predicted the final $c_{B,Lys}$ with an average error of only 5.6% difference compared to the experimental measured values.

Contrary, tcHM1$_{Raverage}$ performed as the worst -excepting tcHM1$_{R1}$, which was a control- , with an average error of 24.8% due to not iteratively adapting the R$_{Lys}$ over time. Finally, if the R$_{Lys}$ was assumed to be 1 (tcHM1$_{R1}$), the average final c$_{B,Lys}$ error for the test sets inside the training space increased to 54.4% (Table 7). This confirmed the importance of correctly calculating the rejection factor of a protein whenever modeling its concentration evolution.

In contrast, for the test sets that were performed partly outside the training space (Figure 30C), tcHM1$_{Raverage}$ performed best, with an average error of 11.4%, followed by tcHM3$_{MnLR}$ with 19.2%. On the contrary, tcHM3$_{ANN}$ yielded worse predictions (25.7%), despite having shown to better predict R$_{Lys}$ than the previous two models for these same test sets (Figure 26C). This switch in the models' performance between final c$_{B,Lys}$ and R$_{Lys}$ predictions was due to several reasons. On the one hand, tcHM3$_{ANN}$ overpredicted R$_{Lys}$ throughout most of the test runs - especially outside the training space (Figure 28)-, which led to higher final c$_{B,Lys}$ predictions and therefore higher errors. Conversely, in tcHM1$_{Raverage}$ and tcHM3$_{MnLR}$, the predicted final c$_{B,Lys}$ were balanced between the R$_{Lys}$ overpredictions at the beginning and the underpredictions at the end of the process, thus yielding smaller final c$_{B,Lys}$ predictions and thereby smaller errors, regardless of their higher error for R$_{Lys}$. In tcHM1$_{Raverage}$, the counterbalancing effect was due its constant R$_{Lys}$ value, while in tcHM3$_{MnLR}$ it was due to the shown parabolic shape (Figure 27 and Figure 28). This difference between models was additionally accentuated by the fact that the final amount of protein accumulated on the membrane was underestimated by all models -which means c$_{B,Lys}$ overestimations-, due to the difference in experimental setup between training and test sets when constructing the correction function. It is important to highlight, that in the case of lysozyme no correction function accounting for its accumulation on the membrane was introduced in the models, due to the observed small missing amount at the end of the training set (1.7%, Figure 14B). Nevertheless, in the test runs, the accumulated amounts ranged from 2.7 to 9.5% of the total initial lysozyme content (Table 5). Hence, as a result of not introducing any correction function, all models overpredicted the final c$_{B,Lys}$, thereby further increasing the errors in tcHM3$_{ANN}$ and decreasing them in tcHM1$_{Raverage}$ and tcHM3$_{MnLR}$ - which were already underpredicting the final c$_{B,Lys}$ in most test sets. Therefore, with a correction function the error differences between tcHM3$_{ANN}$, tcHM1$_{Raverage}$ and tcHM3$_{MnLR}$ for the test sets outside the training space would be expected to be smaller. Nevertheless, these results showed that tcHM1$_{Raverage}$ can work as a better candidate than tcHM3$_{ANN}$ for predictions of final c$_{B,Lys}$ in test sets with different conditions to the training experiment, rendering it a valuable tool for stable process extrapolations. This model is easier to construct and can be a good option for situations where simplicity is preferred over accuracy. For example, it could be used in early downstream filtration units, where the impurity content is high and variable and therefore the main interest is to ensure good flux predictions - rather than accurate predictions of the impurity.

It could also be used in multi-component solutions, where several solutes go to the permeate and thereby make difficult to construct black box models for dynamic rejection factor calculation for each of the components. However, it is important to remark that this model is only able to yield accurate predictions of the impurity concentration at the end of the process, since it does not update its $R_{Lys}$ over time, and as long as the concentration factor is similar to the training set.

Finally, tcHM3$_{MnLR}$, even though it showed superior predictions than its homolog tcHM3$_{MLR}$ and was in fact the model with smaller errors for final $c_{B,Lys}$ predictions (Figure 30A), it exhibited very strong deviations when predicting $R_{Lys}$ close to the end of the process. This was due to model over-fitting, since the model fitted the training data as a linear combination of parameters (see *4.3 Rejection factor prediction*), which gave rise to even time-decreasing $R_{Lys}$ curves (Figure 27 and Figure 28). Due to the relationship between $R_{Lys}$ and input variables was assumed to not be linear, the black box of tcHM3$_{MnLR}$ fitted the training data as best as possible but sacrificing predictability at the borders and outside of the trained region, leading to extreme over- or underestimations. The $R_{Lys}$ black box, however, required exact predictions especially at the training space border, since an error in $R_{Lys}$ prediction at this point has a stronger impact on $c_{B,Lys}$, due to the higher concentrations. Hence, taking all of this into account, tcHM3$_{MnLR}$ was not suited for $c_{B,Lys}$ predictions.

Finally, the final $c_{B,Lys}$ predictions from tcHM3$_{ANN}$ and tcHM1$_{R1}$ for each test set are compared to the experimental measured values in Table 7. As it can be seen, with tcHM3$_{ANN}$ the prediction errors were 4.5 times lower compared to tcHM1$_{R1}$, thereby highlighting the importance of modeling the rejection factor of a protein for predicting its concentration.

Table 7: Measured and predicted final lysozyme concentration in the bulk, $c_{B,Lys}$, and their difference in percentage by tcHM3$_{ANN}$ and tcHM1$_{R1}$ models.

| Test set # | Measured $c_{B,Lys}$ [g/L] | Predicted $c_{B,Lys}$ [g/L] | | Difference [%] | |
|---|---|---|---|---|---|
| | | tcHM3$_{ANN}$ | tcHM1$_{R1}$ | tcHM3$_{ANN}$ | tcHM1$_{R1}$ |
| 2 | 4.4 | 4.3 | 6.8 | 1.8 | 56.4 |
| 3 | 4.4 | 4.1 | 6.6 | 5.8 | 51.6 |
| 5 | 3.5 | 3.8 | 5.9 | 7.3 | 66.8 |
| 6 | 6.2 | 5.7 | 8.8 | 7.4 | 43.0 |
| 7 | 2.0 | 2.5 | 3.6 | 28.4 | 85.8 |
| 8 | 2.8 | 3.5 | 4.6 | 24.4 | 65.2 |
| 9 | 6.8 | 8.8 | 14.8 | 29.4 | 116.2 |
| 10 | 4.7 | 5.6 | 8.7 | 20.5 | 86.6 |
| | | | **Average** | **15.6** | **71.4** |

Summarizing, the hybrid model with an ANN black box model of one neuron, tcHM3$_{ANN}$, yielded the most accurate results for predicting the $c_{B,Lys}$ evolution over time, both in the test sets inside and outside the training space. However, for the latter, the predictions of final $c_{B,Lys}$ were worse than in other models, due to not incorporating the time-dependent BSA accumulation and lysozyme membrane fouling effects in the models. This led tcHM3$_{ANN}$ to overpredict $R_{Lys}$ and consequently the final $c_{B,Lys}$ when extrapolating for the initial concentration of both proteins. These errors were additionally accentuated by the fact of not incorporating any correction function for lysozyme accumulation on the membrane. For BSA, on the other hand, even though a correction function was used the models overestimated the final predicted $c_{B,BSA}$. In contrast, tcHM1$_{Raverage}$ yielded stable final $c_{B,Lys}$ predictions for the test sets outside the training space. This was due to its constant $R_{Lys,}$ obtained from the training experiment fitted all independent generated test data in the same way, since the concentration factor was kept comparable. Nevertheless, this model could not make accurate predictions of $c_{B,Lys}$ over time in any test set, due to not updating its $R_{Lys}$.

# 5. Conclusions and outlook

All the models developed throughout this work for the prediction of flux, rejection factor and concentration in two-component cross-flow ultrafiltration systems are summarized in Table 8, with their pros, cons and suggested possible applications.

Table 8: Summary of all developed models with their pros, cons and suggested possible applications. SFM: stagnant film model; ocHM: one-component-hybrid model; tcHM: two-component-hybrid model; ANN: artificial neural network; MLR: multiple linear regression; MnLR: multiple non-linear regression.

| *Model* | *Pros* | *Cons* | *Possible applications* |
|---|---|---|---|
| $SFM_{comb}$ | -Simple, easy to develop (geometrical solution), interpret and use<br><br>-Robust (no prediction variability) | -For pressure-independent region only<br><br>-Assumptions: CP as (only) flux limiting phenomenon, completely solute rejection, constant solution diffusivity, viscosity, density...<br><br>-For one component only. Extension to more components requires:<br>    -Measurement several physical parameters<br>    -Knowledge interaction mechanisms<br>    -Complex models and calculations<br><br>-Poor interpolation capabilities<br><br>-Fails predicting the flux off the design space | Current workhorse in industry due to its simplicity and shown good correlation with rather experimental data |
| $ocHM_{comb}$ | -Both pressure-dependent and independent regions<br><br>-No quantification of impurity<br><br>-No knowledge of underlaying flux governing mechanism necessary<br><br>-Very good flux predictions inside the design space | -For one component only<br><br>-Little insight into functioning of the system<br><br>-Poor extrapolation capabilities especially for variations in the impurity concentration | -One-component solutions: more precise and flexible model than SFM<br><br>-Complex multi-component solutions of many impurities/difficult quantification |
| $tcHM1_{R1}$ | -Two-components. Easy incorporation second component into the model:<br>    -Simple structure adaptation<br>    -No knowledge interactions between components necessary<br>    -No additional sampling step compared to ocHM<br>(of all tcHMs): Flux<br>→ Excellent interpolation capabilities<br>→ Very good extrapolations for CF and higher $c_{B,BSA}$ | -Not incorporates the time-dependent fouling of lysozyme<br>→ Higher errors in flux extrapolations for higher initial $c_{B,Lys}$<br><br>-No $R_{Lys}$ prediction → $c_{B,Lys}$ overestimation | Only when the solute is completely retained by the membrane |
| $tcHM1_{Raverage}$ | -Simplicity. Easy to construct<br>-Interpretability<br>-Robust. No variation in impurity prediction<br>-Faster computation times | -Only prediction of final protein concentration<br><br>-The concentration factor must be comparable to the training set | -Early filtration units where impurity composition is variable<br><br>-Complex multi-component solutions |
| $tcHM2$ | -Good $R_{Lys}$ predictions<br>-Faster computation times | -Wrong flux predictions | This model is not suggested for utilization |

| | | -Not incorporates the time-dependent BSA CP and lysozyme fouling formation | Most precise developed model |
|---|---|---|---|
| *tcHM3<sub>ANN</sub>* | -Best model for dynamic $R_{Lys}$ and $c_{B,Lys}$ prediction<br><br>-Simple offline analytics (based on UV absorbance record) | -Influence of TMP and CF only on $c_{p,Lys}$<br><br>-No lysozyme correction function<br><br>→ $R_{Lyes}$ and final $c_{B,Lys}$ overestimations outside the design space | Applicable to all UF units and systems as long as there are enough offline analytics for the separate quantification of each component |
| *tcHM3<sub>MLR</sub>/*<br>*tcHM3<sub>MnLR</sub>* | -Easier to interpret<br><br>-Faster computation times | -Failed predicting $R_{Lys}$ (non-linearity between $R_{Lys}$ and input parameters) | These models are not suggested for utilization |

Summarizing, the presented multistep-ahead hybrid models, and particularly tcHM3$_{ANN}$, are promising candidates to build digital twins for virtually designing and optimizing processes by varying the input parameters and the product to impurity ratios. In addition, due to predicting the concentration evolution of the mimicked impurity, lysozyme, tcHM3$_{ANN}$ could also be implemented as soft-sensor for real-time monitoring. By just determining the initial concentration of each of the components and the applied CF and TMP parameters, this model could make forecasts of the duration and final product and impurities concentration. Finally, when combined with closed-loop process controllers, this model could be used for model predictive control, taking the adequate measures depending on changes in the process conditions, thereby ensuring the safety and the quality of the product and reducing the risk of batch rejection.

To conclude, the developed models lay the basis for multi-component systems and next generation bioprocessing towards the PAT and QbD initiative. In addition to its extension for more than two-component solutions, these models could also consider the implementation of a black box model for better predicting the product concentration and its quality based on all CQAs.

