# Advanced score system and automated search strategies for parameter estimation in mechanistic chromatography modeling

William Heymann [a,b], Juliane Glaser [c], Fabrice Schlegel [d], Will Johnson [d], Pablo Rolandi [d], Eric von Lieres [a,*]

[a] *Institute of Geo- and Biosciences 1 (IBG-1), Forschungszentrum Jülich, Wilhelm-Johnen-Str. 1, 52428 Jülich, Germany*
[b] *RWTH Aachen University, 52062 Aachen, Germany*
[c] *Digital Integration and Predictive Technologies (DIPT), Amgen Research Munich, Staffelseestr. 2, 81477 München, Germany*
[d] *Digital Integration and Predictive Technologies (DIPT), Amgen, 360 Binney St, Cambridge, MA 02142, United States*

## ARTICLE INFO

## ABSTRACT

Least squares estimation of unknown parameters from measurement data is a well-established standard method in chromatography modeling but can suffer from critical disadvantages. The description of real-world systems is generally prone to unaccounted mechanisms, such as dispersion in external holdup volumes, and systematic measurement errors, such as caused by pump delays. In this scenario, matching the shape between simulated and measured chromatograms has been found to be more important than the exact peak positions. We have therefore developed a new score system that separately accounts for the shape, position and height of individual peaks. A genetic algorithm is used for optimizing these multiple objectives. Even for non-conflicting objectives, this approach shows superior convergence in comparison to single-objective gradient search, while conflicting objectives indicate incomplete models or inconsistent data. In the latter case, Pareto optima provide important information for understanding the system and improving experiments. The proposed method is demonstrated with synthetic and experimental case studies of increasing complexity. All software is freely available as open source code (https://github.com/modsim/CADET-Match).

## 1. Introduction

Chromatography models can aid the rational design and robust operation of separation processes. However, before a model can be used it must first be calibrated and validated. Some of the parameters involved in chromatography models are not directly measurable, and empirical correlations exist for some but not all of them. This naturally leads to a procedure where the parameters are estimated based on chromatogram data. In this paper parameter estimation is broken up into goals and search strategies. Different combinations of proposed goals and search strategies are introduced and tested on synthetic and industrial data sets. For industrial application, the entire estimation procedure must be mostly automated and robust, which is the major focus of this contribution. We introduce a procedure that is designed to get good-enough answers in reasonable time while dealing with systematic errors and random noise in the data. The presented procedure is general purpose so that it can solve a wide range of problems and is flexible enough to adapt to future needs. It is implemented in the open source software CADET-Match that is freely available on GitHub (https://github.com/modsim/CADET-Match).

The goal of the parameter estimation procedure can be composed of one or more metrics, i.e. specific features that quantify how good or bad a simulated chromatogram matches the corresponding experimental data. Note that the terms goal, metric and objective are used with distinct meaning that will be formally defined in the respective sections of this paper. Each metric mathematically formalizes the goodness of an objective. Typical metrics are based on the sum of squared differences between simulation and measurement data or the shape similarity between these curves independent of a time offset. Optimization problems are usually formulated such as to minimize the metrics. Hence, a suitable metric must have a low value when the objective is good, a high value when the objective is bad, and provide a path from bad to good where increasingly better objectives are indicated by monotonically decreasing values of the metric. These statements might seem simple but can be problematic in practice, as will be

---

* Corresponding author.
  *E-mail address:* e.von.lieres@fz-juelich.de (E. von Lieres).

illustrated in Section 7. Ideally the path between high and low metric values is smooth. Some metrics are easier to compute gradients for and this can influence the choice of an adequate search strategy.

Search strategies are all about tradeoffs. From the no free lunch theorem [1] "any algorithm, any elevated performance over one class of problems is offset by performance over another class." This means that there is no general-purpose best search strategy for all kinds of problems and a proper choice depends on the nature of the goal as well as the starting points. Search strategies that are based on gradient descent can be very efficient for some problems but poorly suited for others. Search strategies that are based on evolution do generally work on a wider class of problems but converge typically much slower than more specialized strategies such as gradient descent. Evolutionary search strategies are black-box algorithms that have very few requirements to function compared to other search strategies. They are often a good first choice when little is known about the goal before more efficient strategies are applied. This paper focuses on gradient descent and evolutionary algorithms.

## 2. Literature review

Parameter estimation is not an end to itself. A calibrated model only matters to the extent it can add value to its defined purpose. Once a model has been calibrated it can be used to design an optimal separation process [2,3]. A model can also be used to see the tradeoffs between productivity and yield to design processes optimized for previously defined goals [4]. Additionally, models can be used to understand the underlying physics and the impact of effects, such as surface diffusion [5].

Early methods for calibrating models used algebraic approximations and frontal experiments in 1983 for the mass transfer and binding [6,7]. As more complex binding models like steric mass-action (SMA) [8] where published the calibration approaches became more complex and a combination of breakthrough and pulse experiments were integrated into the calibration process [9]. Early methods still tended to assume systems are always in equilibrium for computational reasons but did start to integrate non-pore penetrating pulses and pore-penetrating but non-binding pulses to find the interstitial and particle porosities [10]. However, due to the coarse nature of the models used, rate models were mostly insensitive to diffusion parameters, which can be explained by high numerical diffusion.

The next stage of calibrating models used rate-based models with decreased numerical diffusion more frequently and mostly manual fitting with some computer assistance. Early rate-based models for the mass transfer were implemented with equilibrium binding models and while most parameters where still manually calculated some numerical fitting started to be used [11]. As systems became more complicated the fitting process started to become more complex with it, such as handling multi-component competitive binding using Lagrange multipliers [12]. Other methods continued to try to efficiently solve problems given the computer resources of the time, such as using the perturbation method [13].

As computing power further increased fewer assumptions had to be made and it was possible to fit multiple experiments simultaneously with overlapping components [14]. This marked a major change because it showed that directly fitting to experimental data gave superior results to previous methods using breakthrough curves. Direct fitting allowed removing assumptions, for example that the binding capacity for each component was the same. It became more common to use a combination of breakthrough and pulse experiments with rate-based binding and transport models [3,15]. While it appears methods get more effective and efficient in terms of fitting models to data, the fact remains that they require good starting points for the gradient descent algorithms used in the fitting and the models are still quite coarsely discretized due to limited computing resources available.

As more computing power became available it drove the development of other methods to find suitable model parameters. A data-based approach using a support vector machine (SVM) and quantitative-structure-property-relationship (QSPR) is used to predict binding properties [16]. Other methods use batch experiments to directly measure film diffusion and binding parameters instead of having to rely on parameter estimation [17].

Computing power continued to rapidly advance to the point that finer grained models can be used for single column systems. One of the drawbacks of gradient descent is that accurate gradients are needed which can be computationally expensive to approximate or the entire model needs to have analytical derivatives. An iterative approach was used with finite difference using complex numbers. This approach requires modifying a model such that complex numbers can be used [18]. As solvers get more accurate and computers get more efficient, fitting diffusion parameters is used to better understand diffusion inside beads with fitted pore diffusion in the range we measure today [19]. Even as solvers get more accurate for single column systems, for multi-column systems with multiple components older methods of model calibration using algebraic assumptions again are used [20].

With increasing compute power, activities in development of the parameter estimation methods increased. While inverse fitting is difficult and gets more complex as model complexity increases, the results are superior to older methods that use empirical correlations and are also shown to be insensitive to UV noise but extremely sensitive to time offsets [21]. Often sum of squared errors are used as an estimation objective. The drawback of sum of squared errors is that it is extremely sensitive to time offsets of peaks and not as sensitive to the shape of the chromatograms. One way to deal with this is to use a weighted sum of mean, standard deviation, and skewness of the chromatograms which provides more sensitivity to the overall shape and more tolerance of non-overlapping starting points and is an important stepping point towards the methods presented in this paper [22].

Variable transforms for parameter estimation are missing from earlier papers. Parameters of different scales slow optimizer performance. A major step forward was combining variable transforms, fitting simultaneous experiments using a combination of breakthrough and pulses, modeling the external column effects using a combination of PFR and CSTR, and using gradient descent [23]. However, even with all these advances fitting of many models remains difficult [24,25].

Finally, we move towards more advanced and robust but computationally expensive methods. With more computing power available using a global optimization algorithm like a genetic algorithm and using gradients for local refinement further improves parameter estimation by removing the need to have a good initial starting point [26,27]. An alternative way to solve parameter estimation problems that is still in its infancy is using neural networks [28]. The basic idea is to use a model to sample the space and then train a neural network using the model output chromatogram as the input and the output of the network as the parameters. An unknown chromatogram based on the same model can then be provided as input for the neural network and the parameters directly obtained as output without any estimation. Depending on the data used to train the network and the quality of the network this approach has the potential to dramatically shorten estimation times or provide good starting points for refinement.

A common thread that weaves through all the history here is that creating calibrated models is complex but necessary. Many of the parameters are correlated and not intuitively disentangled

from each other and from extra column effects and necessitated not only multiple experiments but also a stagewise estimation process [29]. This is what sets the stage for a robust and automated parameter estimation process.

## 3. Chromatography modelling

This paper addresses packed bed liquid chromatography at preparative scale. Such systems can be described in-silico by combining different models of the governing transport and binding processes. The general rate model (GRM), Eqs. (1), (2), with suitable boundary conditions, Eqs. (3)–(6), and non-equilibrium SMA binding, Eqs. (7)–(10), is a common choice. The GRM describes convection and dispersion in the interstitial column volume, diffusion in the porous particles, and binding to the inner surfaces of these particles.

$$\frac{\partial c_i^b}{\partial t} = -u_c \frac{\partial c_i^b}{\partial z} + D_c \frac{\partial^2 c_i^b}{\partial z^2} - \frac{1 - \varepsilon_c}{\varepsilon_c} \frac{3}{r_p} k_f \left( c_i^b - c_i^p (r = r_p) \right) \quad (1)$$

$$\frac{\partial c_i^p}{\partial t} = D_p \left( \frac{\partial^2 c_i^p}{\partial r^2} + \frac{2}{r} \frac{\partial c_i^p}{\partial r} \right) - \frac{1 - \varepsilon_p}{\varepsilon_p} \frac{\partial c_i^s}{\partial t} \quad (2)$$

$$D_c \frac{\partial c_i^b}{\partial z} (z = 0) = u_c \left( c_i^b (z = 0) - c_i^{in} \right) \quad (3)$$

$$\frac{\partial c_i^b}{\partial z} (z = L_c) = 0 \quad (4)$$

$$\frac{\partial c_i^p}{\partial r} (r = r_p) = \frac{1}{\varepsilon_p D_p} k_f \left( c_i^b - c_i^p (r = r_p) \right) \quad (5)$$

$$\frac{\partial c_i^p}{\partial r} (r = 0) = 0 \quad (6)$$

$$\frac{dc_i^s}{dt} = \tilde{k}_{a,i} c_i^p \left( \frac{\bar{c}_0^s}{c_{r,i}^s} \right)^{\nu_i} - \tilde{k}_{d,i} c_i^s \left( \frac{c_0^p}{c_{r,i}^p} \right)^{\nu_i} \quad (7)$$

$$k_{a,i} = \tilde{k}_{a,i} \left( c_{r,i}^s \right)^{-\nu_i} \quad (8)$$

$$k_{d,i} = \tilde{k}_{d,i} \left( c_{r,i}^p \right)^{-\nu_i} \quad (9)$$

$$\bar{c}_0^s = \Lambda - \sum_{i=1}^{N_c} (\nu_i + \sigma_i) c_i^s \quad (10)$$

The concentrations $c_i^\kappa$ with $\kappa \in \{b, p, s\}$ refer to the interstitial column or bulk liquid ($b$), particle pores ($p$) and stationary phase ($s$) of the system. Chromatography models depend on many parameters, several of which need to be estimated by fitting simulated chromatograms to experimental data. In Eqs. (1)–(6) the column porosity, $\varepsilon_c$, particle porosity, $\varepsilon_p$, axial dispersion, $D_c$, film diffusion, $k_f$, pore diffusion, $D_p$, adsorption rate, $k_{a,i}$, desorption rate, $k_{d,i}$, shielding coefficient, $\sigma_i$ and characteristic charge, $\nu_i$, of component $i = 1, \ldots, N_c$ are typically estimated from measured chromatograms. Film and pore diffusion can generally differ between components but are assumed to be identical for the specific molecules used in this study. Other parameters such as the column length, $L_c$, particle radius, $r_p$, interstitial velocity, $u_c$, and ionic capacity, $\Lambda$, can be controlled or measured in advance of the simulation. The same is true for the initial concentrations, $c_i^\kappa (t = 0)$, and inlet concentration profiles, $c_i^{in}$.

The SMA model, as originally introduced by Brooks and Cramer [8], becomes numerically unstable for molecules with high characteristic charge, $\nu$, such as monoclonal antibodies on high capacity resins. This critically important problem is effectively avoided by

scaling the rate constants by the $\nu^{th}$ power of reference concentrations $c_r^p$ and $c_r^s$ [30]. Recommended values are the highest salt concentration in the feed during elution for $c_r^p$, and the ionic binding capacity of the resin for $c_r^s$. This ensures both terms raised to the $\nu^{th}$ power in Eq. (7) range between 0 and 1, while without scaling they can become extremely large and consequently cause accuracy loss and instability of the numerical solver. Moreover, the units of the scaled rate constants are independent of the characteristic charge, which is beneficial for the physical interpretation of these values and for the performance of search algorithms during parameter estimation.

The tubing is modeled using a dispersive plug flow reactor (DPFR), Eq. (11). The tubing model is solved with the same inlet and outlet boundary conditions, as the column model, Eqs. (3) and (4), with the tubing length, $L_t$, in place of the column length, $L_c$.

$$\frac{\partial c_i^t}{\partial t} = -u_t \frac{\partial c_i^t}{\partial z} + D_t \frac{\partial^2 c_i^t}{\partial z^2} \quad (11)$$

In the tubing, $c_i^t$ denotes concentration, $u_t$ velocity and $D_t$ axial dispersion. A continuously stirred tank reactor (CSTR) is used to model mixers and mixing effects. In addition to the model equations presented here, CADET-Match works with any combination of transport and binding models that is covered by the CADET solver, which is an independent and continuously extended open-source project (https://github.com/modsim/CADET).

## 4. Parameter estimation

Many parameters in chromatography models are highly correlated in the sense that small changes in different parameters impact on the simulated chromatogram in similar ways. This could technically be neglected when accurate model predictions are required only for the exact same column dimensions and operating conditions that were used for parameter estimation. However, this is not the case when the calibrated model is to be applied for guiding rational process design and scale-up.