# 6. Bibliography

1. Q 6 B Specifications: Test Procedures and Acceptance Criteria for Biotechnological/Biological Products. 17 (2006).

2. Research, C. for D. E. and. Q3B(R) Impurities in New Drug Products (Revision 2). *U.S. Food and Drug Administration* https://www.fda.gov/regulatory-information/search-fda-guidance-documents/q3br-impurities-new-drug-products-revision-2 (2020).

3. Hogwood, C. E. M., Tait, A. S., Koloteva-Levine, N., Bracewell, D. G. & Smales, C. M. The dynamics of the CHO host cell protein profile during clarification and protein A capture in a platform antibody purification process. *Biotechnol. Bioeng.* **110**, 240–251 (2013).

4. Guiochon, G. & Beaver, L. A. Separation science is the key to successful biopharmaceuticals. *J Chromatogr A* **1218**, 8836–8858 (2011).

5. Walsh, G. Biopharmaceutical benchmarks 2014. *Nature Biotechnology* **32**, 992–1000 (2014).

6. Kelley, B. DOWNSTREAM PROCESSING OF MONOCLONAL ANTIBODIES: CURRENT PRACTICES AND FUTURE OPPORTUNITIES. in *Process Scale Purification of Antibodies* (ed. Gottschalk, U.) 1–21 (John Wiley &amp; Sons, Inc., 2017). doi:10.1002/9781119126942.ch1.

7. Strube, J., Grote, F., Josch, J. P. & Ditz, R. Process Development and Design of Downstream Processes. *Chemie Ingenieur Technik* **83**, 1044–1065 (2011).

8. Conner, J. *et al.* Chapter 26 - The Biomanufacturing of Biotechnology Products. in *Biotechnology Entrepreneurship* (ed. Shimasaki, C.) 351–385 (Academic Press, 2014). doi:10.1016/B978-0-12-404730-3.00026-9.

9. van Reis, R. & Zydney, A. Bioprocess membrane technology. *Journal of Membrane Science* **297**, 16–50 (2007).

10. Blatt, W. F., Dravid, A., Michaels, A. S. & Nelsen, L. Solute Polarization and Cake Formation in Membrane Ultrafiltration: Causes, Consequences, and Control Techniques. *Membrane Science and Technology* 47–97 (1970) doi:10.1007/978-1-4684-1851-4_4.

11. Wang, Z. *et al.* Membrane cleaning in membrane bioreactors: A review. *Journal of Membrane Science* **468**, 276–307 (2014).

12. Zydney, A. L. Stagnant film model for concentration polarization in membrane systems. *Journal of Membrane Science* **130**, 275–281 (1997).

13. Wijmans, J. G., Nakao, S. & Smolders, C. A. Flux limitation in ultrafiltration: Osmotic pressure model and gel layer model. *Journal of Membrane Science* **20**, 115–124 (1984).

14. van den Berg, G. B. & Smolders, C. A. Flux decline in ultrafiltration processes. *Desalination* **77**, 101–133 (1990).

15. Cross Flow Filtration Method Handbook. https://www.bioprocessonline.com/doc/cross-flow-filtration-method-handbook-0001.

16. Vilker, V. L., Colton, C. K. & Smith, K. A. The osmotic pressure of concentrated protein solutions: Effect of concentration and ph in saline solutions of bovine serum albumin. *Journal of Colloid and Interface Science* **79**, 548–566 (1981).

17. Darcy, H. *Les Fontaines publiques de la ville de Dijon. Exposition et application des principes à suivre et des formules à employer dans les questions de distribution d'eau, etc*. (V. Dalmont, 1856).

18. Wijmans, J. G., Nakao, S., Van Den Berg, J. W. A., Troelstra, F. R. & Smolders, C. A. Hydrodynamic resistance of concentration polarization boundary layers in ultrafiltration. *Journal of Membrane Science* **22**, 117–135 (1985).

19.    Thiess, H., Leuthold, M., Grummert, U. & Strube, J. Module design for ultrafiltration in biotechnology: Hydraulic analysis and statistical modeling. *Journal of Membrane Science* **540**, 440–453 (2017).

20.    Baek, Y., Singh, N., Arunkumar, A. & Zydney, A. Ultrafiltration behavior of an Fc-fusion protein: Filtrate flux data and modeling. *Journal of Membrane Science* **528**, 171–177 (2017).

21.    Den Berg, G. B. V., Hanemaaijer, J. H. & Smolders, C. A. Ultrafiltration of protein solutions; the role of protein association in rejection and osmotic pressure. *Journal of Membrane Science* **31**, 307–320 (1987).

22.    Namila, F. N. U. *et al.* The effects of buffer condition on the fouling behavior of MVM virus filtration of an Fc-fusion protein. *Biotechnology and Bioengineering* **116**, 2621–2631 (2019).

23.    Borujeni, E. E. & Zydney, A. Membrane fouling during ultrafiltration of plasmid DNA through semipermeable membranes. *Journal of Membrane Science* **450**, 189–196 (2014).

24.    Vela, M. C. V., Blanco, S. Á., García, J. L. & Rodríguez, E. B. Analysis of membrane pore blocking models applied to the ultrafiltration of PEG. *Separation and Purification Technology* **62**, 489–498 (2008).

25.    Iritani, E. A Review on Modeling of Pore-Blocking Behaviors of Membranes During Pressurized Membrane Filtration. *Drying Technology* **31**, 146–162 (2013).

26.    Jim, K. J., Fane, A. G., Fell, C. J. D. & Joy, D. C. Fouling mechanisms of membranes during protein ultrafiltration. *Journal of Membrane Science* **68**, 79–91 (1992).

27.    Binabaji, E., Ma, J., Rao, S. & Zydney, A. L. Theoretical analysis of the ultrafiltration behavior of highly concentrated protein solutions. *Journal of Membrane Science* **494**, 216–223 (2015).

28.    Binabaji, E., Ma, J. & Zydney, A. L. Intermolecular Interactions and the Viscosity of Highly Concentrated Monoclonal Antibody Solutions. *Pharm. Res.* **32**, 3102–3109 (2015).

29.    Aimar, P. & Field, R. Limiting flux in membrane separations: A model based on the viscosity dependency of the mass transfer coefficient. *Chemical Engineering Science* **47**, 579–586.

30.    Binabaji, E., Ma, J., Rao, S. & Zydney, A. L. Ultrafiltration of highly concentrated antibody solutions: Experiments and modeling for the effects of module and buffer conditions. *Biotechnol. Prog.* **32**, 692–701 (2016).

31.    Matsuyama, H., Shimomura, T. & Teramoto, M. Formation and characteristics of dynamic membrane for ultrafiltration of protein in binary protein system. *Journal of Membrane Science* **92**, 107–115 (1994).

32.    van den Berg, G. B. & Smolders, C. A. Concentration polarization phenomena during dead-end ultrafiltration of protein mixtures. The influence of solute-solute interactions. *Journal of Membrane Science* **47**, 1–24 (1989).

33.    Teng, M.-Y., Lin, S.-H., Wu, C.-Y. & Juang, R.-S. Factors affecting selective rejection of proteins within a binary mixture during cross-flow ultrafiltration. *Journal of Membrane Science* **281**, 103–110 (2006).

34.    Iritani, E., Mukai, Y. & Murase, T. Separation of binary protein mixtures by ultrafiltration. *Filtration & Separation* **34**, 967–973 (1997).

35.    Müller, C. H., Agarwal, G. P., Melin, T. & Wintgens, T. Study of ultrafiltration of a single and binary protein solution in a thin spiral channel module. *Journal of Membrane Science* **227**, 51–69 (2003).

36.    Saksena, S. & Zydney, A. L. Influence of protein–protein interactions on bulk mass transport during ultrafiltration. *Journal of Membrane Science* **125**, 93–108 (1997).

37.    Shukla, A. A. & Hinckley, P. Host cell protein clearance during protein a

chromatography: Development of an improved column wash step. *Biotechnol Progress* **24**, 1115–1121 (2008).

38.     Huuk, T. C. *et al.* Modeling of complex antibody elution behavior under high protein load densities in ion exchange chromatography using an asymmetric activity coefficient. *Biotechnology Journal* **12**, (2016).

39.     Benner, S. W., Welsh, J. P., Rauscher, M. A. & Pollard, J. M. Prediction of lab and manufacturing scale chromatography performance using mini-columns and mechanistic modeling. *Journal of Chromatography A* **1593**, 54–62 (2019).

40.     Hebbi, V., Roy, S., Rathore, A. S. & Shukla, A. Modeling and prediction of excipient and pH drifts during ultrafiltration/diafiltration of monoclonal antibody biotherapeutic for high concentration formulations. *Separation and Purification Technology* **238**, 116392 (2020).

41.     Haribabu, M., Dunstan, D. E., Martin, G. J. O., Davidson, M. R. & Harvie, D. J. E. Simulating the ultrafiltration of whey proteins isolate using a mixture model. *Journal of Membrane Science* **613**, 118388 (2020).

42.     Glassey, J., Stosch, M. von & Stosch, M. von. *Hybrid Modeling in Process Industries*. (CRC Press, 2018). doi:10.1201/9781351184373.

43.     Solomatine, D. P. & Ostfeld, A. Data-driven modelling: some past experiences and new approaches. *Journal of Hydroinformatics* 20 (2008).

44.     Mitchell, T. M. *Machine Learning*. (McGraw-Hill, 1997).

45.     Kuhn, M. & Johnson, K. *Applied Predictive Modeling*. (Springer New York, 2013). doi:10.1007/978-1-4614-6849-3.

46.     Comparison of Modeling Methods for DoE-Based Holistic Upstream Process Characterization - Bayer - 2020 - Biotechnology Journal - Wiley Online Library. https://onlinelibrary.wiley.com/doi/full/10.1002/biot.201900551.

47.     Stosch, M. von & Willis, M. J. Intensified design of experiments for upstream bioreactors. *Engineering in Life Sciences* **17**, 1173–1184 (2017).

48.     Nagrath, D., Messac, A., Bequette, B. & Cramer, S. A Hybrid Model Framework for the Optimization of Preparative Chromatographic Processes. *Biotechnology progress* **20**, 162–78 (2004).

49.     Chew, C. M., Aroua, M. K. & Hussain, M. A. A practical hybrid modelling approach for the prediction of potential fouling parameters in ultrafiltration membrane water treatment plant. *Journal of Industrial and Engineering Chemistry* **45**, 145–155 (2017).

50.     Grisales Díaz, V. H., Prado-Rubio, O. A., Willis, M. J. & von Stosch, M. Dynamic hybrid model for ultrafiltration membrane processes. in *Computer Aided Chemical Engineering* vol. 40 193–198 (Elsevier, 2017).

51.     Krippl, M., Dürauer, A. & Duerkop, M. Hybrid modeling of cross-flow filtration: Predicting the flux evolution and duration of ultrafiltration processes. *Separation and Purification Technology* **248**, 117064 (2020).

52.     Baeshen, M. N. *et al.* Production of Biopharmaceuticals in E. coli: Current Scenario and Future Perspectives. *J. Microbiol. Biotechnol.* **25**, 953–962 (2015).

53.     Swartz, J. R. Advances in Escherichia coli production of therapeutic proteins. *Current Opinion in Biotechnology* **12**, 195–201 (2001).

54.     Hart, A. & Bailey, E. Protein Composition of Vitreoscillu Hemoglobin Inclusion Bodies Produced in Escherichia coZi. 6.

55.     Veeraragavan, K. Studies on two major contaminating proteins of the cytoplasmic inclusion bodies in *Escherichia coli. FEMS Microbiology Letters* **61**, 149–152 (1989).

56.     Characterization of inclusion bodies in recombinant Escherichia coli producing high

levels of porcine somatotropin.

57.     Rinas, U., Boone, T. C. & Bailey, J. E. Characterization of inclusion bodies in recombinant Escherichia coli producing high levels of porcine somatotropin. *Journal of Biotechnology* **28**, 313–320 (1993).

58.     Dawson, M. For Parenteral Drug Products. 7 (2017).

59.     *Sixty-second report / WHO Expert Committee on Biological Standardization: Geneva, 17 - 21 October 2011*. (World Health Organization, 2013).

60.     Guidance for Industry- Characterization and Qualification of Cell Substrates and Other Biological Materials Used in the Production of Viral Vaccines for Infectious Disease Indications. 50.

61.     Beatson, R. *et al.* Transforming growth factor-β1 is constitutively secreted by chinese hamster ovary cells and is functional in human cells. *Biotechnol. Bioeng.* **108**, 2759–2764 (2011).

62.     Wang, X., Hunter, A. K. & Mozier, N. M. Host cell proteins in biologics development: Identification, quantitation and risk assessment. *Biotechnol. Bioeng.* **103**, 446–458 (2009).

63.     Liu, X. *et al.* Identification and characterization of co-purifying CHO host cell proteins in monoclonal antibody purification process. *Journal of Pharmaceutical and Biomedical Analysis* **174**, 500–508 (2019).

64.     Valente, K. N., Lenhoff, A. M. & Lee, K. H. Expression of difficult-to-remove host cell protein impurities during extended Chinese hamster ovary cell culture and their impact on continuous bioprocessing: Expression of Difficult-to-Remove Host Cell Proteins. *Biotechnol. Bioeng.* **112**, 1232–1242 (2015).

65.     Champion, K., Madden, H., Dougherty, J. & Shacter, E. Defining Your Product Profile and Maintaining Control Over It, Part 2. 5 (2005).

66.     Sauer, D. G. *et al.* A two-step process for capture and purification of human basic fibroblast growth factor from E. coli homogenate: Yield versus endotoxin clearance. *Protein Expression and Purification* **153**, 70–82 (2019).

67.     Omasa, T., Onitsuka, M. & Kim, W.-D. Cell engineering and cultivation of chinese hamster ovary (CHO) cells. *Curr Pharm Biotechnol* **11**, 233–240 (2010).

68.     Krippl, M., Bofarull-Manzano, I., Duerkop, M. & Dürauer, A. Hybrid Modeling for Simultaneous Prediction of Flux, Rejection Factor and Concentration in Two-Component Crossflow Ultrafiltration. *Processes* **8**, 1625 (2020).

# 7. List of Figures

# 8. List of Tables

# 9. Appendix

Table A1 gives a summary of all test sets that were performed to assess the errors of the developed hybrid models. It includes the process parameters (TMP and CF), the initial protein concentrations in the bulk ($c_{0,BSA}$ and $c_{0,Lys}$) as well as the initial ($V_0$) and final bulk volume ($V_f$) and the numbers of samples taken during each test run.

Table A1: Summary of test set parameters.

| Test set number | TMP | CF | $c_{0,BSA}$ [g/L] | $c_{0,Lys}$ [g/L] | $V_0$ [mL] | $V_f$ [mL] | Number of samples taken |
|---|---|---|---|---|---|---|---|
| 1 | 1.8 | 200 | 6.7 | 0.00 | 1031.3 | 38.3 | 12 |
| 2 | 1.8 | 200 | 4.0 | 0.28 | 1031.4 | 41.2 | 12 |
| 3 | 2.8 | 300 | 3.8 | 0.32 | 1081.4 | 51.0 | 4 |
| 4 | 2.1 | 250 | 3.7 | 0.38 | 1081.4 | 50.8 | - |
| 5 | 2.5 | 280 | 4.6 | 0.25 | 1081.4 | 44.6 | 12 |
| 6 | 1.8 | 200 | 3.8 | 0.50 | 1031.4 | 57.1 | 15 |
| 7 | 1.6 | 230 | 6.0 | 0.15 | 1081.4 | 44.2 | 11 |
| 8 | 1.4 | 270 | 8.8 | 0.19 | 1081.4 | 43.8 | 9 |
| 9 | 1.8 | 260 | 2.4 | 0.57 | 1072.9 | 41.1 | 8 |
| 10 | 2.0 | 350 | 3.6 | 0.34 | 1081.4 | 41.8 | 8 |

In Table A2 the mass transfer coefficient k and gel concentration $c_G$ for the flux prediction using the SFM are summarized. Due to the SFM is based on one-component only, and that the BSA concentration was 4 to 46 times higher than the lysozyme concentration, BSA was chosen as the modeled component. k and $c_G$ for BSA were calculated both in the BSA alone (Table A2 left) and the BSA with lysozyme (Table A2 right) training data. k from the two-component solution was generally lower than for one-component, as a result of the fouling of lysozyme, which reduced the solute mass transfer along the membrane.

Table A2: Mass transfer coefficient k and gel concentration $c_G$ for the SFM based on the BSA alone (left) and BSA with lysozyme (right) training data. Both parameters were in general smaller in the two-component solution than for one-component, due to the reduced transmembrane mass transfer as a result of lysozyme fouling.

| | | k based on BSA | | | | | k based on BSA with lysozyme | | |
|---|---|---|---|---|---|---|---|---|---|
| | | Feedflow [mL/min] | | | | | Feedflow [mL/min] | | |
| | | 100 | 200 | 300 | | | 100 | 200 | 300 |
| TMP [bar] | 0.8 | 47.36 | 36.23 | 31.63 | TMP [bar] | 0.8 | 17.61 | 33.84 | 14.03 |
| | 1.3 | 54.67 | 41.97 | 32.33 | | 1.3 | 25.03 | 36.12 | 27.92 |
| | 1.8 | 54.51 | 42.14 | 30.13 | | 1.8 | 27.59 | 36.75 | 39.06 |
| | 2.3 | 53.40 | 41.63 | 28.99 | | 2.3 | 28.11 | 38.22 | 44.73 |
| | 2.8 | 53.75 | 42.97 | 27.95 | | 2.8 | 27.70 | 38.58 | 46.84 |
| | | $c_G$ based on BSA | | | | | $c_G$ based on BSA with lysozyme | | |
| | | Feedflow [mL/min] | | | | | Feedflow [mL/min] | | |
| | | 100 | 200 | 300 | | | 100 | 200 | 300 |
| TMP [bar] | 0.8 | 277.83 | 303.41 | 279.25 | TMP [bar] | 0.8 | 665.40 | 280.12 | 4421.42 |
| | 1.3 | 302.79 | 330.45 | 323.30 | | 1.3 | 355.06 | 312.56 | 887.24 |

| 1.8 | 322.39 | 345.88 | 355.30 | 1.8 | 288.49 | 277.52 | 419.22 |
| 2.3 | 332.29 | 353.74 | 369.46 | 2.3 | 264.69 | 273.21 | 304.56 |
| 2.8 | 323.99 | 327.45 | 378.28 | 2.8 | 256.72 | 252.80 | 263.97 |

In Table A3 the NRMSE for flux prediction by all developed models (excepting $tcHM3_{MLR}$ and $tcHM3_{MnLR}$) for each test set is shown.

Table A3: Summary of flux NRMSE for all test sets by all the hybrid (HM) and stagnant film (SFM) models built during this work.

| | | | | NRMSE flux [%] | | | |
|---|---|---|---|---|---|---|---|
| Test set number | $ocHM_{BSA}$ | $SFM_{comb}$ | $ocHM_{comb}$ | $tcHM1_{R1}$ | $tcHM1_{Raverage}$ | $tcHM2$ | $tcHM3_{ANN}$ |
| 1 | 1.8 ± 0.2 | 9.7 | 9.0 ± 0.1 | 1.6 ± 0.4 | 1.7 ± 0.5 | 11.9 ± 1.4 | 1.7 ± 0.4 |
| 2 | 8.8 ± 0.4 | 5.3 | 2.2 ± 0.2 | 2.0 ± 0.2 | 2.3 ± 0.3 | 2.0 ± 0.2 | 2.1 ± 0.3 |
| 3 | 11.1 ± 0.3 | 3.0 | 2.3 ± 0.2 | 4.1 ± 0.2 | 3.6 ± 0.2 | 4.6 ± 0.3 | 3.2 ± 0.2 |
| 4 | 10.5 ± 0.2 | 3.0 | 1.8 ± 0.1 | 3.6 ± 0.2 | 3.3 ± 0.3 | 3.7 ± 0.6 | 3.0 ± 0.2 |
| 5 | 7.2 ± 0.2 | 5.0 | 4.9 ± 0.2 | 4.1 ± 0.1 | 3.9 ± 0.3 | 3.4 ± 0.4 | 3.7 ± 0.2 |
| 6 | 12.1 ± 0.3 | 7.2 | 4.5 ± 0.1 | 3.8 ± 0.2 | 4.3 ± 0.2 | 4.1 ± 0.6 | 4.1 ± 0.1 |
| 7 | 2.3 ± 0.2 | 8.5 | 8.5 ± 0.2 | 5.0 ± 0.3 | 4.8 ± 0.3 | 3.6 ± 0.4 | 4.9 ± 0.4 |
| 8 | 2.5 ± 0.3 | 9.3 | 6.7 ± 0.2 | 3.3 ± 0.3 | 3.4 ± 0.3 | 5.2 ± 0.9 | 3.3 ± 0.5 |
| 9 | 10.9 ± 0.2 | 6.3 | 5.0 ± 0.1 | 8.0 ± 0.2 | 8.2 ± 0.2 | 9.5 ± 0.4 | 8.0 ± 0.3 |
| 10 | 7.9 ± 0.4 | 4.4 | 2.7 ± 0.3 | 4.9 ± 0.2 | 4.7 ± 0.3 | 5.4 ± 0.4 | 4.5 ± 0.4 |

Table A4 contains the NRMSE for $R_{Lys}$ prediction of all developed two-component-hybrid models (tcHMs).

Table A4: Summary of $R_{Lys}$ NRMSE.

| | | | NRMSE $R_{Lys}$ [%] | | | |
|---|---|---|---|---|---|---|
| Test set number | $tcHM1_{R1}$ | $tcHM1_{Raverage}$ | $tcHM2$ | $tcHM3_{ANN}$ | $tcHM3_{MLR}$ | $tcHM3_{MnLR}$ |
| 2 | 91.2 ± 0.0 | 32.5 ± 0.0 | 4.7 ± 0.1 | 6.2 ± 0.2 | 9.9 ± 0.0 | 8.7 ± 0.0 |
| 3 | 70.2 ± 0.0 | 40.3 ± 0.0 | 25.3 ± 0.6 | 20.4 ± 0.5 | 20.9 ± 0.0 | 21.8 ± 0.0 |
| 5 | 70.1 ± 0.0 | 34.9 ± 0.0 | 4.6 ± 0.4 | 6.9 ± 0.6 | 21.1 ± 0.0 | 10.2 ± 0.0 |
| 6 | 80.0 ± 0.0 | 39.7 ± 0.0 | 20.8 ± 1.4 | 14.7 ± 2.2 | 26.9 ± 0.0 | 30.0 ± 0.0 |
| 7 | 75.3 ± 0.0 | 32.6 ± 0.0 | 19.6 ± 2.6 | 25.7 ± 1.3 | 34.5 ± 0.0 | 29.4 ± 0.0 |
| 8 | 65.9 ± 0.0 | 45.3 ± 0.0 | 34.9 ± 5.6 | 24.5 ± 1.3 | 69.4 ± 0.0 | 23.8 ± 0.0 |
| 9 | 89.8 ± 0.0 | 35.6 ± 0.0 | 45.6 ± 0.8 | 10.0 ± 2.0 | 55.5 ± 0.0 | 89.9 ± 0.0 |
| 10 | 78.3 ± 0.0 | 40.6 ± 0.0 | 11.8 ± 0.5 | 6.2 ± 0.7 | 20.0 ± 0.0 | 34.5 ± 0.0 |
| **Average** | **77.6 ± 0.0** | **37.7 ± 0.0** | **20.9 ± 1.5** | **14.3 ± 1.1** | **32.3 ± 0.0** | **31.0 ± 0.0** |

Table A5 contains the difference, in percentage, between the measured and predicted final $c_{B,Lys}$ by all two-component-hybrid models for each test set.