In these applications it is crucially important that the estimated parameter values correctly represent and accurately quantify the impact of the respective underlying physical mechanism described by the model. Parameter correlations can be avoided by a staged estimation procedure that isolates these parameters using specific experiments whose design and order depends on the mathematical structure of the model equations [11]. Typically, four types of experiments are required with 1) detached column to determine the band broadening effect of extra-column volumes, 2) a non-binding tracer that does not penetrate the particle pores to determine column porosity and axial dispersion, 3) a non-binding but pore-penetrating tracer to determine particle porosity, film diffusion and pore diffusion, and 4) the target molecules to determine the parameters of the binding model.

The first stage is increasingly recognized to be crucial for determining unbiased values of all other model parameters that can be transferred across operating conditions, system configurations and scales. Neglecting extra column effects or even small errors in accounting for them can have a large impact on the binding parameters that are estimated at a later stage [31]. Extra column effects have been previously accounted for by shifting the time scale [24]. More comprehensive models comprise a series or network of DPFR and CSTR [23]. The respective model parameters are isolated by removing the column from the simulated system.

In the second stage, model parameters for characterizing the packed bed are isolated by setting the film diffusion coefficient to zero, which effectively eliminates Eqs. (2)–(6) from the system. Dextran is normally used as non-binding and non-pore penetrating tracer. However, this tracer often behaves non-ideally, which

causes specific challenges in the parameter estimation procedure as will be discussed in Section 7.3.

In the third stage, model parameters for characterizing the porous particles are isolated by setting the adsorption constant to zero, which effectively eliminates Eqs. (3)–(6) from the system. Ideally, the target molecule is used as pore-penetrating tracer under non-binding conditions, such as high salt in ion-exchange chromatography. Smaller tracers, such as acetone, are prone to overestimate the film and pore diffusion coefficients. Occasionally, the film and pore diffusion coefficients remain too correlated to be estimated independent of each other. In these cases, the pore diffusion can be determined from some other type of experiment, such as using confocal laser scanning microscopy [19]. Alternatively, a simpler transport model can be used, such as the lumped rate model with pores.

In the fourth stage, the parameters of the binding model are estimated. Due to inherent non-linearity of the more complex binding models that are required for describing preparative chromatography, this is usually by far the most complicated and time-consuming stage, in particular for complex multi-component systems with competitive binding. Base-line separation is typically not achieved in experimental data available for calibrating such models. Hence, the binding model parameters of multiple chemical components cannot be completely isolated from each other. However, correlations between the binding parameters can be reduced by estimating them from a set of several chromatograms that are measured at specifically designed operating conditions. For large molecules, such as monoclonal antibodies, this typically includes a breakthrough and two or three gradient elution experiments with varying slopes. The design of such experiments can be optimized by evaluating different estimation strategies on synthetic data that are generated using an initial guess of plausible values for the sought model parameters. Similar components, such as charge variants or high/ low molecular weight impurities, are often lumped in groups to reduce model complexity and the number of estimated parameters.

Measurement data are generally prone to random and systematic errors that are in turn propagated to the estimated parameters. In the proposed stagewise procedure, parameters estimated in one stage are fixed in the next and, consequently, parameter errors are carried over to subsequent stages. This can be avoided in a Bayesian approach where the posterior parameter distribution of one stage is used as prior information in the next stage, which will be subject of a separate study. Without that, it is particularly important to be accurate in the early stages.

## 5. Materials and methods

Chromatographic cation exchange runs were conducted on Äkta Avant (GE Healthcare). Blue Dextran 2000 (GE Healthcare) was used for non-pore penetrating pulse experiments. A volume of 1 mL Dextran solution with a concentration of 0.002 mM was injected at a flow rate of 5 mL/min with and without column attached. These experiments were carried out in duplicates. Dextran chromatogram data was measured at UV 280 nm.

Pore-penetrating as well as load, wash and elution steps were performed using a Fractogel $SO_3^-$ (EMD Millipore) resin in a packed bed column with inner diameter 16 cm and length 25 cm at a flow rate of 5 mL/min. The column was pre-equilibrated for 3 CV with 200 mM sodium acetate buffer containing 1 M NaCl at pH 5.0. The column was equilibrated for 3 CV with 100 mM sodium acetate buffer at pH 5.

Pulse injections under non-binding conditions were run using previously purified monomeric antibody as tracer. The antibody was prepared in a 100 mM sodium acetate buffer with 500 mM NaCl at pH 5 to prevent adsorption.

For elution experiments, the column was loaded with 152 mL or 540 mL, respectively, of Filtered Virus Inactivated Pool (FVIP) of monoclonal antibody product. The material was obtained from a previous capturing step. The antibody material was produced by CHO cell culture. The FVIP material was conditioned with arginine, targeting 50 mM arginine concentration in the FVIP material. A wash step was conducted with 100 mM sodium acetate buffer at pH 5. The elution step was conducted with a linear 1 mM/CV gradient between two buffers with 100 mM sodium acetate and 100 mM sodium acetate plus 1 mM NaCl, respectively, at pH 5.0. After elution the column was cleaned using 1 M NaOH solution. Chromatogram data of the load, wash and elution experiments were measured at UV 300 nm and UV 280 nm. For the 152 mL load volume run, fractionation samples were taken at 8 mL fraction volume starting at 0.1 AU of UV 280 signal.
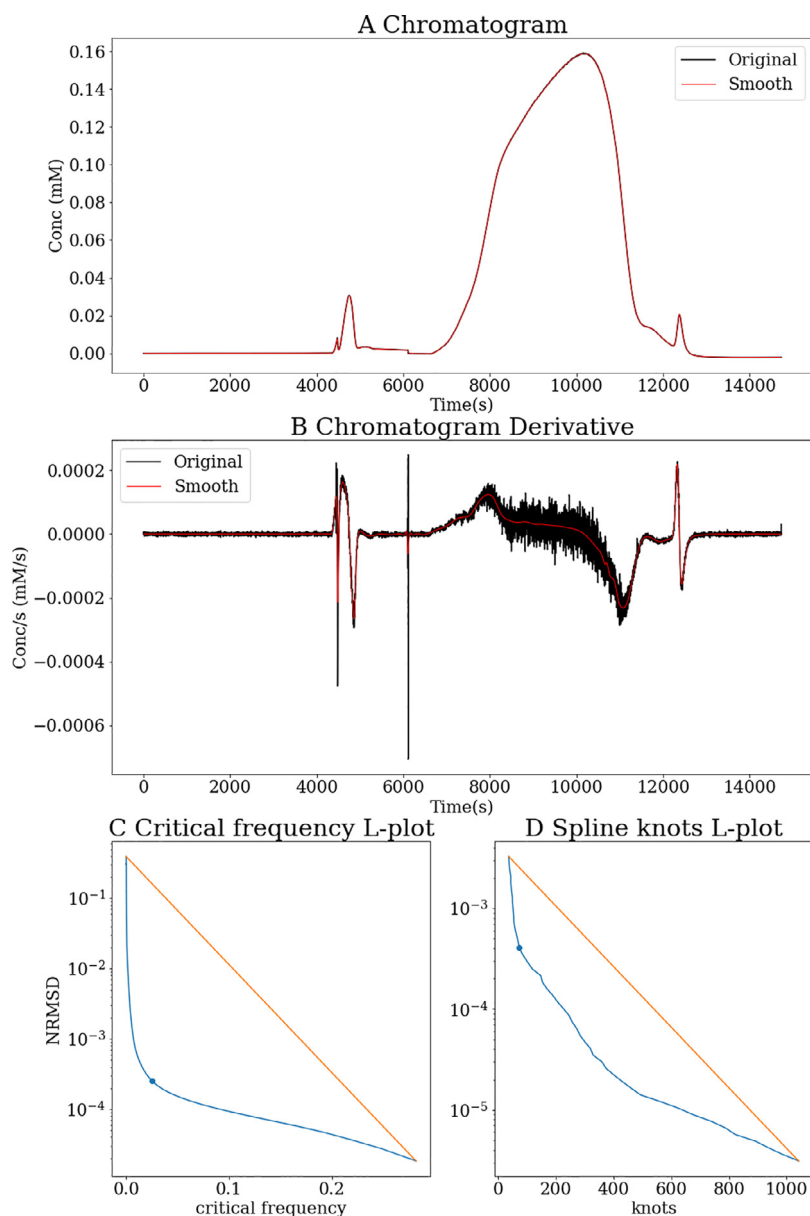
## 6. Data smoothing

Experimental data is often noisy, and the goals introduced in the next section are highly sensitive not only to the shape of the chromatogram but also to such noise. In addition, most of these goals and some search strategies also require a smooth first derivative. Hence, the data needs to be smoothed to reduce the noise. For routine application in industrial workflows, the smoothing needs to be automatic and work robustly, i.e. not attenuate relevant features, without human involvement. This is complicated by the situation that an experiment with detached column typically results in a signal length in the order of seconds, a non-binding but pore penetrating tracer pulse in the order of minutes, and a gradient elution experiment in the order of hours. Thus, the automated smoothing method must work on data of very different time scales and still reliably isolate a signal from the noise that retains all relevant features and is smooth enough to obtain accurate first derivatives.

Most high frequency noise removal strategies fall into a few general categories. The Fast Fourier Transform (FFT) [32] converts a signal to frequency space where high frequency ranges can be removed before the signal is converted back to the original space. FFT filters are fast and suitable for removing high frequency noise. They can be robustly automated and include a simple method for computing the first derivative. Moving average filters and windowed polynomial regression, such as the Savitzky-Golay filter [33], can also locally approximate the signal with reduced noise. These methods depend on several parameters for controlling the window width and smoothing factors. Such filters can work quite well with humans choosing these parameters but are generally more difficult to robustly automate as compared to splines with only one parameter. Splines [34] are another alternative for reconstructing signals from noisy observations, and they can also provide first derivatives. However, splines can become computationally expensive due to an increased number of required knots for signals with high frequency noise. The required knot number can be automatically determined as will be detailed in the next paragraph.

We now introduce an automated smoothing procedure for noisy chromatograms that is implemented in CADET-Match. The procedure is illustrated using a previously published chromatogram [30] shown in Fig. 1A. This example is typical for industrial data and the signal contains particularly small and large features. Preliminary tests have revealed that no single choice of the above described methods can automatically and robustly remove measurement noise from chromatograms without filtering out relevant features over a wide range of time scales. However, satisfying results were achieved by combining an FFT based filter with a spline. Prior to the smoothing procedure, the chromatograms are normalized by dividing the concentrations of each component by the maximum

**Fig. 1.** Original and smoothed signal (A) and first derivative (B) of an example chromatogram taken from [30], critical frequency (C) and spline knots (D).

concentration of that component. This transformation is reversed after the smoothing procedure.

First, a third-order Butterworth low pass filter [35] is applied to the normalized chromatogram using the scipy.signal.butter function [36]. The Butterworth filter is as flat as possible in the passband below the critical frequency and then dampens the signal with 20 decibel per decade. A suitable critical frequency for smoothing a given chromatogram is automatically determined using the so-called elbow point method. Fig. 1C shows the logarithm of the normalized root mean square difference (NRMSD) between the filtered and original signals over the critical frequency of the filter. The elbow point maximizes the distance between that curve and a straight line between the extreme points. It indicates the best compromise between removing as much noise as possible while approximating the signal as accurately as possible.

Second, a 5<sup>th</sup> order spline with non-equidistant knots is applied to the low pass filtered chromatogram using the scipy.interpolate.UnivariateSpline function [36]. Cubic splines would be sufficient for approximating the original signal, but higher order splines are beneficial for computing derivatives. A suitable number of knots is determined using the elbow point method again. Fig. 1D shows the logarithm of the NRMSD between the approximated and original signals over the number of knots. This elbow point indicates the best compromise between using as few knots as possible while approximating the signal as accurately as possible. To the left of the elbow point, the NRMSD drops very quickly with an increasing number of knots. Then, the spline switches from smoothing to interpolating and the NRMSD decreases very slowly. Technically, a smoothness factor is passed to the SciPy function instead of the number of knots. The function internally computes the minimal number of knots for which the NRMSD falls below the specified smoothness factor.

The presented smoothing procedure also allows computing smooth first derivatives, Fig. 1B. It is routinely applied to every chromatogram, measured or simulated. The latter allows saving compute time by using rather coarse simulator tolerances in early stages of the parameter estimation procedure that could otherwise cause numerical problems with some of the applied search algo-

rithms. The smoothing procedure is applied once to each measured chromatogram. Iterative search strategies can involve numerous similar simulations with hardly varying elbow points for the resulting chromatograms. Hence, these points are pre-determined for a representative set of initial values and reused further on. Very smooth input data might not show distinct elbow points for the Butterworth filter and/ or the spline approximation. In these cases, the respective methods are simply disabled.

## 7. Goal system

We introduce a new goal system for estimating chromatography model parameters. Here, goal means a set of shape-sensitive metrics. Each metric is a single scalar value such as the time difference between simulation and measurement at peak maximum, the height difference at peak maximum, etc. Metrics are defined on the basis of specific knowledge of the modeled process and of typical errors in the measurement data. The metrics in a goal can be passed to a multi-objective search algorithm or they can be combined into one objective and passed to a single-objective search algorithm.

Multiple metrics can guide (multi-objective) search strategies much better to the desired optimum than the commonly applied sum of squared differences (SSD) can guide (single objective) search strategies, as will be demonstrated in the results section. The metrics are grouped into scores to organize the specification of goals for different parameter estimation procedures in CADET-Match. A suitable goal must have the property that as the fit quality improves the value of at least one metric must decrease and as the fit quality worsens, the value of at least one metric must increase. A goal that does not have this property can guide search algorithms in the wrong direction. This might appear trivial but is critically important and at the core of why new goals were designed. Due to competitive binding and other complex mechanisms, many model parameters influence the simulated chromatograms in non-linearly coupled ways. Therefore, some customarily applied metrics such as SSD can increase while the model parameters move closer to their correct values.

In addition, measured chromatograms from industrial large-scale applications are often affected by systematic errors such as pump delays that can cause a time offset between the measured and simulated signals, unless the model captures the cause of the delay which is often not possible in practice. Pump delays occur when there is a difference between when a pump is given the signal to start and when it starts. This data is usually not available and thus can't be modeled. To further complicate matters pump delays may not be consistent between runs or within a run. A good goal needs to account for this, since otherwise the simulated peaks end up in the correct location but with the wrong shape. Wrong shapes generally indicate errors in the underlying physics of the model. Hence, good metrics should prefer peaks with nearly fitting shape but small offsets rather than peaks without offset but with wrong shapes.