Table A5: Summary of the difference, in percentage, between measured and predicted final $c_{B,Lys}$ for each test set by the different tcHM.

| | Final $c_{B,Lys}$ difference [%] | | | | | |
|---|---|---|---|---|---|---|
| Test set number | tcHM1$_{R1}$ | tcHM1$_{Raverage}$ | tcHM2 | tcHM3$_{ANN}$ | tcHM3$_{MLR}$ | tcHM3$_{MnLR}$ |
| 2 | 56.4 | 25.8 | 3.9 | 1.8 | 7.0 | 6.7 |
| 3 | 51.6 | 25.3 | 1.7 | 5.8 | 0.9 | 5.5 |
| 5 | 66.8 | 20.3 | 1.6 | 7.3 | 13.2 | 1.1 |
| 6 | 43.0 | 27.6 | 8.4 | 7.4 | 15.0 | 15.5 |
| 7 | 85.8 | 11.2 | 33.6 | 28.4 | 22.3 | 23.0 |
| 8 | 65.2 | 21.3 | 46.5 | 24.4 | 31.3 | 10.3 |
| 9 | 116.2 | 0.9 | 67.7 | 29.4 | 85.2 | 30.6 |
| 10 | 86.6 | 12.3 | 26.0 | 20.5 | 36.6 | 13.0 |
| **Average** | **71.4** | **18.1** | **23.7** | **15.6** | **26.4** | **13.2** |

Figure A1 gives the entire training data set of Figure 9, for the recorded fluxes of each training experiment for the combination of $c_{B,i}$ and TMP at 200 mL/min



Figure A1: Training data sets including different protein concentration and TMPs at CF 200 mL/min. (A), (B) Multi-component training set containing BSA and lysozyme (blue) in the same solution and single component solution of (C) BSA (red) and (D) lysozyme (green). Figure from [68].

## 9.1 Manuscript

The results obtained in this work were published as a research article in the journal *Processes,* in the volume 8, issue 12, article 1625 (2020), as "*Hybrid modeling for Simultaneous Prediction of Flux, Rejection Factor and Concentration in Two-Component Crossflow Ultrafiltration*", by the authors Maximilian Krippl, Ignasi Bofarull-Manzano, Mark Duerkop and Astrid Dürauer.

https://doi.org/10.3390/pr8121625

The manuscript of such publication is attached in the following.

# Hybrid Modeling for Simultaneous Prediction of Flux, Rejection Factor and Concentration in Two-Component Crossflow Ultrafiltration

**Maximilian Krippl [1], Ignasi Bofarull-Manzano [1], Mark Duerkop [1,2] and Astrid Dürauer [1,*]**

[1]   Department of Biotechnology, Institute of Bioprocess Science and Engineering, University of Natural Resources and Life Sciences, 1190 Vienna, Austria; maximilian.krippl@boku.ac.at (M.K.); ignasi.bofarull-manzano@boku.ac.at (I.B.-M.); mark.duerkop@boku.ac.at (M.D.)

[2]   Novasign GmbH, 1190 Vienna, Austria

*   Correspondence: astrid.duerauer@boku.ac.at; Tel.: +43-1476-5479-095

**Abstract:** Ultrafiltration is a powerful method used in virtually every pharmaceutical bioprocess. Depending on the process stage, the product-to-impurity ratio differs. The impact of impurities on the process depends on various factors. Solely mechanistic models are currently not sufficient to entirely describe these complex interactions. We have established two hybrid models for predicting the flux evolution, the protein rejection factor and two components' concentration during crossflow ultrafiltration. The hybrid models were compared to the standard mechanistic modeling approach based on the stagnant film theory. The hybrid models accurately predicted the flux and concentration over a wide range of process parameters and product-to-impurity ratios based on a minimum set of training experiments. Incorporating both components into the modeling approach was essential to yielding precise results. The stagnant film model exhibited larger errors and no predictions regarding the impurity could be made, since it is based on the main product only. Further, the developed hybrid models exhibit excellent interpolation properties and enable both multi-step ahead flux predictions as well as time-resolved impurity forecasts, which is considered to be a critical quality attribute in many bioprocesses. Therefore, the developed hybrid models present the basis for next generation bioprocessing when implemented as soft sensors for real-time monitoring of processes.

**Keywords:** semi-parametric model; neural network; tangential flow filtration; downstream processing; advanced process monitoring

## 1. Introduction

Membrane separation is a unit operation used in virtually all bioprocesses. One prominent type, crossflow ultrafiltration, is widely used from cell harvest and virus clearance approaches to product concentration steps. In downstream processing of biopharmaceuticals, ultrafiltration (UF) is commonly applied for concentration and buffer exchange after the capture step. It is also applied after virus filtration in single-pass mode to concentrate the sample before it is loaded onto the polishing chromatography, or after polishing to reach the final conditions for product formulation [1]. These process steps entail varying ratios of process and impurities to product concentration.

Modeling of process steps is of increasing importance for bioprocesses. Such process models increase understanding of processes, facilitate the discovery of optimal process conditions and are indispensable for model predictive control. The latter is a cornerstone of Quality by Design and Process Analytical Technology, which is recommended by authorities for biopharmaceutical production. The right balance of model complexity and usability is crucial to employ such models effectively for different unit operations.

To simplify the modeling of downstream processes, a common assumption is to reduce the overall sample composition down to a single target molecule. Coefficients and parameters used in mechanistic models, such as mass transfer models, are often approximated, taking only the target molecule into account. Such models may be limited if the sample contains high levels of impurity.

For some process steps, such as polishing chromatography [2] or ultra/diafiltration [3,4] before formulation, this assumption of one-component solutions is realistic, since the product is already of high purity at this process stage. For earlier process steps, however, this simplification deviates substantially from reality and can lead to erroneous models, e.g., for filtration steps after the capture step. Here, the neglected presence of host cell proteins [5], DNA [6], or protein aggregates [7] can strongly distort the prediction of the model, since effects like membrane fouling and interactions between the product and impurities are not considered. In more complex mechanistic models, if the impurities are well characterized, such effects can be considered. For example, for crossflow filtration, a hard sphere-based mixture model, including multiphase computational fluid dynamics and concentration polarization, was applied to a whey protein solution, leading to a permeate flux prediction error within 20% [8]. Other work has shown that mechanistic models of pore blockage and cake filtration can also predict filter fouling during virus filtration, as a function of the protein of interest, virus and membrane [9]. The initial and late stage of the filtration, however, was dominated by different mechanisms, rendering it difficult to build a valid model for the entire process. The influence of two components on (crossflow) UF was found to affect the process in different ways, from strong [10] to weak [11] to varying [5,12–14] protein-protein (or protein-membrane) interactions. To account for the highly different effects of all components on the process, the experimental part of data generation to estimate the parameters for mechanistic models might become very labor-intensive and the calculations rather complex. Further, if the overall behavior of the process changes because of varying concentrations of impurities, the assumptions of mechanistic models might not hold, to the detriment of the prediction.

One advantage of machine learning supported modeling approaches is that the effects of the impurity on flux and membrane fouling do not need to be fully quantified by the operator [15]. The quantification of these phenomena is performed by machine learning tools, such as an artificial neural network (ANN) [16]. Hybrid models combine the advantages of data-driven black box models (such as ANNs), correlating input with output variables (such as the concentration of impurity with the decrease in flux) with knowledge-based mechanistic models (white box models) derived from conservation of kinetic laws [17]. Hybrid models have been applied to bioprocesses for upstream [18] and downstream applications [19,20].

To compare the predictive power of a model concerning the training space, two terms are often used: interpolation and extrapolation. Interpolation allows the model to make predictions for parameters that lie within the range of training experiments. A model with good interpolation capabilities can make predictions with fewer training observations, since it is able to make reliable estimates of the spaces between the observations. A model with poor interpolation capabilities requires more granular coverage of the training space to make accurate predictions of test experiments. Extrapolation (also called range extrapolation) describes the extent to which a model can make predictions if the tested input parameters are outside the training space. A model with good extrapolation capabilities can make accurate predictions for parameters beyond the training space. A detailed explanation of interpolation and extrapolation in hybrid modeling is given in [21].

Recently, we have shown the benefits of hybrid modeling for the prediction of UF flux evolution. However, this previous model was only established for a one-component system [22]. In the present study, we extended the hybrid model to describe the impact of a modeled protein impurity on the decrease of the permeate flux over time in crossflow UF including the rejection behavior of the product and the impurity. This enables the operator to gain a more detailed understanding of the process. In addition, the impurity concentration is a critical quality attribute (CQA) in almost all manufacturing bioprocesses and if it is too high, the produced batch must be discarded. The presented hybrid models can predict the impurity concentration up front, and potentially minimizes the risk of batch rejection.

Product and impurity were mimicked with different ratios of bovine serum albumin (BSA) to lysozyme concentrations in the starting solution. BSA and lysozyme exhibit different physicochemical properties to facilitate separation and quantification. While BSA was fully retained by the membrane, lysozyme was only partially retained, rendering the predictions of the permeate flux over time even more complex. In a first assessment, we compared the abilities of the well-established mechanistic stagnant film model (SFM) and the recently established one-component hybrid to predict the filtration progress of a two-component solution. Further, we presented two hybrid model structures to predict the evolution of permeate flux and protein concentration of product and impurity by multi-step ahead predictions. One hybrid model included a static lysozyme rejection factor ($R_{Lys}$), while the other updated $R_{Lys}$ dynamically in an iterative way. These model outputs were influenced by the transmembrane pressure (TMP), crossflow velocity (CF), the initial BSA concentration $c_{B,BSA}$ and lysozyme bulk concentration $c_{B,Lys}$. Finally, these novel hybrid model structures were compared to the SFM regarding flux and concentration prediction.

## 2. Materials and Methods

### 2.1. Equipment and Chemicals

All UF experiments were performed on an ÄKTA Crossflow system (Cytiva, Marlborough, MA, USA) controlled by UNICORN 5.31 software. The reservoir tank held up to 1100 mL of bulk solution. The system featured an inline pH probe and UV monitor on the permeate side and a pressure-based reservoir level sensor. The experiments were performed with a Sartocon Slice Hydrosart Cassette hydrophilic, stabilized cellulose-based membrane (Sartorius AG, Göttingen, Germany) with a membrane area of 200 cm$^2$. The model proteins were BSA and lysozyme (A2153 and L6876, both purchased from Sigma-Aldrich, St. Louis, MO, USA). The molecular weight cutoff (MWCO) of the membranes was 30 kDa, chosen so that BSA (66 kDa) was fully retained and lysozyme (14 kDa) was partially retained. BSA and lysozyme were chosen to mimic the protein of interest and process-related impurities, respectively. A filtration buffer of 50 mM phosphate-buffered saline (PBS), pH 8, was used.

### 2.2. Training and Test Data Generation

For the training experiments, the bulk reservoir was filled with 1000 mL of the lowest bulk BSA and/or lysozyme concentration $c_{B,BSA}$ and $c_{B,Lys}$ (see Table A1). The following two steps were then alternated. First, the TMP and CF were increased stepwise, while the permeate was redirected to the feed reservoir to keep the protein concentration $c_B$ constant. For each combination of TMP and CF, the permeate flux was recorded. Second, the sample was concentrated until the next desired $c_B$ was reached. These two steps were repeated at all concentrations given in Table A1. A total of 90 equilibrium fluxes were recorded for different concentrations and combinations of TMP and CF. Our previous work with a one-component system [22] showed that this training set size was sufficient to develop a well-trained hybrid model with accurate flux predictions. A detailed summary of all scouted TMPs, CFs, $c_B$s and recorded fluxes is given in Table A1. Samples were taken after each concentration step for offline measurement. A more detailed description of the methodology for the training experiment is given in an earlier publication [22].

During the test experiments, samples were taken from the retentate and permeate. The measured retentate and permeate concentrations were used to calculate the rejection factor R of the model proteins. A summary of the performed test sets is provided in Table A2.

### 2.3. Concentration Polarization Correction

When concentrating the sample throughout the training experiments, we observed that the measured $c_{B,BSA}$ was lower than the expected concentration calculated from permeate volume ($V_P$) and mass balance. The difference between observed and calculated concentration increased with concentration (see Figure A3B). This was because the concentration polarization (CP) layer—the

protein gradient that forms on the surface of the membrane—increased with $c_{B,BSA}$. This deviation was considered for the test experiments by employing a quadratic polynomial function (Equation A1) and used to correct the calculated $c_{B,BSA}$.