### 7.1. Sum of square difference

For the SSD, the squared differences between simulated and measured chromatograms are summed up over the time points, Eq. (12). For the NRMSD, the SSD is divided by the number of time points before taking the square root and dividing the result by the maximum of the measurement data, Eq. (13). Due to the monotonicity of this transformation, SSD and NRMSD have the same minima. These metrics can be applied to the entire chromatogram, $J = \{1; \ldots; N_d\}$, or a subset of the data, $J \subset \{1; \ldots; N_d\}$. The SSD is most commonly applied with gradient descent search algorithms.

**Table 1**
Synthetic data with resulting NRMSD for illustrating alignment issue.

|  | Ground Truth | Scenario 1 | Scenario 2 |
|---|---|---|---|
| $k_a$ | 2.00 | 2.9e02 | 2.00 |
| $k_d$ | 10.0 | 3.7e03 | 10.0 |
| $\nu$ | 7.00 | 9.60 | 6.00 |
| $\sigma$ | 50.0 | 99.0 | 50.0 |
| SSD |  | 4.3e+00 | 1.5e+01 |
| NRMSD |  | 4.2e-03 | 7.7e-03 |

Hence, it is included here for comparison. The theory is well established in the framework of maximum likelihood estimation for independent and identically distributed random measurement errors. However, these preconditions are generally not valid for modeling large-scale preparative chromatography where systematic errors such as feed variations, pump delays and flow rate variations typically dominate the detector noise. The NRMSD is better suited than the SSD for interpreting the results, because the numerator has the same unit as the data and is related to the maximum concentration by the denominator.

$$SSD(X_i, Y_i)_J = \sum_{j \in J} \left( X_{i,j} - Y_{i,j} \right)^2 \tag{12}$$

$$NRMSD(X_i, Y_i)_J = \frac{\sqrt{\frac{1}{|J|} \sum_{j \in J} \left( X_{i,j} - Y_{i,j} \right)^2}}{\max\limits_{j \in J} |X_{i,j}|} \tag{13}$$

The SSD requires a sufficient overlap between the simulated and measured chromatograms to be sensitive to parameter changes and guide the search algorithm towards the optimum. This can complicate the choice of suitable starting points, in particular for sharp and/ or small peaks. A further disadvantage of the SSD is illustrated in Fig. 2 using a synthetic example with parameters shown in Table 1. The parameters of scenario 2 are much closer to the ground truth, with only a relatively small deviation in the characteristic charge, $\nu$, even though Scenario 1 has a smaller SSD and would hence normally be considered a better fit. In addition, the peak shape of Scenario 2 is more similar to the ground truth but out of alignment.

In real experiments, such time offsets are often caused by pump delays that cannot be explained by the mechanistic model. In this case, the SSD favors peaks that are in the right position even though it is obvious to the human eye that the peak shape is wrong. The model can also reproduce the correct peak shape but not in the right position with a much larger SSD. As the peak shape is predominantly determined by the binding model parameters, the SSD would lead to unphysical parameter values. Hence, we will now introduce alternative metrics that favor peak shape over position and are less demanding on the choice of suitable starting points.

### 7.2. Alternative metrics

The shape and position of a chromatogram are determined by mass transport through the entire system, including the column and external volumes, and binding to the functionalized resin. The disadvantages of the SSD are avoided by separately measuring the shape, position, and height of individual peaks without requiring base line separation. Metrics for peak position are sensitive to changes of the respective model parameters, independent of peak overlaps between simulation and measured data. This provides flexibility and robustness with respect to the choice of starting points for the search algorithms, which is critically important for automation in industrial applications. Focusing on individual
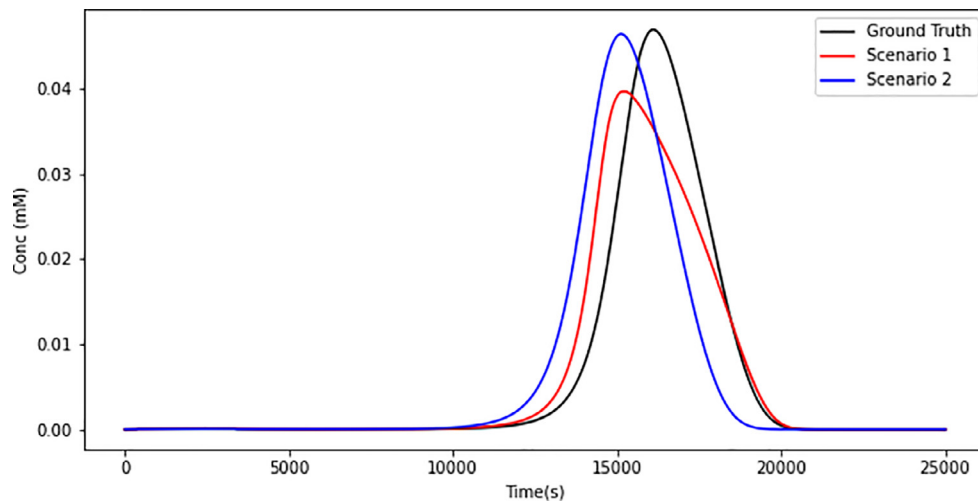
**Fig. 2.** SSD chromatogram alignment issue illustrated by synthetic data.

peaks allows to reduce the impact of process variations and further components that are not fully included in the model. For example, pump washes or pressure alarms can cause spurious peaks, and industrial feeds usually contain large numbers of more or less uncharacterized impurities. In such cases, separate metrics can be assigned to distinct but partially separated peaks of target components and impurities of high and low molecular weight. Separate metrics also help to provide (multi-objective) search algorithms with more precise information on which component impacts which peak, as will be detailed in Section 9. All metrics are designed for minimization and yield zero for a perfect match between simulation and experiment.

### 7.2.1. Peak shape

The shape metric is the most innovative of the new metrics and a core component of nearly all scores. It is the difference between one and the maximum of the Pearson correlation [37] between measured and simulated chromatograms over a continuous range of time offsets, Eq. (14). For evaluating this metric, the simulated chromatogram is shifted in time, $Y_i^\tau(t) = Y_i(t - \tau)$. The maximum in Eq. (14) is determined by an initial grid search followed by Powell's method. While it is advisable for the SSD to simulate the chromatogram on the same grid as the measurement data, continuous offsets require interpolating the simulated data. This is implemented in CADET-Match using the 5th order spline from the smoothing procedure described in Section 6. Allowing for continuous time offsets that are independent of the discrete measurement grid is crucial for creating a smooth metric. The shape metric is typically applied to individual peaks that are sliced out of the chromatogram. By design, this metric only accounts for shape similarity and requires two other metrics to measure the time offset and height difference between simulation and measurement data.

$$Shape(X_i, Y_i)_J = 1 - \max_\tau \left( \frac{cov(X_i, Y_i^\tau)_J}{\sigma_{X_i, J} \sigma_{Y_i^\tau, J}} \right) \tag{14}$$

### 7.2.2. Peak position

The position metric can be more complex than it might first appear. It is based on the time offset, $t_s$, obtained from maximizing Eq. (15).

$$t_s(X_i, Y_i)_J = \arg \max_\tau \left( \frac{cov(X_i, Y_i^\tau)_J}{\sigma_{X_i, J} \sigma_{Y_i^\tau, J}} \right) \tag{15}$$

The standard position metric gives an immediate penalty for a time offset with a linear ascent to one when out of alignment by $t_r$, Eq. (16). Here, $t_r$ is the length of the measurement time interval.

It can be replaced by the retention time of a non-binding tracer if sufficient starting points are provided to the search algorithm. As will be shown in the results section, this metric it a good choice for estimating column and particle porosity. However, it requires great care in running experiments to ensure there are as few delays as possible and alarms are immediately canceled. Such delays affect the chromatogram almost exactly like changes in the column and particle porosities. As previously discussed, in the presence of such delays it can be advantageous for the parameter estimation procedure to compromise on the alignment of simulated and measured peaks while matching their shape and height. Hence, an alternative position metric is introduced that initially reduces the penalty by $1/2$ in a range of less than $1/10 \, t_r$ and then linearly ascends to one when out of alignment by $t_r$, Eq. (17). Fig. 3 illustrates the difference between the standard and initially reduced position penalty metrics. The initial reduction, $1/2$, and range, $1/10$, are chosen by experience and can be changed by the user.

$$Position(X_i, Y_i)_J = \frac{t_s(X_i, Y_i)_J}{t_r} \tag{16}$$

$$Position^*(X_i, Y_i)_J = \begin{cases} \frac{1}{2} \frac{t_s(X_i, Y_i)_J}{t_r}, & \frac{t_s(X_i, Y_i)_J}{t_r} \leq \frac{1}{10} \\ \frac{19}{18} \frac{t_s(X_i, Y_i)_J}{t_r} - \frac{1}{18}, & \frac{t_s(X_i, Y_i)_J}{t_r} > \frac{1}{10} \end{cases} \tag{17}$$

### 7.2.3. Peak height

The peak height metric relates the maximal concentration of the simulated chromatogram to that of the measured data, Eq. (18). This metric ascends to one when the difference in either direction is larger than 100%.

$$Height(X_i, Y_i)_J = \left| 1 - \frac{\max\limits_{j \in J} Y_{i,j}}{\max\limits_{j \in J} X_{i,j}} \right| \tag{18}$$

### 7.3. Combined scores

The previously introduced metrics serve as building blocks for creating scores that quantify the difference between simulated chromatograms and measurement data. Scores are defined for individual components and can target the full chromatogram, individual peaks, or parts thereof, such as only the front of a peak. Each score is a set of metrics that depend on the index of the component, $i$, and on the set of considered time points, $J$. Goals will be composed of one or several scores that can then be combined into one objective or passed to a multi-objective search algorithm. The following scores have been defined per component this may
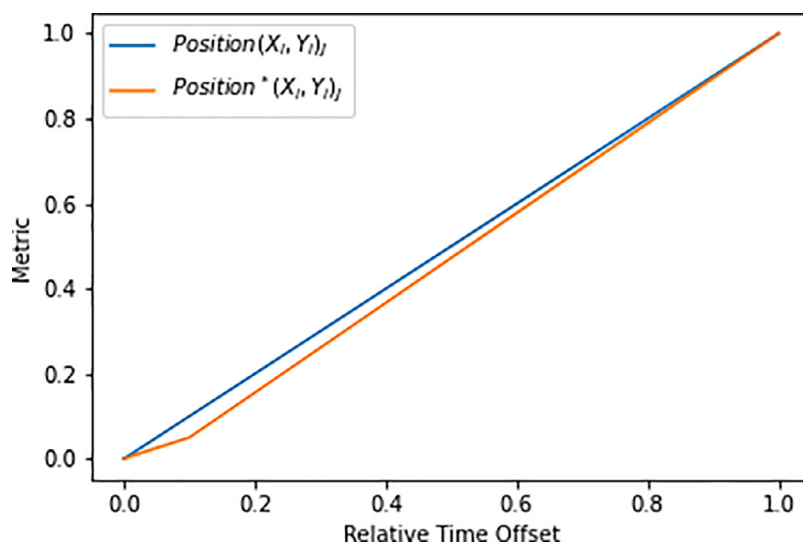
**Fig. 3.** Comparison of standard position penalty and initially reduced position penalty metrics.

include using a virtual component created by summing any combination of components.

### 7.3.1. Sum of square difference

The SSD score is the set of differences between simulated and measured chromatogram data, Eq. (19). For technical reasons, each difference is interpreted as separate metric. In Section 8, a goal will be defined as sum of squares of these metrics.

$$S_{SSD} = \left\{ X_{i,j} - Y_{i,j} \,\middle|\, j \in J \right\} \tag{19}$$

### 7.3.2. Full peak

The most straightforward of the new scores is the combination of peak shape, position, and height, $S_{Gauss}$. This score is usually applied to a time interval, specified by the index set, $J$, that contains a single peak of nearly Gaussian shape. This time interval is not automatically detected but needs to be specified by the user. A more elaborated score, $S_{Peak}$, additionally accounts for the shape, minimum and maximum of the time derivative and is better suited for fitting non-Gaussian peaks. By combining the peak shape with the shape of its derivative, this score is highly sensitive to the curvature of the chromatogram. The time offsets in the peak and in the slope are technically not constrained to be equal. In practice they hardly differ, except for the very first iterations of search algorithms with poor starting points. The scores $S_{Gauss}^*$ and $S_{Peak}^*$ are analogously defined with $Position^*(X_i, Y_i)_J$ in place of $Position(X_i, Y_i)_J$. As will be demonstrated in the results section, the score $S_{Peak}$ with standard position penalty is particularly useful for estimating transport parameters, while the score $S_{Peak}^*$ with initially reduced position penalty is more suitable for estimating binding parameters.

$$S_{Gauss} = \left\{ Shape(X_i, Y_i)_J; Position(X_i, Y_i)_J; Height(X_i, Y_i)_J \right\} \tag{20}$$

$$S_{Slope} = \left\{ Shape(\dot{X}_i, \dot{Y}_i)_J; Height\left(-\dot{X}_i, -\dot{Y}_i\right)_J; Height\left(\dot{X}_i, \dot{Y}_i\right)_J \right\} \tag{21}$$

$$S_{Peak} = S_{Gauss} \cup S_{Slope} \tag{22}$$

### 7.3.3. Peak front

In some cases, only the front of a peak can be used for parameter estimation while other parts of the peak are deteriorated by unspecific interactions of a tracer molecule with the column or tubing. Dextran is a prominent example for such non-ideal behavior that leads to strong tailing and a reduced peak height. On the other hand, Dextran is commonly applied as tracer that does not penetrate the particle pores. Errors in the execution of an experiment can also render the back of a peak unusable for parameter estimation. These situations are addressed by a score, $S_{Front}$, that considers shape and position but not height of the peak. This score is typically used with rather short time intervals and few data points.

$$S_{Front} = \left\{ Shape(X_i, Y_i)_J; Position(X_i, Y_i)_J \right\} \tag{23}$$

The peak front score is designed to extract as much usable information as possible from the chromatogram. Unsupervised application of this score requires to automatically determine the usable time interval while robustly removing the non-ideal parts with high precision on the cut points. For dextran data, the back end of this interval is chosen at the first inflection point of the measured chromatogram, i.e. the upper cut point is at the first maximum of the time derivative. By experience, this is a good choice as non-ideal interactions mainly impact on the height and tailing of the peak. The lower cut point is chosen where the measured chromatogram starts to differ from the baseline by more than 0.1% of the concentration at the upper cut point. By experience, 0.1% is a robust choice for this threshold. The exact positions of these cut points are determined using Powell's method on the continuous spline approximation from Section 5. The nearest time points of the discrete measurement data are then used as boundaries of the time interval specified by $J$.

### 7.3.4. Fractionation data

Optical detectors that are typically applied for measuring chromatograms can usually not distinguish between different chemical components. Instead, they deliver a single sum signal where the contributions of the individual components are weighted by their extinction coefficients. Such signals alone cannot be used for parameter estimation unless the peaks of the relevant components are sufficiently separated. For instance, the acidic, main, and basic components of a monoclonal antibody often completely overlap in a single peak. This situation is normally addressed by fractionation, i.e. pooling the efflux of the column into a series of vials. Each of these vials is then analyzed offline to quantify the components of interest, which provides additional information for setting up a dedicated parameter estimation score.

The previously introduced metrics can generally be applied to the concentrations in each vial using the centers of the corresponding collection intervals as time points. For precise comparison, the corresponding per component simulations are averaged over the same collection intervals when a metric is applied to fractionation data. The resulting information is often sparse, with 5 to 10 fractions per peak, and can be afflicted with additional errors in the fractionation times and volumes. Small shifts in the collection intervals can cause major changes in the distribution of components between the analyzed fractions, particularly for sharp peaks. A spline is applied to the original simulated data before it is shifted and virtually fractionated to determine the time offset, $t_s$, in order to maintain sub-grid accuracy. Based on this offset, the scores $S^*_{Gauss}$ and $S^*_{Peak}$ as well as their immediately penalized versions can be computed. Analogously, the SSD score can be applied to fractionation data by averaging the simulations over the collection intervals.

## 8. Search strategies

CADET-Match uses two alternative search strategies, gradient descent, and a multi-objective genetic algorithm. For gradient descent, all metrics need to be combined into a single scalar value, while the genetic algorithm can operate on multiple metrics.

### 8.1. Gradient descent

Gradient descent algorithms search for a local optimum of the goal function using derivative information with respect to the sought parameters [38]. Gradient descent has long been used for parameter estimation in chromatography. It is very efficient near the sought optimum but can fail if the goal function is not smooth or the Jacobian becomes singular. Moreover, this algorithm is prone to becoming trapped in local optima, which can be far from the global optimum. This can be avoided by basin hopping or multi-start strategies. The latter is often applied when refining the results of population-based search strategies.

### 8.2. Genetic algorithm (GA)

Genetic algorithms (GA) where first published by Holland in 1975 [39] and are an example of biomimicry. At their core they work like a colony of bacteria adapting to an outside environment and share many of the same features. GAs are embarrassingly parallel. An initial population is created, often using quasi-random methods such as Latin hypercube sampling [40] or Sobol sequences [41]. Each member of the population is then evaluated based on one or more objectives. At the end of each generation the fittest members survive and reproduce to form the next generation. The next generation is created by a combination of breeding and mutation on the surviving members. There are variations in this procedure that maintain population diversity, choose members to be included in the next generation and change how breeding and mutation are implemented. These variations result in different algorithms such as NSGA2 [42], NSGA3 [43] and SPEA2. In view of the no free lunch theorem [1], different GA variants were tested and optimized on a variety of problems before settling on NSGA2 for single-objective problems and NSGA3 for multi-objective problems.

### 8.3. Progress monitoring

Building complex models correctly, properly processing experimental data, and determining suitable starting points can be difficult and tedious tasks. Based on experience, new models or concepts are unlikely to be correctly implemented on the first attempt. However, such issues can only be tested by attempting to fit model to data. Since errors can often be identified in early phases of the parameter estimation process, CADET-Match provides functionality to monitor the progress of specific indicators such as peak height, shape, mass, etc. This allows observing if the starting points yield reasonable results and if the search algorithm continuously improves the goal. Online monitoring enables early aborting if progress is poor or if results are already good enough. This is essential for rapid testing of models, goals, starting points and stopping criteria. Since suitable starting points can be hard to determine, a GA with rather large population size is generally a good choice for initial testing. Multi-start gradient search is not a good alternative, as parallel iterative processes are more difficult to monitor.

### 8.4. Parameter transformation

Most search strategies struggle when the parameters to be estimated are spread over orders of magnitude or correlated with each other. Parameter transformations can help to soften these challenges. CADET-Match provides several transformation rules, i.e. bi-unique maps between model parameters, $p$, that are passed to the chromatography simulator and estimated parameters, $p'$, that are passed to the search algorithm. These transformations are based on upper bounds, $\hat{p}$, and lower bounds, $\check{p}$, of the model parameters.

The linear transformation, Eq. (24), maps the original range $[\check{p}, \hat{p}]$ to [0,1]. This is usually sufficient when the upper and lower parameter bounds are less than three orders of magnitude apart from each other. For wider ranges, a nonlinear transformation, Eq. (25), is advisable. The latter automatically adapts the step width of the search algorithm to the magnitude of the respective model parameter. Otherwise, the same step could be huge for one parameter but tiny for another.

$$p = \left(\hat{p} - \check{p}\right) \cdot p' + \check{p} \tag{24}$$

$$p = \exp\left(\left(\log\left(\hat{p}\right) - \log\left(\check{p}\right)\right) \cdot p' + \log\left(\check{p}\right)\right) \tag{25}$$

Nonlinear parameter correlations are hard to detect and need to be specifically addressed. For instance, the adsorption and equilibrium constants, $k_a$ and $k_{eq}$, are usually much less correlated than the adsorption and desorption constants, $k_a$ and $k_d$. The relation $k_{eq} = k_a/k_d$ allows to pass $k_a$ and $k_d$ to the simulator while the search algorithm operates on $k_a$ and $k_{eq}$. The corresponding transformation, Eqs. (26) and 27, also accounts for large parameter ranges. This decouples the binding rate from the concentration equilibrium.

$$k_a = \exp\left(\left(\log\left(\hat{k}_a\right) - \log\left(\check{k}_a\right)\right) \cdot k'_a + \log\left(\check{k}_a\right)\right) \tag{26}$$

$$k_d = \frac{\exp\left(\left(\log\left(\hat{k}_a\right) - \log\left(\check{k}_a\right)\right) \cdot k'_a + \log\left(\check{k}_a\right)\right)}{k'_{eq}} \tag{27}$$

## 9. Practical application

The goal system and search strategy are first verified on synthetic examples of increasing complexity (Section 10) and then validated on experimental data (Section 11). The parameter estimation procedure is highly automated and the same for all case studies.

Depending on the current stage in the parameter estimation procedure (Section 4) and on the quality of the data, the goals are based on one of the $S_{Front}$, $S_{Peak}$ or $S^*_{Peak}$ scores and separately on the $S_{SSD}$ score for comparison. While SSD is used for search, the results of different scores and search algorithms are compared using NRMSD. The $S_{SSD}$ scores are usually taken over

the whole chromatogram, $J = \{1, \ldots, N_d\}$. For single-objective optimization, the scores of multiple components are simply merged into one set. In case only a sum signal can be measured, the corresponding simulations are also summed, weighted by the respective extinction coefficients if necessary, and both sum signals are compared as in the single-component case. The corresponding goal, $G_{SSD}$, is created by adding the squared metrics in this set, Eq. (28). Single-objective goals for the alternative scores are similarly created by merging the scores of different components and adding the squared metrics in the resulting set. CADET-Match offers other single objective goals, such as the average or the maximum of the metrics, but they are not used in this study.

$$G_{SSD} = \sum_{i=1}^{N_c} \sum_{j=1}^{N_d} \left(X_{i,j} - Y_{i,j}\right)^2 \qquad (28)$$

Multi-objective optimization is more demanding on the search strategy but can benefit from much richer information on the impact of individual parameters or parameter groups on the path from the starting points to the sought optimum. For example, column porosity affects the peak position metric more strongly than the peak shape metric and column dispersion affects the peak shape metric more strongly than the peak position metric. Providing these metrics to a multi-objective search algorithm practically decouples the parameters, i.e. progress can be made in the objectives independent of each other. If these objectives are not in conflict, which is most often the case, the search algorithm still converges to a unique optimum. In more complex settings with numerous parameters and metrics, progress can be made in several objectives while other objectives are temporarily sacrificed even when they are not conflicting in the global optimum. The multi-objective approach improves the performance of genetic algorithms, as these are population-based by nature, while gradient search is more efficient on single objectives.

In addition, multi-objective search automatically detects conflicting metrics, in case they exist, and provides detailed information on the resulting Pareto front [44]. This usually indicates a tradeoff, e.g., between matching shape and position as illustrated in *Fig. 2*. In such cases, the user can manually select the preferred optimum. To fully automate this process, the Pareto optimal parameter set with lowest mean of the involved metrics is selected. Eq. (29) defines the mean, $\bar{S}$, of a given set, $S$, of metrics, $M$, with $|S|$ elements. The metrics range from 0 to 1 with 0 indicating a perfect match and 1 a very poor match. This scale is reversed before and after taking the geometric mean, as otherwise one very low metric could potentially dominate the entire mean. For comparison, the mean, $\bar{S}$, is reported on the final result even when the optimization was based on the $S_{SSD}$ score.

$$\bar{S} = 1 - \left(\prod_{M \in S} (1 - M)\right)^{\frac{1}{|S|}} \qquad (29)$$

The search space is confined by box constraints that are specified in the case studies. A Sobol sequence is used to create a starting population with 100 $N_p$ individuals, where $N_p$ is the number of concurrently estimated parameters. When gradient descent is used with $S_{SSD}$, starting points with less than 5% peak overlap between simulation and experiment are removed from the population without replacement to guarantee sufficient sensitivity of the goal with respect to changes in the estimated parameters. This is not required for the new scores $S_{Front}$, $S_{Peak}$ and $S^*_{Peak}$ that include the peak position metric. This metric is naturally sensitive to parameter changes, independent of the peak overlap.

For gradient descent, the scores are always combined into a single objective. The trust-region reflective algorithm [45] from the scipy.optimize.least_squares function [36] is started at each point of the population. These searches are independent of each other

**Table 2**
Fixed parameters for simulating synthetic ground truth.

| Parameter | Value | Unit |
|---|---|---|
| $Q$ | 2.88e-08 | $m^3/s$ |
| $A_c$ | 1.04e-04 | $m^2$ |
| $L_c$ | 2.50e-01 | $m$ |
| $r_p$ | 4.50e-05 | $m$ |
| $\Lambda$ | 2.25e+00 | $mol/m^3$ |

and run in parallel to save time. They are stopped with a tolerance of $x_{tol} = 10^{-10}$ or if the current simulation fails. Looser tolerances were tested on synthetic examples and found to terminate too soon in some cases.

The non-dominated sorting genetic algorithm (NSGA) is used with a single objective (NSGA2) for $S_{SSD}$ and multiple objectives (NSGA3) for $S_{Front}$, $S_{Peak}$ and $S^*_{Peak}$. Distributed Evolutionary Algorithms in Python (DEAP) [46] is used for both algorithms. Crossover and mutation rates of 1.0 are applied to ensure that in every generation each member will mutate at least one parameter and swap parameters with at least one other member. The GA terminates after 30 generations without a new point being added to the Pareto front. Each entry on the final Pareto front is locally optimized in parallel using the trust-region reflective algorithm from above. This is a single-objective search, but the results are again analyzed with respect to Pareto optimality based on the individual metrics.

The reported wall clock times should be understood as approximate due to the non-deterministic nature of the GA, and small changes in stopping criteria can cause substantial differences in runtime. Moreover, the GA utilizes the available compute cores more efficiently than gradient search. All simulations where run on a dual socket Intel(R) Xeon(R) CPU E5–2683 v4 @ 2.10 GHz with a total of 32 cores and 64 threads with Ubuntu Linux 18.04.1, Python 3.7.7, Intel MKL 2019.3.199, CADET 4.0.1 and CADET-Match 0.6.13. Version information on other used software packages is shown in Table S1. CADET-Match can be installed from the Python Package Index (PyPI). The full code, including the scripts for the following case studies, is freely available on GitHub (https://github.com/modsim/CADET-Match).

## 10. Synthetic case study

The synthetic case study was designed using the general rate model, based on experience with experimental data so that the parameters are within plausible ranges. They have not been tuned to make the output particularly easy or hard to match. All parameters are reported to two decimal places. The fixed parameters for simulating the synthetic ground truth data are shown in Table 2. In practice, these parameters are specified or separately determined by other means. The ground truth of the estimated parameters is reported in the following subsections. All synthetic data has a time spacing of 1 value per second. The column and tubing start off with a concentration of 0.0 except in the case of the SMA binding model where the column begins with bound salt equal to the ionic capacity (Eq. (10)) and a liquid phase salt concentration equal to the loading salt concentration. Independent and identically distributed Gaussian noise with mean 0.0 and standard deviation 0.1% of peak maximum is added to the simulated chromatogram of each component.

The synthetic cases are designed to verify the functioning of the goal system and search strategies. They are presented in the same staged order that the parameters are normally estimated in (Section 4). However, the results of one stage are not carried over to the next stage, but previously estimated parameters are fixed at

**Table 3**

Non-binding and non-pore penetrating tracer pulse experiment. Synthetic ground truth, parameter bounds, estimated parameters and performance indicators.

| | Ground Truth | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|---|
| | | Lower | Upper | $S_{Front}$ | $S_{SSD}$ | $S_{Front}$ | $S_{SSD}$ |
| $\varepsilon_c$ | 0.40 | 0.20 | 0.50 | 0.40 | 0.40 | 0.40 | 0.40 |
| $D_c$ | 3.0e-07 | 1.0e-12 | 1.0e-5 | 3.0e-07 | 3.0e-07 | 3.0e-07 | 3.0e-07 |
| NRMSD | | | | 9.2e-05 | 1.4e-07 | 7.6e-06 | 8.5e-08 |
| $\bar{S}$ | | | | 3.5e-07 | 3.3e-06 | 3.3e-06 | 3.3e-06 |
| Wall Time | | | | 0:03:10 | 0:03:37 | 0:04:44 | 0:02:24 |

their nominal values, to test all goal and search algorithm combinations with a defined ground truth. In the synthetic case study, multiple objectives of the new scores are generally not in conflict. The mean of the involved metrics, $\bar{S}$, is shown for comparison with the experimental case study, where it is used to select a Pareto optimal parameter set, but not further discussed for the synthetic case study. A match is considered successful if the sought parameters are correctly estimated to two digits.
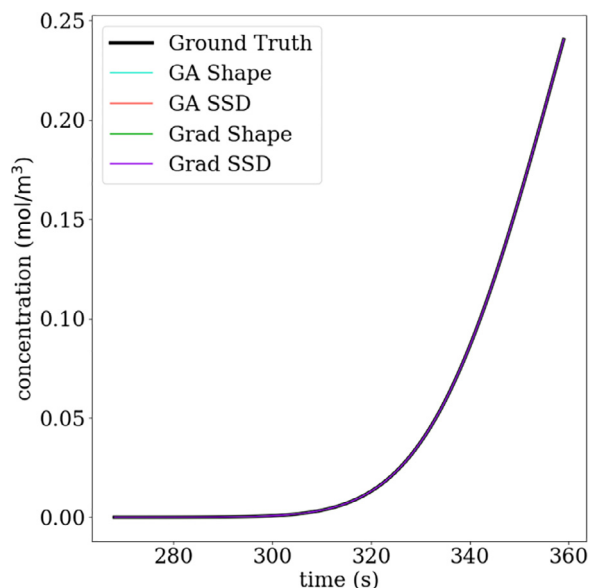
The parameter transformations from Section 8.4 are automatically applied. If the upper and lower bounds are less than 3 orders of magnitude apart, the linear transformation is used, Eq. (24). Otherwise, the logarithmic transformation is used, Eq. (25), except for the custom transformation for $k_a$ and $k_d$, Eqs. (26) and (27). The reference concentrations of the SMA model are $c_r^s = 225 \ mol/m^3$ and $c_r^p = 450 \ mol/m^3$.