*2.4. Protein Quantification*

BSA and lysozyme concentrations were determined with an analytical high-performance size-exclusion chromatography (SEC-HPLC) using a TSKgel G3000SWXL column (5 μm, 7.8 × 300 mm; TOSOH, Shiba, Tokyo, Japan). The separation was performed under isocratic conditions with 50 mM sodium phosphate, 200 mM NaCl, pH 6.5 as running buffer at a flow rate of 0.4 mL/min. Samples were diluted to a final concentration of 0.1 to 1.0 g/L using 50 mM PBS, pH 8 and filtered through a 0.22 μm Millex-GV Filter (Merck Millipore, Billerica, MA, USA) prior to analysis. The injection volume was 10 μL per sample. Due to the difference in the size of BSA and lysozyme, the peaks were fully separated and could be quantified independently, using standard calibrations from BSA and lysozyme stock solutions.

*2.5. Hybrid Modeling*

2.5.1. Black Box Model

The black box inside the first hybrid model (HM 1) aimed to predict the flux based on the combination of inputs parameters: TMP, CF and the bulk protein concentrations of BSA and lysozyme, $c_{B,BSA}$ and $c_{B,Lys}$, respectively. In the second hybrid model (HM 2), an additional black box was employed to predict the rejection factor of lysozyme $R_{Lys}$ (Figure 1B). An ANN was utilized for this purpose and optimized by varying the number of hidden nodes from 1 to 7. The ANN was set up with the feedforwardnet function and trained with the trainbr function, using MATLAB 2018b. A detailed description is the ANN structure and optimizer function is given in the Appendix A.

2.5.2. White Box Model

The white box model is the mechanistic part of the hybrid model and consisted of a mass balance. The incrementally decreasing bulk volume ($dV_B$ in Equation (1)) was derived from the permeate flux (J), which is the output of the black box, and the membrane area (A). The rejection factor R for component i was calculated with Equation (2), considering the concentration of i in both the retentate ($c_R$; in crossflow filtration $c_R$ is equal to $c_B$) and the permeate ($c_P$). Equation (1) and Equation (2) were used to predict $c_B$ of each component and Equation (3) to calculate the $V_B$ after dt.

$$\frac{dV_B}{dt} = -A \cdot J \tag{1}$$

$$R_i = 1 - \frac{c_{P,i}}{c_{R,i}} \tag{2}$$

$$\frac{d(c_{B,i} \cdot V_B)}{dt} = (A \cdot J \cdot c_{B,i})^{R_i} \tag{3}$$

2.5.3. Training and Test Data

$R_{Lys}$ was calculated from the training set with the UV absorbance at 280 nm on the permeate side. A separate lysozyme training run was performed to correlate the UV signal at 280 nm with the permeate concentration determined by SEC-HPLC. The correlation curve (Figure A3A) with an $R^2$ of 0.97 was used to calculate $c_{P,Lys}$, and subsequently $R_{Lys}$ for all observations of the training set was used to train the black box.

The observed flux and $R_{Lys}$ were compared to the predictions of the hybrid models using the normalized root-mean-square error (NRMSE)

$$\text{NRMSE} = 100 \cdot \frac{\sqrt{\frac{1}{n} \sum_{i=1}^{n} (y_i - \hat{y}_i)^2}}{y_{max} - y_{min}} \tag{4}$$

where $n$ is the number of overserved fluxes $y_i$ and the corresponding predicted fluxes $\hat{y}_i$. The normalization $y_{max} - y_{min}$ allows a fair comparison of various fluxes due to different concentrations and process parameters.
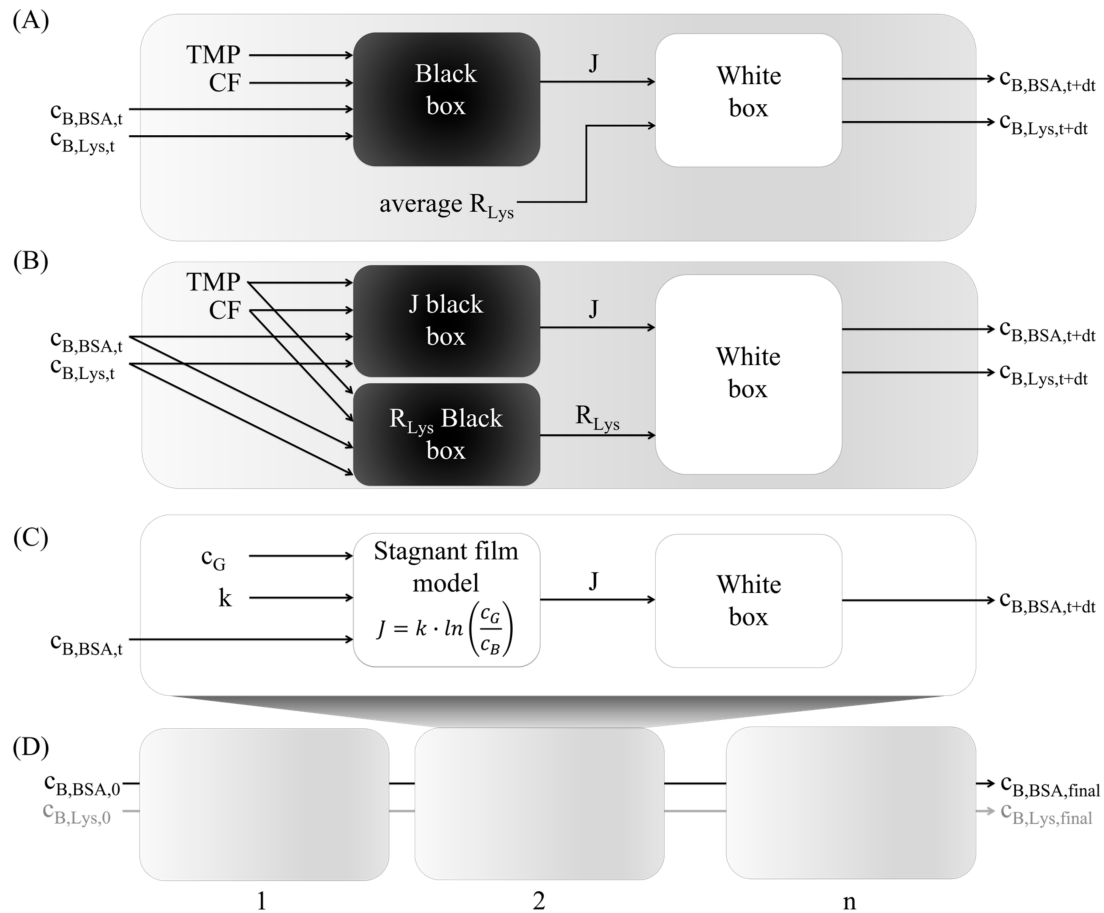


**Figure 1.** Schematic representation of the two hybrid model and mechanistic model structures, with implementation in the multi-step ahead model. (**A**) Hybrid model 1 (HM 1) using static, average $R_{Lys}$ from the training set, (**B**) hybrid model 2 (HM 2) with two separate black boxes for flux and dynamic $R_{Lys}$ prediction, (**C**) stagnant film model (SFM). (**D**) Multi-step ahead hybrid model structure.

2.5.4. Multistep-Ahead Hybrid Model

The structures of the investigated hybrid model are given in Figure 1. The first and simplest HM 1 (Figure 1A) assumed a constant $R_{Lys}$ of 0.77 for all test sets based on the weighted average of all permeate and retentate concentrations samples taken throughout the training experiment. The weighted average considered the variation in $c_{B,Lys}$ and sample intervals using trapezoid rule integration. For the second HM 2 structure (Figure 1B), the flux and $R_{Lys}$ were predicted separately, using two different black box models. The flux and $R_{Lys}$ were fed into the same white box model, which yielded the predicted $c_{B,BSA}$ and $c_{B,Lys}$ after a defined time increment. The developed hybrid model is capable of predicting multiple steps ahead, as depicted in Figure 1D. The multi-step ahead structure uses HM 1, HM 2 or the SFM to predict $c_{B,BSA}$ and $c_{B,Lys}$ for a time increment (dt). The concentrations of the first iteration

were used to calculate future fluxes and $c_B$s of the second iteration, and so on. Multiple iterations were performed until the desired stop criterion was reached. In our case, the stop criterion was the final retentate volume.

The presented hybrid models were used to predict the evolution of flux and $R_{Lys}$ throughout the UF process. Furthermore, the models yielded a prediction for the final $c_{B,BSA}$ and $c_{B,Lys}$. The final $c_{B,BSA}$ and $c_{B,Lys}$ predictions were compared to the final $c_{B,BSA}$ and $c_{B,Lys}$ measured by SEC-HPLC. The model errors were compared using the NRMSE.

Figure 2 shows a flowchart of the hybrid model methodology applied for crossflow filtration. Training experiments were performed by variations in the parameters that are expected to influence the flux. Following this, the model was trained on this training set with a defined experimental design space. The established models were applied to a validation data set that was not used for training. The model structure was optimized by varying the tuning parameters, e.g., number of nodes in an ANN and adding or removing training parameters. The model with the tuning parameters that led to the lowest error in the validation set was then applied to independent test runs with static process conditions.
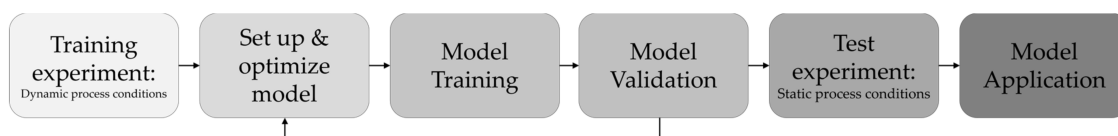


**Figure 2.** Flowchart of the hybrid model methodology for application in crossflow filtration.

2.5.5. Stagnant Film Theory

The presented hybrid models were compared to the established SFM. The SFM derives predictions from the mass transfer model described by convective transport toward the membrane and back-diffusion caused by the concentration gradient [23]. According to the SFM, the flux J is related to the bulk concentration $c_B$ of a single component by

$$J = k \cdot \ln\left(\frac{c_G}{c_B}\right) \tag{5}$$

where $c_G$ is the gel layer concentration at the membrane surface and k is the mass transfer coefficient that depends on the diffusion coefficient and the thickness of the gel layer [23]. The SFM is valid in the pressure-independent region of the filtration. Since k and $c_G$ cannot be adjusted directly during the filtration, a correlation between the adjustable parameters TMP and CF, and k and $c_G$ was required. When plotting the flux versus log($c_B$) for a constant TMP and CF, k and $c_G$ are estimated by the slope of linear regression and $c_G$ was estimated by extrapolating the regression line to the intersection with the abscissa (Figure A6). It has been shown that this way of calculating k yields more accurate results than the Sherwood correlation [24–26] and more solid predictions compared to the osmotic pressure model [27] for similar settings. To compare the SFM to the hybrid models, the black box was replaced by the SFM in Equation (5) using the parameters k and $c_G$ instead of TMP and CF (Figure 1C). In test runs, where the TMP and CF conditions were not covered in the training set, k and $c_G$ were estimated using linear interpolation.