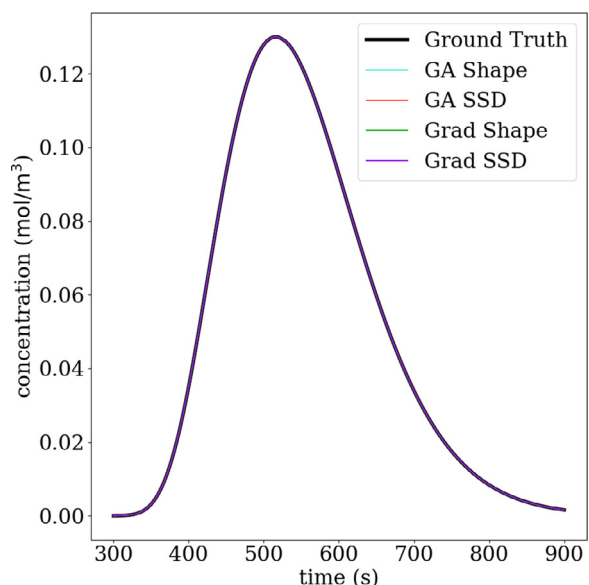### 10.1. Non-Binding and non-pore penetrating tracer pulse

The first parameter estimation stage is omitted here, as the synthetic data is generated without tubing and other holdup volumes that are external to the column. In the second stage, column porosity, $\varepsilon_c$, and axial dispersion, $D_{ax}$, are estimated from a pulse experiment with a non-binding and non-pore penetrating tracer molecule. In theory, these model parameters do hardly depend on the chosen tracer molecule. In real experiments, dextran is typically used for this purpose. This can be problematic due to non-ideal behavior of this tracer that often leads to strong tailing and a reduced peak height, as has been introduced in Section 7.3. Only the peak front is used because the rest of the peak is often deteriorated by unspecific interactions of dextran with the column. Since this non-ideal behavior is not covered by the applied chromatography model, it cannot be reproduced in the peak tailing. However, an idealized setting is used for verifying the respective score, $S_{Front}$. The time interval $J$ of this score is determined as described in Section 7.3. It ranges from 269 to 359 s (Figure S1). For comparison, the $S_{SSD}$ score is applied to the same slice of the chromatogram. Table 3 shows the parameter bounds that were specified in the search algorithms. The parameter estimation results are shown in Table 3 and Fig. 4. Any combination of tested search algorithm and score were able to recover the estimated parameters to prescribed accuracy in less than five minutes. The NRMSD is about 0.01% of the peak maximum for GA and $S_{Front}$ and lower in the other cases. For the human eye, the fitted chromatograms in Fig. 4 are practically indistinguishable from the synthetic ground truth, even 100x magnified (Fig. S1).

### 10.2. Non-Binding but pore-penetrating tracer pulse

In the third stage, particle porosity, $\varepsilon_p$, film diffusion, $k_f$, and pore diffusion, $D_p$, are estimated from a non-binding but pore-penetrating tracer pulse. These parameters depend on the size and shape of the chosen tracer molecule more strongly than porosity and axial dispersion in the column. Hence, it is advisable to use the target protein under non-binding conditions, such as high salt



**Fig. 4.** Non-binding and non-pore penetrating tracer pulse experiment. Synthetic ground truth and estimated chromatograms.
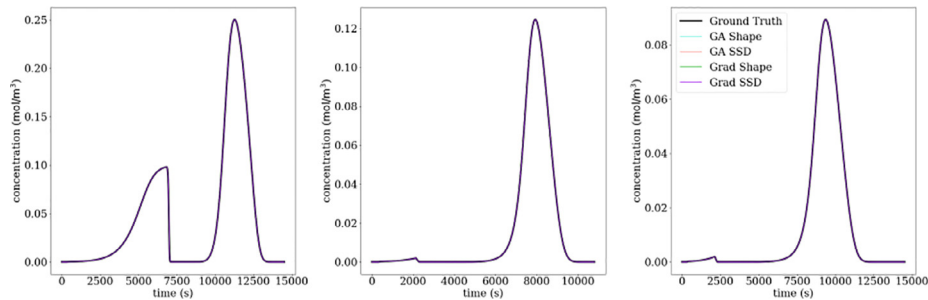


**Fig. 5.** Non-binding but pore penetrating tracer pulse experiment. Synthetic ground truth and estimated chromatograms.

in ion-exchange chromatography. This setting is used for verifying the score $S_{Peak}$ in comparison with $S_{SSE}$. The time interval $J$ covers the whole peak (Fig. S2). Results are shown in Table 4 and Fig. 5. Also in this case, any combination of tested search algorithm and score were able to recover the sought parameters to pre-

**Table 4**

Non-binding but pore penetrating tracer pulse experiment. Synthetic ground truth, parameter bounds, estimated parameters and performance indicators.

| | Ground Truth | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|---|
| | | Lower | Upper | $S_{Peak}$ | $S_{SSD}$ | $S_{Peak}$ | $S_{SSD}$ |
| $\varepsilon_p$ | 0.30 | 0.2 | 0.5 | 0.30 | 0.30 | 0.30 | 0.30 |
| $k_f$ | 2.0e-07 | 1.0e-09 | 1.0e-05 | 2.0e-07 | 2.0e-07 | 2.0e-07 | 2.0e-07 |
| $D_p$ | 5.0e-11 | 1.0e-14 | 1.0e-06 | 5.0e-11 | 5.0e-11 | 6.0e-11 | 5.0e-11 |
| NRMSD | | | | 7.5e-07 | 7.5e-07 | 4.7e-04 | 7.5e-07 |
| $\bar{S}$ | | | | 2.2e-08 | 1.8e-08 | 2.3e-05 | 1.6e-08 |
| Wall Time | | | | 0:04:08 | 0:04:16 | 0:05:06 | 0:07:48 |



**Fig. 6.** Single component gradient elution experiments with different loading times and gradient slopes. Synthetic ground truth and estimated chromatograms.

scribed accuracy, now in under eight minutes. The NRMSD is 0.05% of the peak maximum for gradient search and $S_{Peak}$ and lower in the other cases. For the human eye, the fitted chromatograms in Fig. 5 are practically indistinguishable from the synthetic ground truth. When 100x magnified, a slight difference can be observed for gradient search with $S_{Peak}$ (Fig. S2).

### 10.3. Single component gradient elution

In the fourth stage, the parameters of the binding model are estimated from gradient elution data. The parameters of the non-equilibrium SMA binding model, i.e. scaled adsorption and desorption rates, $\tilde{k}_a$ and $\tilde{k}_d$, characteristic charge, $\nu$, and shielding coefficient, $\sigma$, are estimated from gradient elution data. Only the ionic capacity, $\Lambda$, is set to the ground truth, as this parameter is usually determined separately by titration. Estimating the binding parameters is particularly difficult due to strong non-linearity of the SMA model. Hence, synthetic data of three linear gradient elution experiments with different gradient slopes are used, and the loading phase of the first experiment is extended to a full breakthrough, as shown in Fig. 6. The minor peaks in the other two experiments indicate complete loading.

The column is loaded at inlet salt and protein concentrations of $c_0^{in} = 180$ $mol/m^3$ and $c_1^{in} = 0.10$ $mol/m^3$ for 6500 $s$ (Fig. S3) in the first experiment and for 1800 $s$ in the second (Fig. S4) and third experiment (Figure S5). It is washed at an inlet salt concentration of $c_0^{in} = 70$ $mol/m^3$ for 2000 $s$ in all three experiments. The elution gradients start at an inlet salt concentration of $c_0^{in} = 70$ $mol/m^3$ and have slopes of 0.08 $mol/(m^3 \cdot s)$, 0.06 $mol/(m^3 \cdot s)$ and 0.04 $mol/(m^3 \cdot s)$. They are applied for 6000 $s$, 7000 $s$ and 11000 $s$. The time interval of the $S_{Peak}$ score on the full breakthrough in the first experiment was chosen from 1000 $s$ to 7300 $s$. The minor peaks in the second and third experiments indicate complete loading but are not utilized for parameter estimation. The time intervals of the $S_{Peak}$ score on the elution peaks were chosen from 8500 $s$ to 14000 $s$, from 5000 $s$ to 10000 $s$ and from 6000 $s$ to 12000 $s$.

The sought parameters are estimated by fitting separate model instances with respective boundary conditions simultaneously to all synthetic experiments. The $S_{Peak}^*$ score is applied to the four ma-

jor peaks and is composed of six metrics, resulting in $4 \cdot 6 = 24$ objectives for the genetic algorithm. A single-objective goal for gradient search is created by adding the squared metrics. The $S_{SSE}$ score is applied to the same slices of the chromatograms. The parameter transformation in Eqs. (26) and (27) is applied with $\tilde{k}_{eq} = \tilde{k}_a/\tilde{k}_d$. That is, the search algorithms operate on $\tilde{k}_a'$ and $\tilde{k}_{eq}'$ instead of $\tilde{k}_a$ and $\tilde{k}_d$. Results are shown in Table 5 and Fig. 6. The values of $\tilde{k}_d$ are reported for completeness. Search bounds are not needed for this parameter. Any combination of search algorithm and score were able to recover the sought parameters in less than 10 h. For the $S_{Peak}^*$ score, the GA was about two times faster than gradient search and about three times faster for the $S_{SSD}$ score. The NRMSD is less than 0.01% of the peak maximum for gradient search and $S_{Peak}^*$ and lower in the other cases.

### 10.4. Multi component gradient elution

This case study is analogous to the previous one but with two components. The synthetic ground truth mimics charge variants of a protein without baseline separation. The transport parameters of the first three stages are the same for both components. SMA binding parameters of both components are simultaneously estimated from three gradient elution experiments, including one full breakthrough. Full chromatogram data is assumed to be known for each component. This is not normally given in practice but an important intermediate step in verifying the goal system and search strategies. The assumption will be dropped in the next case study. In this case, the breakthrough peaks at the end of all three loading phases are large enough to provide useful information. Hence, the $S_{Peak}^*$ score is applied to two components and six peaks, resulting in $2 \cdot 6 \cdot 6 = 72$ objectives.

The inlet concentrations and durations of the load, wash and elution phases are the same as in the single component case (Figs. S6-S8), with the same inlet concentrations for both proteins. The time interval of the $S_{Peak}^*$ score on the full and partial breakthroughs were chosen from 0 $s$ to 7300 $s$, from 0 $s$ to 2500 $s$ and from 0 $s$ to 2500 $s$. The time intervals of the $S_{Peak}^*$ score on the elution peaks were chosen from 9000 $s$ to 14000 $s$, from 5000 $s$ to 10000 $s$ and from 5000 $s$ to 12000 $s$ for the first protein and
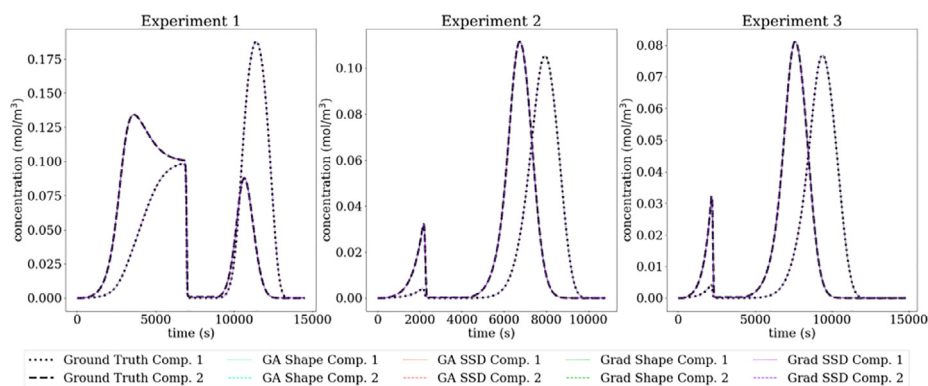
**Table 5**

Single component gradient elution experiment. Synthetic ground truth, parameter bounds, estimated parameters and performance indicators.

| | Ground Truth | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|---|
| | | Lower | Upper | $S^*_{Peak}$ | $S_{SSD}$ | $S^*_{Peak}$ | $S_{SSD}$ |
| $\tilde{k}_a$ | 0.30 | 1.0e-02 | 1.0e+02 | 0.30 | 0.30 | 0.30 | 0.30 |
| $\tilde{k}_d$ | 1.50 | – | – | 1.50 | 1.50 | 1.50 | 1.50 |
| $\tilde{k}_{eq}$ | 0.20 | 1.0e-02 | 1.0e+02 | 0.20 | 0.20 | 0.20 | 0.20 |
| $\nu$ | 7.00 | 1.00 | 50.0 | 7.00 | 7.00 | 7.00 | 7.00 |
| $\sigma$ | 50.0 | 1.00 | 100 | 50.0 | 50.0 | 50.0 | 50.0 |
| NRMSD | | | | 1.8e-06 | 1.8e-06 | 6.5e-05 | 1.8e-06 |
| $\bar{S}$ | | | | 1.4e-04 | 1.4e-04 | 1.7e-05 | 1.4e-04 |
| Wall Time | | | | 5:58:09 | 3:29:55 | 9:29:12 | 9:00:45 |

**Table 6**

Two component gradient elution experiment. Synthetic ground truth, parameter bounds, estimated parameters and performance indicators.

| | Ground Truth | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|---|
| | | Lower | Upper | $S^*_{Peak}$ | $S_{SSD}$ | $S^*_{Peak}$ | $S_{SSD}$ |
| $\tilde{k}_{a,1}$ | 2.00 | 1.0e-02 | 1.0e+02 | 2.00 | 2.00 | 1.90 | 2.00 |
| $\tilde{k}_{d,1}$ | 10.0 | – | – | 10.0 | 10.0 | 9.70 | 10.0 |
| $\tilde{k}_{eq,1}$ | 0.20 | 1.0e-02 | 1.0e+02 | 0.20 | 0.20 | 0.20 | 0.20 |
| $\nu_1$ | 7.00 | 1.00 | 50.0 | 7.00 | 7.00 | 7.00 | 7.00 |
| $\sigma_1$ | 50.0 | 1.00 | 100 | 50.0 | 50.0 | 50.1 | 50.0 |
| $\tilde{k}_{a,2}$ | 2.00 | 1.0e-02 | 1.0e+02 | 2.00 | 2.00 | 2.40 | 2.00 |
| $\tilde{k}_{d,2}$ | 10.0 | – | – | 10.0 | 10.0 | 12.0 | 10.0 |
| $\tilde{k}_{eq,2}$ | 0.20 | 1.0e-02 | 1.0e+02 | 0.20 | 0.20 | 0.20 | 0.20 |
| $\nu_2$ | 5.00 | 1.00 | 50.0 | 5.00 | 5.00 | 5.00 | 5.00 |
| $\sigma_2$ | 50.0 | 1.00 | 100 | 50.0 | 50.0 | 49.8 | 50.0 |
| NRMSD | | | | 3.3e-05 | 3.3e-05 | 1.8e-03 | 3.3e-05 |
| $\bar{S}$ | | | | 9.9e-06 | 9.9e-06 | 6.1e-04 | 9.9e-06 |
| Wall Time | | | | 2 days: 21:23:58 | 1 day: 20:08:52 | 5 days: 12:20:55 | 3 days: 23:09:31 |



**Fig. 7.** Two component gradient elution experiment with different loading times and gradient slopes. Synthetic ground truth and estimated chromatograms.

from 9000 *s* to 13000 *s*, from 4000 *s* to 9000 *s* and from 4000 *s* to 11000 *s* for the second protein.
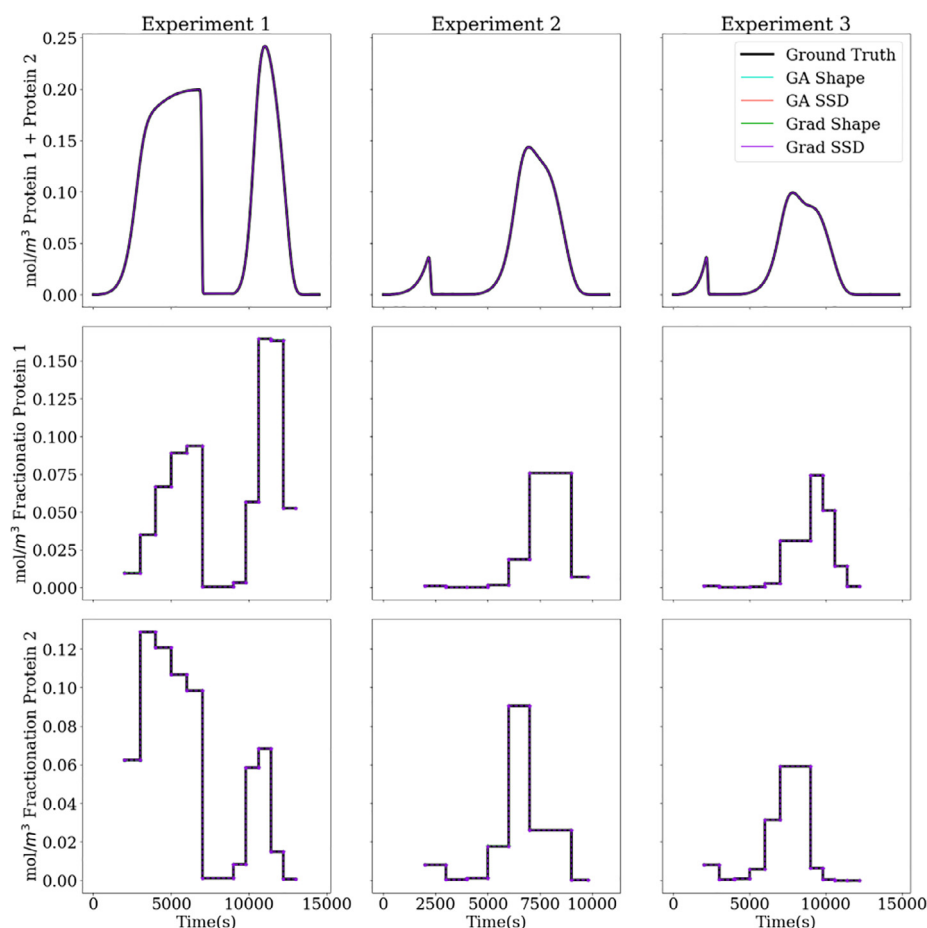
Results are shown in Table 6 and Fig. 7. Any combination of search algorithm and score were able to recover the estimated parameters, except for the scaled adsorption rates, $\tilde{k}_a$, and shielding coefficients, $\sigma$, when estimated using $S^*_{Peak}$ and gradient search. The shielding coefficients deviate from the ground truth by less than 0.5%. Even though the scaled adsorption rates are off by 5% and 20% for the first and second component, the corresponding NRMSD is below 0.2% of the peak maximum, which is still hard to observe by the human eye. This exemplifies that parameters with little impact on the model prediction are harder to estimate, which can be tolerated if accurate model predictions are sufficient for the application. However, it becomes problematic when accurate parameter values are required, e.g., for tabulation and application in scenarios with higher impact on model predictions.

The scaled desorption constants, $\tilde{k}_d$, which are not directly estimated but deduced from the respective recovered equilibrium constants, $\tilde{k}_{eq}$, are off by 3% and 20%. The increased compute time

of more than 5 days also indicates a flat optimum which is more challenging for gradient search than for the GA. We speculate that the combination of $S^*_{Peak}$ and gradient search is particularly sensitive to the smoothing procedure, and this issue is subject of ongoing research. The other combinations improved the NRMSD by about 2 orders of magnitude, and the same is true for the mean, $\bar{S}$, of the 72 metrics in the optimum. The other combinations ran between 2 and 4 days with GA on $S_{SSD}$ being the fastest.

### 10.5. Multi component gradient elution with fractionation

The last synthetic case is close to real experimental data, as the assumption that the full chromatogram of each component can be separately observed is dropped. Instead, only the sum signal and fractionation data are available. The basic setup is parallel to the previous section but with additional fractionation data. Fractionation is equivalent to an extremely coarse discretization of the signal and can yield less than 10 data points compared to thousands in a normal chromatogram. In practice, the column efflux is parti-

**Fig. 8.** Two component gradient elution experiment with different loading times and gradient slopes. Synthetic ground truth and estimated chromatograms: sum signal (top), fractionation data of component 1 (mid) and 2 (bottom).

tioned into a series of samples that are analyzed offline to determine the concentration of each component. In industrial settings, parameter estimation is further complicated by additional sources of error from the fractionation process. The impact of such errors is currently studied but not in the scope of this publication.

Here, each fraction is approximately 1000 s in length. Precise collection times and pool concentrations are shown in Tables S2-S4. The parameters of the SMA binding model are estimated using the $S^*_{Peak}$. score on the sum chromatogram and the $S^*_{Gauss}$ score on the fractionation data. $S^*_{Peak}$ with six metrics on six peaks and $S^*_{Gauss}$ with three metrics to three data sets and two components results in $1 \cdot 6 \cdot 6 + 2 \cdot 3 \cdot 3 = 54$ objectives. The $S_{SSE}$ score is applied to the sum chromatogram and to the fractionation data of each component.

Results are shown in Table 7 and Fig. 8. Similar to the previous case, any combination of search algorithm and score were able to recover the estimated parameters, except the combination of $S^*_{Peak}$ and gradient search. The shielding coefficients again deviate from the ground truth by less than 0.5%, and the scaled adsorption rates by 10% and 20% for the first and second component. Potential reasons for this have been explained in the previous section. The NRMSD is below 0.1% of the peak maximum. The NRMSD is more than two orders of magnitude smaller for the other algorithm and sore combinations. Differences are hardly visible to the human eye. This is remarkable, since this case is based on much sparser experimental data than the previous one. However, systematic errors that typically occur in real experiments are not considered here. They will be present in the next case. The runtimes mainly differ between the algorithms and are similar between the scores,

with the GA needing less than 2.5 days and gradient search about 5 days.

## 11. Experimental case study

An experimental dataset for two charge variants of a monoclonal antibody has been measured at Amgen. Available are two dextran pulses with detached column, two dextran pulses with attached column, two protein pulses under non-binding conditions, a gradient elution with fractionation and a gradient elution with extended loading phase but without fractionation. Evaluating industrial data is more complicated than synthetic, since there is no ground truth available. Moreover, the detector noise is typically dominated by systematic errors such as feed variations, pump delays and flow rate variations. These issues are addressed by our new score system. The experimental case study comprises the four parameter estimation stages described in Section 4. Parameters estimated in one stage are fixed in the next, with the results separately passed on for each search algorithm and score combination using the selected result shown in each stage's table. Orthogonality in the applied models and procedures avoids lumping of fundamental mechanisms and minimizes parameter correlations. This greatly improves predictivity of the calibrated model across operating conditions and scales. Modeling and propagation of errors is actively researched but not in the scope of this publication.

In the experimental case study, multiple objectives of the new scores happen to be in conflict due to imperfections in the model and data. For a new optimum to be added to the Pareto front, at least one parameter and at least one metric need to differ from

**Table 7**

Two component gradient elution experiment with fractionation. Synthetic ground truth, parameter bounds, estimated parameters and performance indicators.

| | Ground Truth | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|---|
| | | Lower | Upper | $S^*_{Peak}$ | $S_{SSD}$ | $S^*_{Peak}$ | $S_{SSD}$ |
| $\tilde{k}_{a,1}$ | 2.00 | 1.0e-02 | 1.0e+02 | 2.00 | 2.00 | 1.80 | 2.00 |
| $\tilde{k}_{d,1}$ | 10.0 | – | – | 10.0 | 10.0 | 8.80 | 10.0 |
| $\tilde{k}_{eq,1}$ | 0.20 | 1.0e-02 | 1.0e+02 | 0.20 | 0.20 | 0.20 | 0.20 |
| $v_1$ | 7.00 | 1.00 | 50.0 | 7.00 | 7.00 | 7.00 | 7.00 |
| $\sigma_1$ | 50.0 | 1.00 | 100 | 50.0 | 50.0 | 50.0 | 50.0 |
| $\tilde{k}_{a,2}$ | 2.00 | 1.0e-02 | 1.0e+02 | 2.00 | 2.00 | 2.40 | 2.00 |
| $\tilde{k}_{d,2}$ | 10.0 | – | – | 10.0 | 10.0 | 12.0 | 10.0 |
| $\tilde{k}_{eq,2}$ | 0.20 | 1.0e-02 | 1.0e+02 | 0.20 | 0.20 | 0.20 | 0.20 |
| $v_2$ | 5.00 | 1.00 | 50.0 | 5.00 | 5.00 | 5.00 | 5.00 |
| $\sigma_2$ | 50.0 | 1.00 | 100 | 50.0 | 50.0 | 49.8 | 50.0 |
| NRMSD | | | | 4.2e-06 | 4.2e-06 | 8.7e-04 | 4.2e-06 |
| $\bar{S}$ | | | | 4.5e-05 | 4.5e-05 | 4.6e-04 | 4.5e-05 |
| Wall Time | | | | 2 days: 9:38:44 | 2 days: 3:40:50 | 4 days: 15:23:18 | 5 days: 2:19:39 |

**Table 8**

Fixed parameters for experimental case study.

| Parameter | Value | Unit |
|---|---|---|
| $Q$ | 8.33e-08 | $m^3/s$ |
| $A_c$ | 2.01e-04 | $m^2$ |
| $L_c$ | 2.50e-01 | $m$ |
| $r_p$ | 4.50e-05 | $m$ |
| $L_t$ | 1.46e-01 | $m$ |
| $\Lambda$ | 2.23e+00 | $mol/m^3$ |

**Table 9**

Dextran pulse experiment with detached column. Parameter bounds, estimated parameters and performance indicators.

| | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|
| | Lower | Upper | $S_{Front}$ | $S_{SSD}$ | $S_{Front}$ | $S_{SSD}$ |
| $A_t$ | 1.5e-05 | 2.5e-05 | 1.9e-05 | 1.9e-05 | 1.9e-05 | 1.9e-05 |
| $D_t$ | 1.0e-09 | 1.0e-05 | 2.7e-06 | 2.4e-06 | 2.5e-06 | 2.4e-06 |
| NRMSD | | | 3.8e-02 | 3.8e-02 | 3.8e-02 | 3.8e-02 |
| $\bar{S}$ | | | 4.6e-04 | 4.2e-04 | 4.1e-04 | 4.2e-04 |
| Wall Time | | | 0:02:43 | 0:01:07 | 0:01:10 | 0:00:55 |
| Results | | | 2 | | 1 | |

the existing optima by at least 1%. The mean of the involved metrics, $\bar{S}$, is used to select the final result. The model parameters in Table 8 are determined independently of the staged parameter estimation procedure. Parameter transformations are applied analogously to the synthetic case study. Reference concentrations of the SMA model are $c_r^s = 225\ mol/m^3$ and $c_r^p = 450\ mol/m^3$.

### 11.1. Dextran pulses with detached column

A DPFR is used to describe the impact of the tubing and other holdup volumes that are external to the chromatography column. For this dataset, more complex models with multiple CSTR and DPFR units were unable to better reproduce the observed behavior but suffered from high parameter correlations and poor identifiability (data not shown). Hence, the external holdup volumes are lumped, and the DPFR model parameters are not meant to reflect the real dimensions of the tubing. It was further observed that the given tubing length, $L_t$, can be used for the DPFR without deteriorating the match between model and data. In the first parameter estimation stage, the cross-section area, $A_t$, and dispersion, $D_t$, of the DPFR are estimated from a dextran pulse experiment with detached column, i.e. with the column replaced by a zero-volume connector. The $S_{Front}$ score accounts for the non-ideal behavior of dextran tracer. The model is simultaneously fitted to two experimental replicates by combining the scores as described in Section 9. The time intervals $J$ are separately determined for both chromatograms. They range from 19 to 32 s for the first experiment and from 15 to 32 s for the second (Fig. S9). As in the synthetic case study, the $S_{SSD}$ score is applied to the same time intervals for comparison.