## 3. Results and Discussion

### 3.1. Training Data Description

The data sets for training the hybrid models were generated from filtering BSA and lysozyme with a 30 kDa MWCO cellulose-based membrane (Hydrosart). A total of three training sets were generated covering three CFs (100, 200 and 300 mL/min) and five TMPs (0.8, 1.3, 1.8, 2.3 and 2.8 bar). The three training sets containing BSA, lysozyme and a combination of both are shown in Figure 3. In the

combined training set, the protein concentration of BSA $c_{B, BSA}$ ranged from 3.77 g/L to 77.93 g/L and of lysozyme $c_{B, Lys}$ from 0.28 g/L to 3.81 g/L. The concentration ranges for all training sets are summarized in Table A1. For a better comparison of Figure 3A–D, the *x*-axis of Figure 3C,D are reduced. The entire graphs are given in Figure A2.
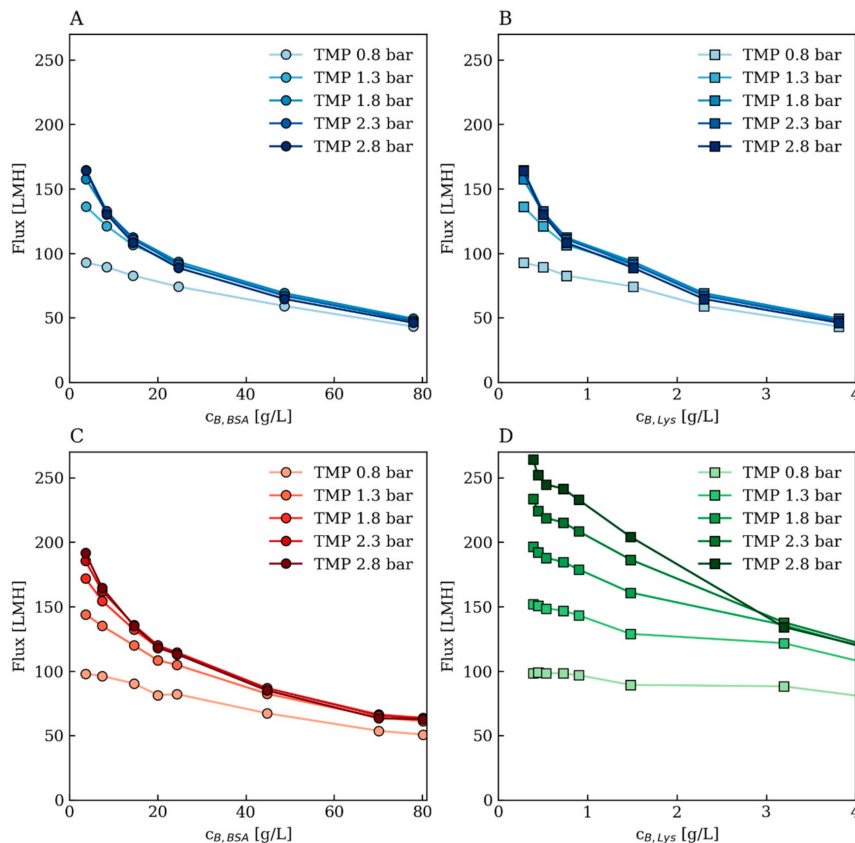


**Figure 3.** Training data sets including different protein concentrations and different TMPs at CF 200 mL/min: two-component training set containing (**A**) BSA and (**B**) lysozyme in the same solution (blue); one-component solution of (**C**) BSA (red) and (**D**) lysozyme (green).

Generally, increasing bulk concentrations $c_B$ led to lower fluxes, while increasing TMP and CF led to higher fluxes in all training sets. This is in accordance with the underlying mechanisms: higher bulk concentrations lead to higher concentrations in the boundary layer and a more prominent effect of the back diffusion along the concentration gradient. An increased TMP leads to higher convective flow towards the membrane, but also to a faster accumulation of the protein at the boundary layer. High CF decreases the thickness of the concentration polarization layer by rectangular displacement. The training set obtained from experiments using only BSA exhibited higher fluxes then the two-component training set. Additionally, the flux decreased faster during filtration of the two-component solution (Figure 3B) compared to the filtration of lysozyme only (Figure 3D). This indicated an increased membrane resistance through the fouling effect on the Hydrosart membrane caused by lysozyme. Being smaller than the pores, lysozyme adsorbed at the inner pore channels [28,29] and reduced its diameter and subsequently the flux through the membrane and the membrane's selectivity.

The two-component training set (Figure 3A,B) was used to train the black box of the hybrid models and to obtain the mechanistic model parameters k and $c_G$. The data set with lysozyme solely (Figure 3D) was used for two reasons: first, to investigate the effect of TMP and CF on the permeability of lysozyme and whether $R_{Lys}$ had to be recalculated for varying input parameters (Figure A5); second, to correlate the permeate lysozyme concentration with the UV signal on the permeate side. This correlation was used to calculate $R_{Lys}$ (Equation (2)) for each observation of the combined training set (Figure 3A,B),

using solely the permeate UV signal. Another training experiment was performed with BSA solely (Figure 3C). The observed fluxes and estimated SFM parameters k and $c_G$ were used to investigate model behavior and error when lysozyme was present in the test set but absent in the training set.

## 3.2. Comparison of the Hybrid Models to the Stagnant Film Theory

The optimal ANN structure in the hybrid models was determined by varying the number of hidden nodes from one to seven and recording the average error of 20 repetitions on the training set. The ANN with four hidden nodes yielded the lowest NRMSE for both HM 1 and HM 2, with an average of 3.4% NRMSE. Higher numbers of hidden nodes led to an error increase due to training set overfitting (Figure A1).

With the SFM, the flux can only be modeled for a one-component system; no adaptations for a two- or multi-component system have been published in the literature so far. In the following, BSA was assumed to be the only component since its concentration was four to 46 times higher than lysozyme in the test runs (Table A2). The k and $c_G$ values of BSA, however, change in the presence of lysozyme. To allow a fair comparison between the hybrid models (which can incorporate multiple components as inputs) and the SFM, both sets of k and $c_G$ were evaluated. Both experiments were carried out with BSA alone. The combination of BSA with lysozyme was used for flux prediction and the results were compared to the prediction of the hybrid models.

The hybrid model trained solely on BSA (Figure 4, red dotted line) and the SFM using k and $c_G$ based solely on BSA (Figure 4, dark grey dot-dashed line) were able to predict a UF process with only BSA present (Figure 4A, black line), but failed to predict the UF flux of BSA and lysozyme (Figure 4B, black line). The latter failed due to membrane fouling by lysozyme and therefore the reduced flux and prolonged process times could not be described by any of these models.
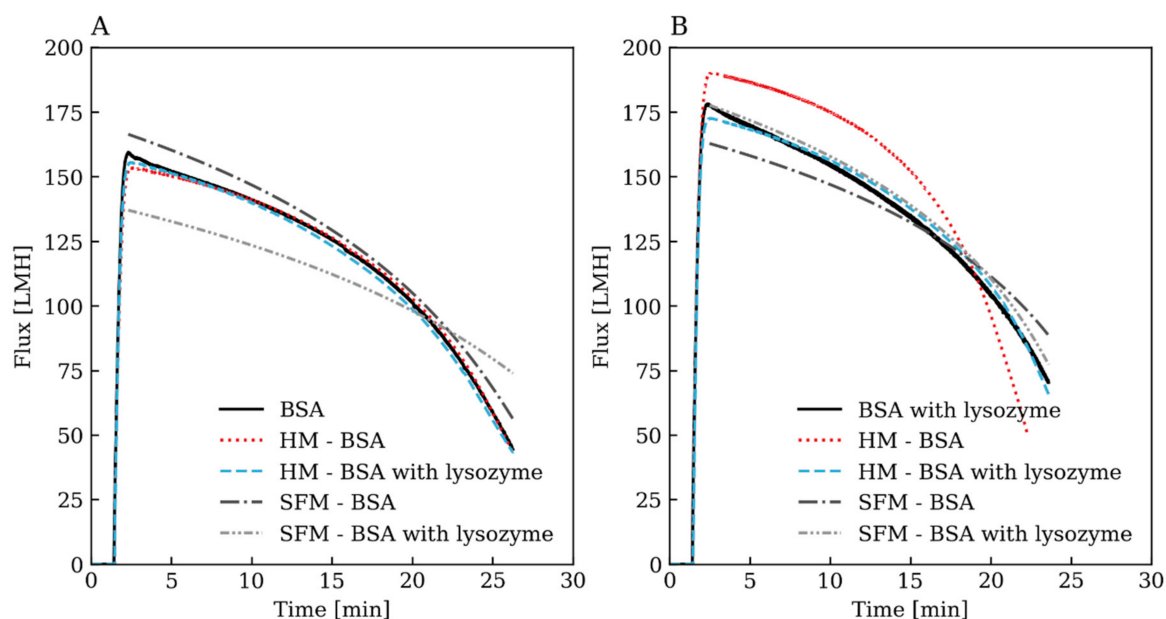


**Figure 4.** Comparing flux prediction of the test set containing (**A**) BSA (TMP 1.8 bar, CF 200 mL/min, initial $c_{B,BSA}$ 6.68 g/L) and (**B**) BSA with lysozyme (TMP 2.1 bar, CF 250 mL/min, initial $c_{B,BSA}$ 3.71 g/L, initial $c_{B,Lys}$ 0.38 g/L), with: hybrid model HM 1 trained on BSA solely (red dotted line) and the BSA and two-component training set (blue dashed line); SFM based on BSA solely (dark grey dot-dashed line) and two-component training set (light grey dot-dot-dashed line).

In contrast, the hybrid model trained with BSA (Figure 3C) and BSA with lysozyme (Figure 3A,B) training runs (Figure 4, blue dashed line) were able to predict both UF processes: BSA solely and BSA with lysozyme (Figure 4, black lines). These results showed that already low amounts of lysozyme

drastically changed the initial flux and flux evolution of the UF process and that incorporating both components in the model was essential for accurate predictions. On the contrary, SFM based on the training run containing BSA with lysozyme was also able to predict the two-component test run well (Figure 4B, light grey dot-dot-dashed line), but showed a drastic offset when predicting a test run with only BSA (Figure 4A, light grey dot-dot-dashed line). The k values from the two-component training set (Table A5) were generally lower than those calculated from solely BSA, since membrane fouling due to lysozyme was assumed. In the absence of lysozyme, however, no membrane fouling occurred and the flux for the same $c_{B,BSA}$ was higher.

In summary, the HM could predict both scenarios, since the varying concentration of lysozyme and its influence on the membrane fouling was integrated into the black box. However, the SFM only predicted one scenario well, depending on which k and $c_G$ were used. For the following two-component predictions, the SFM parameters were based on the two-component training set.

### 3.3. Comparison of Hybrid Model Performance

To further investigate both the interpolation and extrapolation capability of both HMs and the SFM model, a series of test runs were conducted under conditions that were partially not covered by the training sets. To test the hybrid models based on the two-component training set, additional test runs on BSA solutions spiked with lysozyme were performed. The two established hybrid model structures were compared for their $R_{Lys}$, flux and final $c_B$ predictions individually. $R_{Lys}$ effects the in-process $c_{B,Lys}$ prediction and subsequently the flux and final $c_{B,Lys}$. Additionally, the two hybrid models were compared to the SFM in terms of flux and $c_{B,BSA}$ prediction. $c_{B,Lys,}$ and $R_{Lys}$ could not be compared, since SFM can be applied to one-component only.

The test data consisted of nine UF runs performed at different TMP, CF, initial $c_{B,BSA}$ and $c_{B,Lys}$ conditions. Test runs 1–4 were performed within the training space. This meant that TMP and CF were within the training parameters (Figure 5A, blue area) and the initial $c_{B,BSA,}$ and $c_{B,Lys}$ was higher than the initial training concentrations (Figure 5B, blue area). The test runs 1, 2 and 9 were performed in the center of the TMP and CF training space (Figure 5A), with test run 9 containing no lysozyme. Test run 3 was performed at the outer limit of the TMP and CF training space, to investigate how the predictions of the hybrid models changed at the border. Test run 4 was performed under TMP and CF conditions not covered by the training set but within the training space, to investigate the interpolation capabilities of the model. Test runs 5, 6, 7 and 8 were performed under conditions that were partially outside the training space, such as initial $c_{B,Lys}$ (8), initial $c_{B,Lys}$ (5, 6) and CF (7), to test the extrapolation capabilities. The test run parameters are summarized in Figure 5 and Table A2.