Results are shown in Table 7 and Fig. 9. Obviously, the quality of this real-world data is much worse than in the synthetic case study. This is reflected in larger NRMSD values of about 4% of the highest concentration in the selected time interval. Runtimes vary between one and three minutes. The algorithm and score combinations reach the same NRMSD and terminate at the same estimated areas, $A_t$, with two digits precision. The dispersion estimates de-

viate 8% or less from the mean of the compared values, with the largest deviation for $S_{Front}$ with the GA. It is important to understand that these differences do not allow to assess the certainty of the estimates. Measurement errors clearly exist and are propagated to the estimated parameters. However, different methods are required for quantifying their nature and impact.

Interestingly, the GA has found conflicts between the metrics in the $S_{Front}$ score, even though they are designed to be complementary to each other. Two different parameter sets are optimal in the Pareto sense, i.e. one metric can only be improved at the cost of deteriorating another. In this example, multiple optima are likely caused by relatively large errors in the experimental data. However, Pareto optima should always be carefully analyzed, as they indicate inconsistencies in the data that would remain undetected with the $S_{SSD}$ score. In Table 9, the selection of a final result is based on the mean, $\bar{S}$, of the involved metrics. The parameters, metrics and chromatograms of both Pareto optima are shown in Table S5. The parameter estimates are identical within two digits precision, and differences in the corresponding chromatograms are hardly visible to the human eye without magnification. The mean of the metrics is ca. 3x smaller for the selected optimum, while the NRMSD ca. 8 times larger. However, the NRMSD can be misleading as exemplified in Section 7.1.

### 11.2. Dextran pulses with attached column

In the second stage, column porosity, $\varepsilon_c$, and column dispersion, $D_{ax}$, are estimated from a dextran pulse experiment with attached column. The model is fitted to two experimental replicates following the same procedure as above. The time intervals range from 52 to 218 s for the first experiment and from 24 to 218 s for the second (Fig. S10). Results are shown in Table 10 and Fig. 10. The NRMSD is similar and well below 1% for all algorithm and score combinations. Runtimes vary between about two and nine minutes with an advantage for the GA. The NRMSD values are slightly bet-
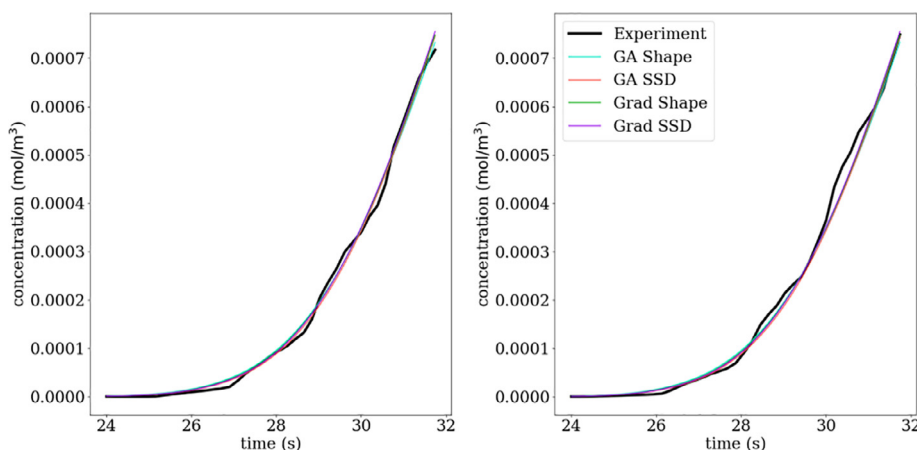
**Fig. 9.** Two replicates of dextran pulse experiment with detached column. Experimental data and estimated chromatograms.
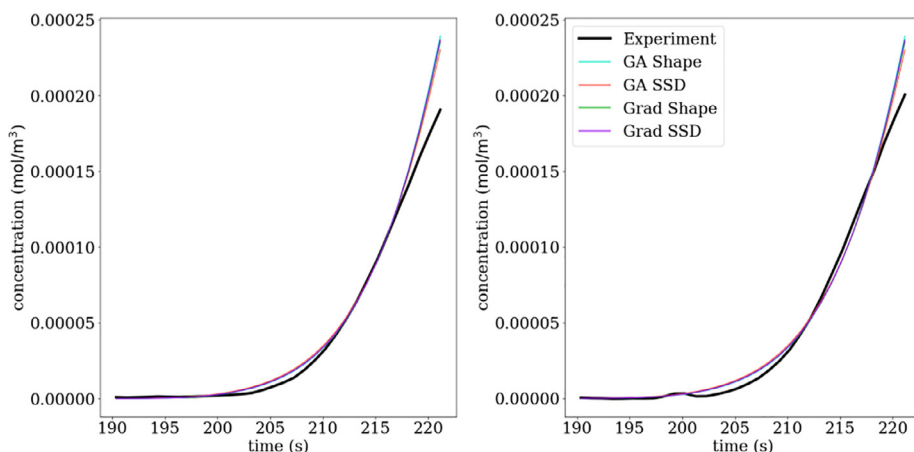


**Fig. 10.** Two replicates of dextran pulse experiment with attached column. Experimental data and estimated chromatograms.

**Table 10**
Dextran pulse experiment with attached column. Parameter bounds, estimated parameters and performance indicators.

| | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|
| | Lower | Upper | $S_{Front}$ | $S_{SSD}$ | $S_{Front}$ | $S_{SSD}$ |
| $\varepsilon_c$ | 2.0e-01 | 4.0e-01 | 3.3e-01 | 3.4e-01 | 3.3e-01 | 3.3e-01 |
| $D_c$ | 1.0e-12 | 1.0e-05 | 3.2e-07 | 3.8e-07 | 3.4e-07 | 3.4e-07 |
| NRMSD | | | 7.0e-03 | 7.0e-03 | 6.6e-03 | 6.6e-03 |
| $\bar{S}$ | | | 1.6e-03 | 1.6e-03 | 1.4e-03 | 1.5e-03 |
| Wall Time | | | 0:04:29 | 0:01:46 | 0:06:02 | 0:08:34 |
| Results | | | 1 | | 2 | |

**Table 11**
Protein pulse experiment under non-binding conditions. Parameter bounds, estimated parameters and performance indicators.

| | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|
| | Lower | Upper | $S_{Peak}$ | $S_{SSD}$ | $S_{Peak}$ | $S_{SSD}$ |
| $\varepsilon_p$ | 2.0e-01 | 5.0e-01 | 3.3e-01 | 3.3e-01 | 3.3e-01 | 3.3e-01 |
| $k_f$ | 1.0e-12 | 1.0e-05 | 2.6e-06 | 1.3e-06 | 2.0e-06 | 1.3e-06 |
| $D_p$ | 1.0e-12 | 1.0e-05 | 2.2e-11 | 3.1e-11 | 2.5e-11 | 3.0e-11 |
| NRMSD | | | 2.1e-02 | 1.9e-02 | 2.1e-02 | 1.9e-02 |
| $\bar{S}$ | | | 1.1e-02 | 1.6e-02 | 1.1e-02 | 1.5e-02 |
| Wall Time | | | 0:21:16 | 0:06:16 | 0:16:02 | 0:14:08 |
| Results | | | 4 | | 2 | |

ter for gradient search with both scores. The porosity estimates are almost identical, while the Dispersion estimates deviate 10% or less from the mean of the compared values, with larger deviations for the GA but in opposite directions for $S_{Front}$ and $S_{SSD}$. Fig. 10 shows a systematic mismatch between the slopes of model and data towards the right end of the considered time intervals, i.e., the mechanistic model does not fully capture the observed process. This is currently analyzed in detail. Uncertainty analysis and model extensions will be subject of separate publications.

For this data set, gradient search has found two Pareto optima for $S_{Front}$ with details shown in Table S6 and Fig. S11. Here, NRMSD and $\bar{S}$ agree on the best choice. Note that the applied multi-start strategy for gradient search is not specifically designed for finding Pareto optima, as the metrics are combined into one objective for each optimizer run. However, Pareto optima can be found as a side

effect when the results of several runs are compared at the level of individual metrics.

### 11.3. Non-Binding protein pulses

In the third stage, particle porosity, $\varepsilon_p$, film diffusion, $k_f$, and pore diffusion, $D_p$, are estimated from a protein pulse under non-binding conditions, i.e. high salt. Using the target protein instead of salt as non-binding but pore penetrating tracer is more reliable, as salt diffuses faster and has a better pore accessibility. Parameter estimation results are shown in Table 11 and Fig. 11. Runtimes are between six and 22 min for GA and are about 15 min for gradient search. The NRMSD is very similar for the compared methods. The fitted chromatograms are also similar with the largest variations around the peak maximum. The porosity estimates are identical.
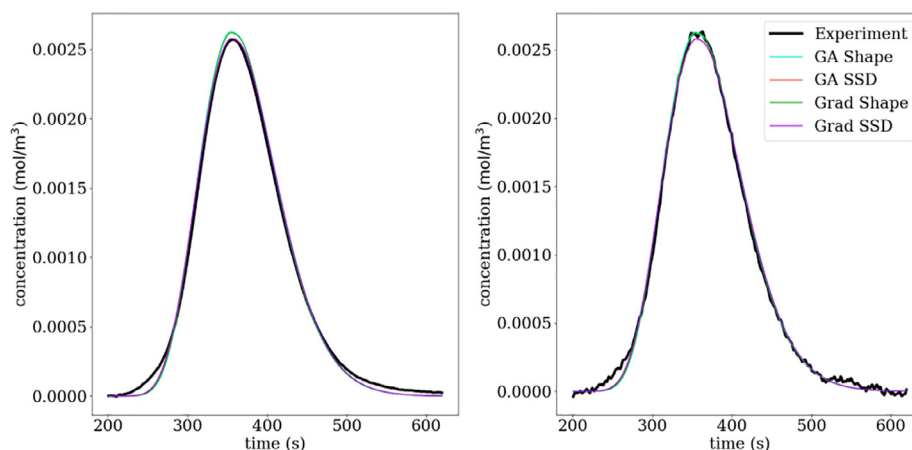
**Fig. 11.** Two replicates of protein pulse experiment under non-binding conditions. Experimental data and estimated chromatograms.

The estimates for particle porosity and pore diffusion show relatively large differences between the scores, partly above 100%, while similar results are obtained by both algorithms. Moreover, both algorithms found multiple Pareto optima for the $S_{Peak}$ score. Table S7 shows the Pareto optima found with GA and $S_{Peak}$. They differ by more than 100% in some parameters, and the lowest mean is in conflict with the lowest NRMSD. Table S8 shows the Pareto optima found with gradient search and $S_{Peak}$. They are closer to each other, and the lowest mean is not in conflict with the lowest NRMSD. Fig. S12 compares the simulated chromatograms of the four Pareto optima of GA and $S_{Peak}$ and the single optimum of gradient search and $S_{SSD}$, which is the standard method for parameter estimation in chromatography. One Pareto optimum matches the peak position in one experiment better and one Pareto optimum matches the peak position in the other experiment better, while the other two Pareto optima have time offsets but describe the peak shape in both experiments better. This inherent conflict is caused by a 0.4 second signal offset between the two experiments which is likely caused by slightly different pump delays in each experiment. As introduced and comprehensively discussed in Section 7, the $S_{SSD}$ cannot handle such inconsistencies in the data, which was a major motivation for designing the new score system. In Table S7, the Pareto optimum with the lowest NRMSD has the largest mean, $\bar{S}$. This Pareto optimum is very similar to single optima found by both algorithms with $S_{SSD}$. Hence, the results of the new score system include the result of the standard method, but the standard method is not able to reveal deficiencies in the model or data. Even when the final optimum is automatically selected, the existence of several Pareto optima indicates potential issues and should generally trigger manual inspection of the model and data. Alternatively, a probability distribution could be passed to the following stage. This is not in the scope of this publication, but future work will address uncertainty propagation in the context of Bayesian optimization.

### 11.4. Gradient elution with partial fractionation

In the fourth and last stage, the parameters of a two component SMA model are estimated from two gradient elution experiments with different loading times and gradient slopes. One experiment has fractionation data available (Fig. S14) while the second features an extended loading time (Fig. S15). Experimental details are described in Section 5. Each fraction is 96 s in length. Precise collection times and pool concentrations are shown in Table S9. The scaled adsorption rates, $\tilde{k}_a$, the scaled equilibrium constant,

$\tilde{k}'_{eq}$, characteristic charge, $\nu$, and shielding coefficient, $\sigma$, are estimated. The adsorption and desorption rates, $\tilde{k}_a$ and $\tilde{k}_d$, are determined from the parameter transformation in Eqs. (26), (27). The ionic capacity, $\Lambda$, was separately determined by titration. The parameters are estimated using the $S^*_{Peak}$ score on the sum chromatogram and the $S^*_{Gauss}$ score on the fractionation data. $S^*_{Peak}$ with six metrics on three peaks and $S^*_{Gauss}$ with three metrics on two peaks and two components results in $1 \cdot 3 \cdot 6 + 2 \cdot 2 \cdot 3 = 30$ objectives. The $S_{SSD}$ score is applied to the sum chromatogram and to the fractionation data of each component.
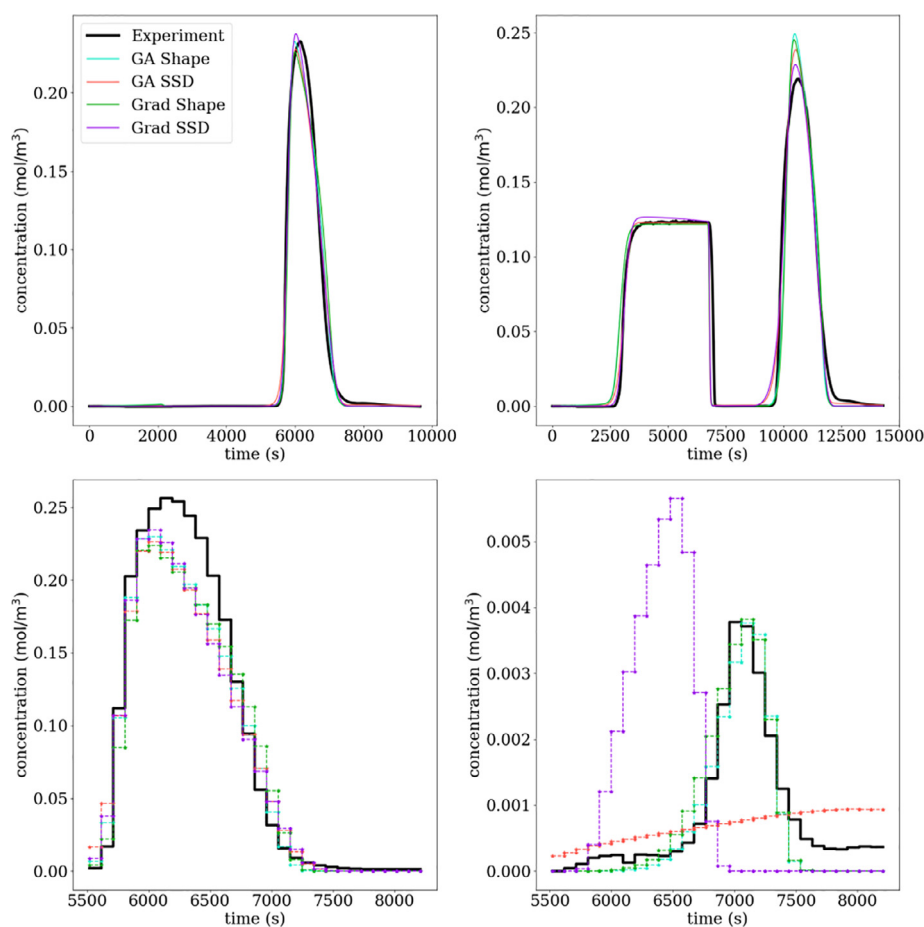
Results are shown in Table 12 and Fig. 12. The compute times of the different algorithm and score combinations range from ca. 1 to 18 days. The GA appears slower than gradient search, but the reported GA runtimes include local refinement by gradient search. For the $S^*_{Peak} \cup S^*_{Gauss}$ score, 3 days of GA runtime were followed by 5 days of local refinement, and for the $S_{SSD}$ score, 5 h of GA runtime were followed by 17 days of local refinement. In both cases, local refinement hardly improved the final result, and the same has been observed in previous stages. With synthetic data local refinement can improve the solution up to the limits of numerical precision due to the model being able to perfectly explain the synthetic data. With experimental data the model can't perfectly explain the data and many small steps are taken to make marginal improvements to the results. The GA can quickly get close enough to an optimal solution that local refinement is left with making small refinements to match an imperfect model to imperfect data. In conclusion, the total runtime can be substantially reduced by skipping the local refinement.