### 3.3.1. Flux Prediction

Regarding the prediction of the flux evolution, the two hybrid models performed similarly (Figures 6A,C,E, 7A,C,E and A4A,C,E). Most test run predictions exhibited a small initial offset. At the beginning of the test experiments, the membrane was clean, while during the training set the membrane exhibited some lysozyme fouling and equilibrium of the concentration polarization layer due to the long training process time. This led to an initially underestimated flux. The offset became more pronounced when initial $c_{B,Lys}$ was higher than 0.3 g/L (test runs 2, 3, 4, 5 and 8; Figures 6C,E, 7A and A4A,C), indicating a stronger membrane fouling at this concentration. Even though all hybrid models were trained with $c_{B,Lys}$ higher than 0.3 g/L, the timely increasing membrane resistance due to fouling reached an equilibrium only after several minutes. After this point, the flux was predicted correctly. The highest initial offset was given in test run 8 (Figure A4A), which exhibited the highest initial $c_{B,Lys}$ and therefore more fouling. Test run 7 (Figure 7E) was performed at CF 350 mL/min, which was outside the training space. Both hybrid models predicted the flux of test run 7 (Figure 7E) well, indicating that the models were not necessarily limited by the training space and showed good extrapolation capabilities of the input parameter CF. Test runs 4, 5 and 6 (Figures 6C,E and 7C) exhibited TMPs and

CFs within the training space parameters and all predicted well. The good flux predictions of these test runs showed the excellent interpolation capabilities of the ANN-aided hybrid models.
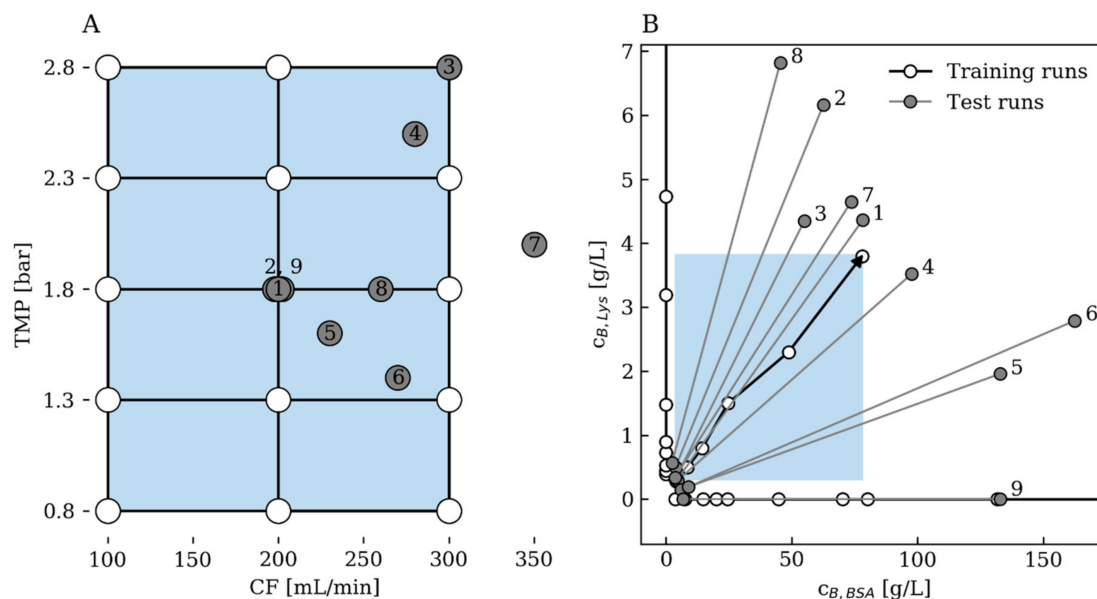


**Figure 5.** Schematic depiction of the training space (blue area) for: (**A**) TMP and CF of training runs (white dots) and test runs (grey dots); (**B**) initial to final $c_{B,BSA}$ and $c_{B,Lys}$ of the test runs (grey dots with grey solid lines) and the covered concentration range of the three training runs (white dots with black solid lines).

The SFM predicted the initial flux and flux evolution inside the training space well (test runs 1, 3 and 4; Figures 6A,C, and A4C). However, for the test runs outside the training space, higher errors were exhibited (test runs 5, 6 and 8; Figures 6E, 7C and A4A). Outside the training space, k and $c_G$ were extrapolated from the training data, which potentially increased flux prediction uncertainty. Furthermore, high lysozyme concentrations also led to higher errors due to stronger fouling over time and not being able to incorporate the second component in the SFM. Here, the SFM underestimated the initial flux drastically (test runs 2 and 8; Figures 7A and A4A). For test run 9 (Figure A4E)—only BSA, no lysozyme—the SFM with k and $c_G$ were exceptionally based on BSA training data (Figure 3C) to allow fair comparison. In this case, the SFM yielded good initial flux predictions, but deviations at the end of the process, while HM 1 and 2 both showed excellent flux prediction over the entire process. On average, the flux prediction error of SFM was 6% NRMSE, while the error of the two hybrid models was 4.1% and 3.9% NRMSE (Figure 8A).

3.3.2. Rejection Factor Prediction for Lysozyme

The rejection factor for lysozyme $R_{Lys}$ increased throughout the UF run, from around 0.6 to almost 1.0, as shown in Figures 6, 7 and A4. The pores became increasingly blocked throughout the UF process, most probably because lysozyme was absorbed in their inner wall, increasing the rejection factor. Results showed that there was no consistent correlation between the TMP and $R_{Lys}$, or CF and $R_{Lys}$ (Figure A5). Therefore, the influence of TMP and CF on $c_{P,Lys}$ was neglected when creating the calibration between UV absorbance and lysozyme permeate concentration. The rejection factor of BSA was 1 for all experiments. The model errors are given in Figure 8B.
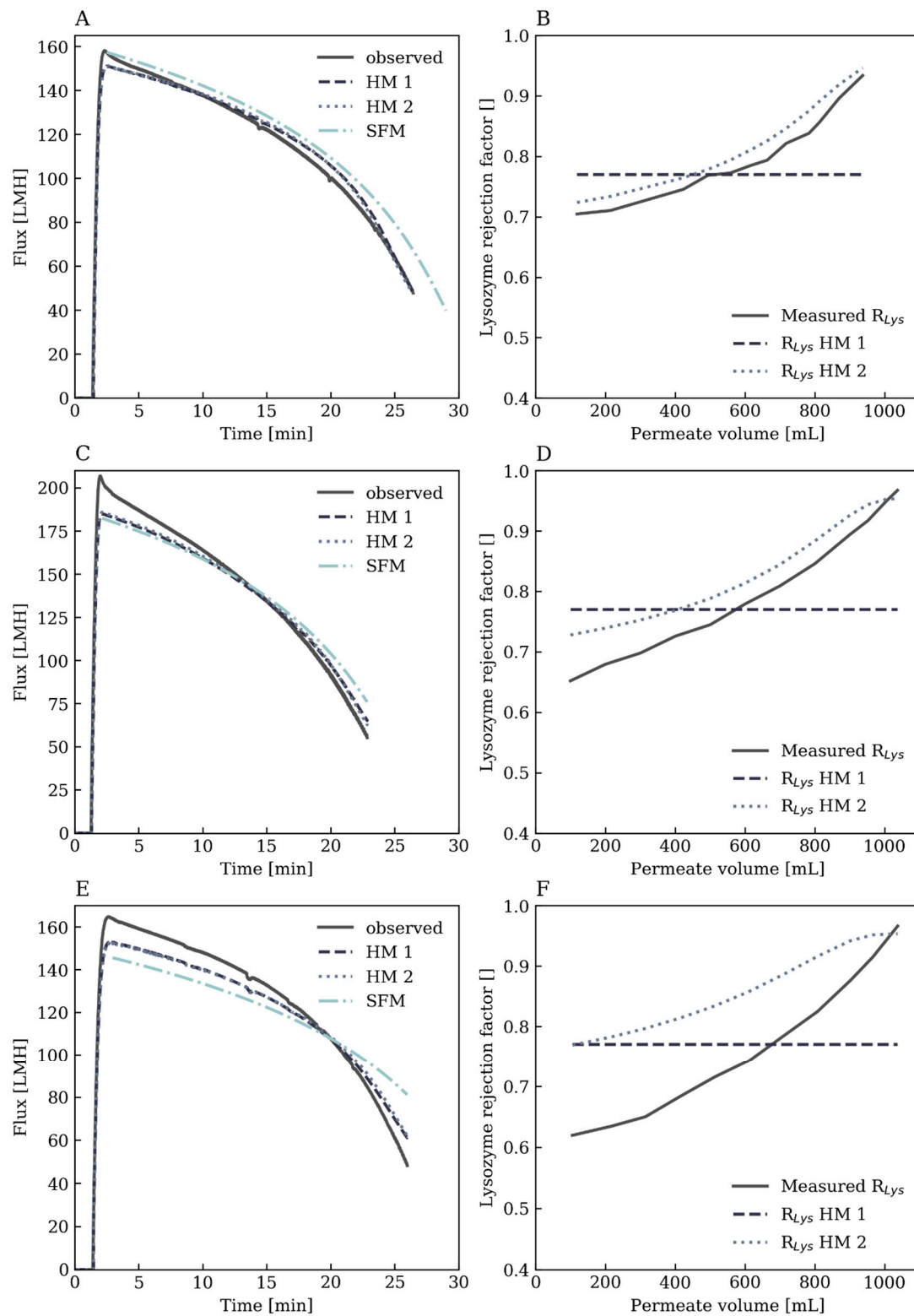
**Figure 6.** Comparison of observed and predicted flux and $R_{Lys}$. (**A**) The flux over time of test run 1, (**B**) $R_{Lys}$ over permeate volume of test run 1, (**C**) the flux over time of test run 4, (**D**) $R_{Lys}$ over permeate volume of test run 4, (**E**) the flux over time of test run 5, (**F**) $R_{Lys}$ over permeate volume of test run 5.
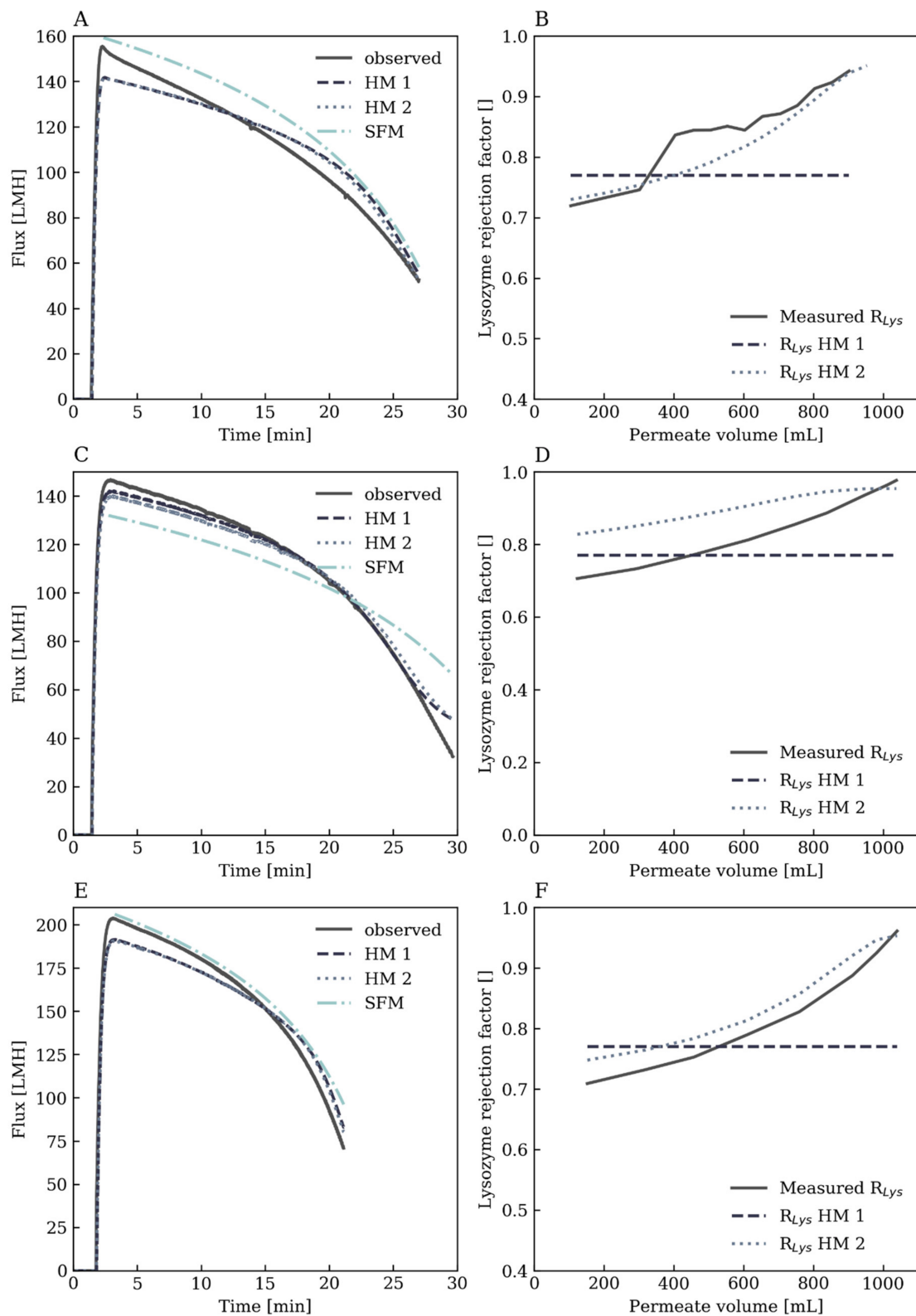
**Figure 7.** Comparison of observed and predicted flux and $R_{Lys}$. (**A**) The flux over time of test run 2, (**B**) $R_{Lys}$ over permeate volume of test run 2, (**C**) the flux over time of test run 6, (**D**) $R_{Lys}$ over permeate volume of test run 6, (**E**) the flux over time of test run 7, (**F**) $R_{Lys}$ over permeate volume of test run 7.
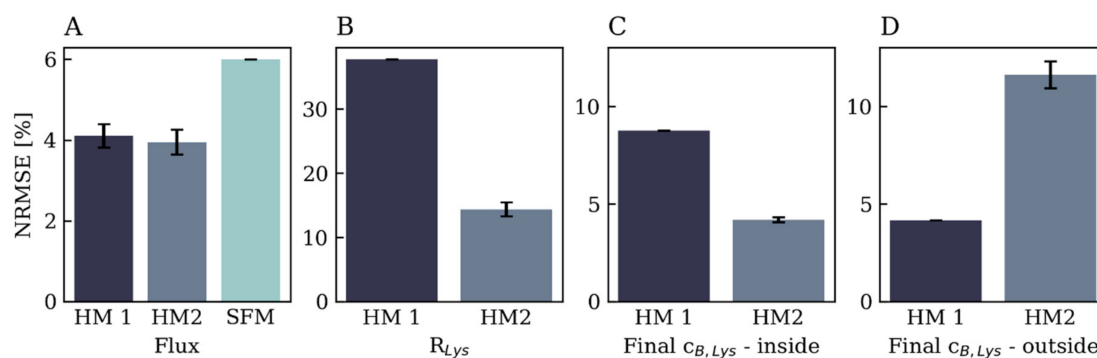
**Figure 8.** Summary of the prediction errors of HM 1, HM 2 and SFM in terms of (**A**) flux, (**B**) $R_{Lys}$, (**C**) final $c_{B,Lys}$ for test runs 1 to 4 with all parameters—TMP, CF, initial $c_{B,BSA}$ and $c_{B,Lys}$—inside the training space and (**D**) final $c_{B,Lys}$ for test runs 5 to 8 performed partly outside the training space.

Hybrid Model 1: Constant Lysozyme Rejection Factor

In HM 1 the rejection factor for lysozyme $R_{Lys}$ was assumed to be constant over time for all test runs, where lysozyme was present and therefore exhibited the largest $R_{Lys}$ error (38% NRMSE) compared to HM 2 (see Figure 8A). All test runs (Figures 6, 7 and A4) show that HM 1 overestimated $R_{Lys}$ at the beginning of all UF runs and underestimated it at the end. The average $R_{Lys}$ based on training data fitted all independently generated test data very well but lacked the ability to adjust to the increasing $R_{Lys}$.

Hybrid Model 2: Dynamic Lysozyme Rejection Factor

In contrast to keeping the rejection factor constant, as in HM 1, a second black box was introduced in HM 2 to predict $R_{Lys}$ dynamically. This prediction was independent of the flux prediction but was based on the same four input parameters, namely TMP, CF, initial $c_{B,BSA}$ and initial $c_{B,Lys}$. The NRMSE of the newly introduced black box was evaluated by comparing the observed $R_{Lys}$ values to the predicted $R_{Lys}$. Since the correlation of $R_{Lys}$ and $V_P$ is quite simple, an ANN with one hidden node was used for $R_{Lys}$ prediction (Figure A1C). For comparison, a multiple linear regression (MLR) model was also tested as an alternative black box, resulting in a less complex hybrid model that required less computation time and facilitated easier interpretability. However, the ANN with one node was chosen instead of the MLR, because of the lower prediction error regarding $R_{Lys}$ and final $c_{B,Lys}$ (see Table A4).

HM 2, with an average $R_{Lys}$ NRMSE of 14%, performed better than HM 1. The improvement was achieved as HM 2 considered the increasing $R_{Lys}$ over the process, which subsequently strongly influenced the final $c_{B,Lys}$ prediction (Section 3.3.3). In test runs 5 and 6 (Figures 6F and 7D) the prediction from HM 2 overestimated $R_{Lys}$. These test runs exhibited a low TMP and high $c_{B,BSA}$. The hybrid model assumed that the CP layer of BSA and fouling due to lysozyme were at an equilibrium, at which the lysozyme transmission was lower than in the test runs, where the CP layer was still building up. Low TMP additionally prolonged the time to reach flux steady state. $R_{Lys}$ of the other test runs 1, 2, 4, 7 and 8 (Figures 6B,D, 7B,F and A4B) were predicted accurately with HM 2.

Even though HM 2 performed better than HM 1 in $R_{Lys}$ prediction, the flux predictions were almost identical (NRMSE 3.9 and 4.1 %). This indicated that they were not affected by small variations or changes in $c_{B,Lys}$.

3.3.3. Endpoint Bulk Concentration

Since $R_{BSA}$ was 1, all models—HM and SFM—predicted the same $c_{B,BSA}$ at the final retentate volume, with an average error of 4.2% (Table A3). BSA did not show membrane fouling and was quantitatively recovered at the end of the process. The predictions of the final $c_{B,Lys}$ varied because

of the different $R_{Lys}$ predictions. The discussion for $c_{B,Lys}$ prediction was divided into the test runs performed strictly inside and outside the training space, since the hybrid models performed differently.

Within the training space—test runs 1–4—HM 1 and HM 2 performed in accordance with the $R_{Lys}$ predictions (Figure 8C). HM 1 exhibited the highest error of 9% since $R_{Lys}$ was not adjusted over the processing time. HM 2 recalculated $R_{Lys}$ with every iteration; its $c_{B,Lys}$ predictions were in good accordance with the measured concentrations, with an NRMSE of 4% and superior to HM 1. Similarly to the $R_{Lys}$, the accuracy of the final $c_{B,Lys}$ prediction benefited from two separately trained black boxes.

In cases where at least one input parameter was outside of the training space—test runs 5–8—HM 1 performed best with an average NRMSE of 4% (Figure 8D). In comparison, HM 2 yielded worse final $c_{B,Lys}$ predictions, exhibiting a three-fold increase in NRMSE (12%). Even though $R_{Lys}$ was updated in HM 2, it was overestimated throughout most of the test runs, leading to higher $c_{B,Lys}$ prediction and a cumulated NRMSE that increased with the duration of the process. In contrast, using HM 1 the initial $R_{Lys}$ over-prediction and under-prediction balanced out and yielded acceptable final $c_{B,Lys}$ predictions.

In summary, the more complex HM 2 showed superior performance within the trained space, which is the case for most modeling applications. It can predict the flux, $R_{Lys}$ and therefore the concentration, of both components at any time point of the process. For predictions outside the trained space, the simpler and more robust HM 1 performed better, giving accurate predictions on flux and the fully retained main component BSA. It can offer valuable insights when exploring parameter ranges if the desired optimal process conditions are not met in the trained space, before it is expanded and used to retrain new hybrid models.

## 4. Conclusions

UF modeling increases process understanding which is key for predicting process performance. The interactions of various components means that mechanistic modeling approaches for multi-component solutions might become very complex and require many experiments.

We developed and compared hybrid models to predict flux, rejection behavior and concentrations for UF of two-component solutions. The models were trained on training experiments that were generated in less than eight hours and tested on independently performed UF runs with varying product and impurity concentrations, TMPs and CFs. We showed that the hybrid model HM 2, with a dynamic impurity rejection factor containing two black boxes, exhibited the best predictions for impurity rejection behavior and final concentration within the trained parameter space and had excellent interpolation properties. The simpler HM 1 yielded stable predictions beyond the trained space, rendering it a valuable tool for extrapolation. Both hybrid models performed similarly well in predicting flux and mimicked product concentration. The SFM with mechanistic parameters exhibited higher flux prediction errors than both hybrid models and could not predict the lysozyme rejection factor and final concentration, since it can only assume a one-component system. Our results show that it is crucial to quantify and incorporate all components, including the impurities, to gain accurate and reliable process models. These variations can be included more easily in the hybrid model approach than in mechanistic models such as SFM, with low experimental effort and no mechanistic parameter adaption required.

A limitation of the presented models is the time-dependent fouling of the mimicked impurity at high initial concentrations. However, at the expected concentration ranges, e.g., after the chromatography capture step, the effect can be neglected.

The proposed hybrid model structure can be used not only for the reliable prediction of final product concentrations, but also of the concentration of various quantifiable classes of impurities. Since impurities are a critical quality attribute (CQA) in many manufacturing bioprocesses, time-resolved concentration predictions help to better understand the process's outcome upfront. Furthermore, by taking adequate measures a potential batch rejection due to high impurity concentration can be avoided. The product and impurities can be measured with online sensors or correlated with

offline analytics using soft sensors. In combination, with closed-loop process controllers, these hybrid models are a valuable tool for increased process understanding and advanced process control.

## Symbols and Abbreviations

| | |
|---|---|
| ANN | artificial neural network |
| BSA | bovine serum albumin |
| CF | cross-flow velocity |
| CP | concentration polarization |
| HM | hybrid model |
| MWCO | molecular weight cutoff |
| NRMSE | normalized root-mean-squared error |
| SEC | size exclusion chromatography |
| SFM | stagnant film model |
| TMP | transmembrane pressure |
| UF | ultrafiltration |
| | |
| $A$ | membrane area [m$^2$] |
| $c_B$ | bulk concentration [g/L] |
| $c_G$ | gel layer concentration [g/L] |
| $c_P$ | permeate concentration [g/L] |
| $c_R$ | retentate concentration [g/L] |
| $dt$ | time increment [s] |
| $J$ | permeate flux [LMH] or [m/s] |
| $k$ | mass transfer coefficient [LMH] |
| $R_{Lys}$ | lysozyme retention coefficient [-] |
| $V_B$ | bulk/reservoir volume [mL] |
| $V_P$ | permeate volume [mL] |

## Appendix A

### *Appendix A.1. Neural Network Model Optimization*

To choose the best-suited ANN structure, varying numbers of hidden nodes were tested. Each ANN was trained on the combined training set and validated. Figure A1A gives an overview of the optimal ANN structure including the inputs TMP, CF, $c_{B,BSA}$, $c_{B,Lys}$, and the output permeate flux (in HM 2 a second ANN with $R_{Lys}$ as output was added) with four hidden neurons. The input and output parameters were scaled between 0 and 1 before optimizing the ANN. This step is necessary to have the parameters on the same scale rendering them comparable. Each node in the hidden and output layer in Figure A1A forms a linear equation. As an example, the first hidden node $x_{21}$ is the sum of each multiplication of an input (TMP, CF, $c_{B,BSA}$, $c_{B,Lys}$) and the corresponding weight ($w^1_{11}$, $w^1_{21}$, $w^1_{31}$, $w^1_{41}$) multiplied with the bias ($b_1$) of the entire hidden layer.

$$x_{21} = b_1\big(w^1_{11}TMP_{scaled} + w^1_{21}CF_{scaled} + w^1_{31}c_{B,BSA,scaled} + w^1_{41}c_{B,\,Lys,scaled}\big)$$

To determine the values for the weights and biases that result in the desired prediction the model is optimized in multiple epochs. As a first step, the weights and biases are randomly chosen and the first prediction with inputs from a given training set is performed. Since the weights and biases are not optimized the flux prediction will be of poor quality and the prediction error will be high. Using the desired output from the training set, the ANN is calculated backward which results in inputs parameters that fit the prediction. The error between the real and the backward calculated inputs is estimated and used to update the according to weights and biases.