The algorithm and score combinations differ not only in their runtimes but also in the achieved parameter estimates. The NRMSD is relatively low for all algorithm and score combinations. Even though it is slightly smaller for $S_{SSD}$ than for $S^*_{Peak} \cup S^*_{Gauss}$, for both algorithms, some alarming issues are observed for $S_{SSD}$. For GA, the characteristic charge of the second component, $\nu_2$, is at the lower bound, and for gradient search, the shielding factor of the second component, $\sigma_2$, is at the upper bound. These values indicate unrealistic parameter values, as wide search intervals were chosen around typically observed values for monoclonal antibodies. Despite the similar NRMSD, the scaled adsorption constants of the first component, $\tilde{k}_{a,1}$, differ by more than 5 orders of magnitude between the search algorithms, and the same is true for the scaled equilibrium constant of the second component, $\tilde{k}_{eq,2}$. This indicates poor identifiability of the estimated parameters with the $S_{SSD}$ score. Moreover, the corresponding chromatograms do not at

**Table 12**
Gradient elution experiment with partial fractionation. Parameter bounds, estimated parameters and performance indicators.

| | Bounds | | GA | | Gradient | |
|---|---|---|---|---|---|---|
| | Lower | Upper | $S^*_{Peak} \cup S^*_{Gauss}$ | $S_{SSD}$ | $S^*_{Peak} \cup S^*_{Gauss}$ | $S_{SSD}$ |
| $\tilde{k}_{a,1}$ | 1.0e-06 | 1.0e+06 | 5.1e+00 | 9.8e+06 | 1.2e+01 | 7.1e-01 |
| $\tilde{k}_{d,1}$ | – | – | 3.4e+04 | 7.6e+08 | 2.3e+05 | 1.7e+01 |
| $\tilde{k}_{eq,1}$ | 1.0e-06 | 1.0e+06 | 1.5e-04 | 1.3e-02 | 5.4e-05 | 4.1e-02 |
| $v_1$ | 1.0e+00 | 2.0e+02 | 1.5e+01 | 8.2e+00 | 1.7e+01 | 6.7e+00 |
| $\sigma_1$ | 1.0e+00 | 2.0e+02 | 5.0e+01 | 4.4e+01 | 5.1e+01 | 4.1e+01 |
| $\tilde{k}_{a,2}$ | 1.0e-06 | 1.0e+06 | 7.4e-03 | 1.5e+04 | 8.9e-03 | 1.0e+06 |
| $\tilde{k}_{d,2}$ | – | – | 2.6e+03 | 1.6e+03 | 1.0e+03 | 1.0e+12 |
| $\tilde{k}_{eq,2}$ | 1.0e-06 | 1.0e+06 | 2.9e-06 | 9.6e+00 | 8.7e-06 | 1.0e-06 |
| $v_2$ | 1.0e+00 | 2.0e+02 | 2.3e+01 | 1.0e+00 | 2.1e+01 | 2.1e+01 |
| $\sigma_2$ | 1.0e+00 | 2.0e+02 | 9.5e+01 | 1.0e+02 | 3.9e+01 | 2.0e+02 |
| NRMSD | | | 3.9e-01 | 3.7e-01 | 3.9e-01 | 3.4e-01 |
| $\bar{S}$ | | | 9.0e-02 | 1.9e-01 | 9.1e-02 | 1.0e+00 |
| Wall Time | | | 8 days 1:33:00 | 17 days 9:42:53 | 0 days 19:14:13 | 6 days 12:06:01 |
| Results | | | 13 | | 3 | |



**Fig. 12.** Gradient elution experiment with different loading times and gradient slopes. Experimental data and estimated chromatograms. Sum signal of first experiment (top left) and second experiment (top right). Fractionation data of component 1 (bottom left) and component 2 (bottom right) in first experiment.

all match the fractionation data of the second component, as can be seen in Fig. 12.

In stark contrast, the $S^*_{Peak} \cup S^*_{Gauss}$ score was able to guide both search algorithms towards satisfying matches between simulation and experiment for all peaks of both components in both experiments. The corresponding parameter estimates of both search algorithms are relatively similar. Again, rigorous uncertainty analysis is not in the scope of this publication. Moreover, the GA with $S^*_{Peak} \cup S^*_{Gauss}$ result even features an initial breakthrough peak in the first experiment, which is tiny but important, as it indicates saturation of the column.

For the $S^*_{Peak} \bigcup S^*_{Gauss}$ score, 13 Pareto optima were found by the GA, Table S10, and three by gradient search, Table S11. For both search algorithms, the NRMSD of the Pareto optima are in conflict with the respective mean, $\bar{S}$. The GA results are compared in Fig. S13. All Pareto optima, except one, match all peaks of both components in both experiments. These twelve Pareto optima indicate a tradeoff between the two experiments. However, none of

them perfectly matches one experiment or the other. This indicates that the mechanistic model does not fully capture all features of the process. The observed deviations are likely caused by a combination of non-ideal hydrodynamics, complex binding processes and experimental errors. Fig. S13 provides rich information on the Pareto optima that can be related to expert knowledge for selecting the most suitable parameters for specific applications of the calibrated model. Specific applications of the model described here can entail process robustness investigation. Simulations of scenarios with varying process input parameters that are known to vary between lots, such as total protein concentration in the feed solution or ionic strength of elution buffers, can inform decision making in case of process performance deviations. Further applications can range from screening process parameters, such as stop collect criteria, to inform process optimization experiments, to risk-analyses for technical transfer activities.

## 12. Conclusions

A novel score system for estimating chromatography model parameters has been introduced and demonstrated using both synthetic as well as experimental case studies. In contrast to least squares estimation, which is the de facto standard approach and mostly combined with single-objective gradient search, the new score system provides multiple objectives (metrics) that are simultaneously optimized. Typical objectives are the shape, position and height of individual peaks. Even when these objectives are not in conflict, which is typically the case for synthetic data, they can help to improve convergence of the search algorithm, particularly for poor start values and in initial phases of the optimization, by allowing progress in one objective while sacrificing another. However, multiple Pareto optima are often found for experimental data, indicating deficiencies in the model or data. For example, peaks might be slightly shifted by pump delays or small changes in the elution buffer that are unknown and hence not covered by the chromatography model. Least squares penalize position offsets much stronger than peak shape variations and would consequently lead to large errors in the estimated parameters of the binding model. The proposed new score system is designed to handle such situations better by allowing tradeoffs between conflicting objectives. It has been found that the mean of involved metrics provides a better measure for the quality of the match between simulation and experiment than the least squares residual. At first, the diversity and uncertainty that is introduced by multiple Pareto optima might appear as a shortcoming in comparison to the simplicity of the standard approach. However, this uncertainty is not created but only revealed by the presented parameter estimation approach. The ground truth is not changed but made accessible to proper analysis and visualization. However, this cannot replace a rigorous analysis of error sources and uncertainty propagation, which will be subject of another publication. Multiple pareto optima can also be post-processed by weighting different metrics or experiments to select an optimum appropriate to the problem at hand. The advantage of handling this in post-processing is that it does not require a priori selection of weights and it does not negatively impact the parameter estimation process.

Another advantage of the new score system is that multiple experiments and fractionation data can be integrated without needing to weight different objectives. The objectives can be combined into one for a gradient search algorithm, but a multi-objective GA has been observed to be much more robust for complex problems. Pareto optima can indicate inconsistencies between multiple experiments, and the Pareto front carries rich information on potential causes of failure, as demonstrated in the experimental case study. This information is important for understanding the system and designing better experiments. The genetic algorithm as

well as the multi-start strategy for gradient search are parallelized with progress monitoring for computational efficiency on compute clusters but also on multi-core processors in personal computers. Due to the modular nature of the new score system, new metrics can be merged into the existing framework without changing the search algorithm. This can be particularly useful for real-world and large-scale industrial data where increasing system complexity might require consideration of additional features. The current metrics have already been shown to properly address non-ideal tracer retention and pump delays.

CADET-Match has been designed as a monolithic approach for fully automated data processing with minimal human intervention. Hence, pragmatic choices were made for some meta-parameters such as approximation order in the smoothing procedure, concentration thresholds for time interval selection or solver tolerances in the chromatography simulation. Countless numerical tests were performed, and much care has been taken to ensure robustness of the presented algorithm on real-life industrial data. Even though the meta-parameters can be changed, due to the open source nature of the code, this is not recommended unless the implications are fully understood. The open code provides full transparency of the applied procedures and allows to adapt and integrate the software in operational workflows. CADET-Match can be used with all model variants that are available in the parent project, CADET, and cover a wide range of transport and binding models. The scripts for running all case studies in this publication are freely available and can easily be adjusted to other scenarios. Moreover, some aspects and code parts of the presented work are not even specific to chromatography or parameter estimation. For example, the smoothing procedure was recently applied to radioactive tracer signals in plant science.

## Declaration of Competing Interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

## CRediT authorship contribution statement

**William Heymann:** Methodology, Software, Formal analysis, Investigation, Data curation, Writing – original draft, Visualization. **Juliane Glaser:** Validation, Data curation, Writing – review & editing, Supervision. **Fabrice Schlegel:** Conceptualization, Data curation, Writing – review & editing, Supervision. **Will Johnson:** Conceptualization, Resources, Data curation, Writing – review & editing, Funding acquisition. **Pablo Rolandi:** Conceptualization, Resources, Writing – review & editing, Funding acquisition. **Eric von Lieres:** Conceptualization, Resources, Writing – original draft, Writing – review & editing, Supervision, Project administration.

## Supplementary materials

Supplementary material associated with this article can be found, in the online version, at doi:10.1016/j.chroma.2021.462693.

## References

[1] D.H. Wolpert, W.G. Macready, No free lunch theorems for optimization, IEEE Trans. Evol. Comput. 1 (1997) 67–82, doi:10.1109/4235.585893.

[2] J. Schmölder, M. Kaspereit, A modular framework for the modelling and optimization of advanced chromatographic processes, Processes 8 (2020) 65, doi:10.3390/pr8010065.

[3] D. Karlsson, N. Jakobsson, A. Axelsson, B. Nilsson, Model-based optimization of a preparative ion-exchange step for antibody purification, J. Chromatogr. A 1055 (2004) 29–39, doi:10.1016/j.chroma.2004.08.151.

[4] F. Ojala, M. Max-Hansen, D. Kifle, N. Borg, B. Nilsson, Modelling and optimisation of preparative chromatographic purification of europium, J. Chromatogr. A 1220 (2012) 21–25, doi:10.1016/j.chroma.2011.11.028.

[5] O. Khanal, V. Kumar, F. Schlegel, A.M. Lenhoff, Estimating and leveraging protein diffusion on ion-exchange resin surfaces, Proc. Natl. Acad. Sci. U.S.A. 117 (2020) 7004–7010, doi:10.1073/pnas.1921499117.

[6] S. Yamamoto, K. Nakanishi, R. Matsuno, T. Kamijubo, Ion exchange chromatography of proteins?predictions of elution curves and operating conditions. II. Experimental verification, Biotechnol. Bioeng. 25 (1983) 1373–1391, doi:10.1002/bit.260250516.

[7] S. Yamamoto, K. Nakanishi, R. Matsuno, T. Kamikubo, Ion exchange chromatography of proteins?prediction of elution curves and operating conditions. I. Theoretical considerations, Biotechnol. Bioeng. 25 (1983) 1465–1483, doi:10.1002/bit.260250605.

[8] C.A. Brooks, S.M. Cramer, Steric mass-action ion exchange: displacement profiles and induced salt gradients, AIChE J. 38 (1992) 1969–1978, doi:10.1002/aic.690381212.

[9] S.D. Gadam, G. Jayaraman, S.M. Cramer, Characterization of non-linear adsorption properties of dextran-based polyelectrolyte displacers in ion-exchange systems, J. Chromatogr. A 630 (1993) 37–52, doi:10.1016/0021-9673(93)80440-j.

[10] T. Gu, Modeling of affinity chromatography, mathematical modeling and scale-up of liquid chromatography. (1995) 81–94. https://doi.org/10.1007/978-3-642-79541-1_8.

[11] U. Altenhöner, M. Meurer, J. Strube, H. Schmidt-Traub, Parameter estimation for the simulation of liquid chromatography, J. Chromatogr. A 769 (1997) 59–69, doi:10.1016/s0021-9673(97)00173-8.

[12] F. James, M. Sepúlveda, F. Charton, I. Quinónes, G. Guiochon, Determination of binary competitive equilibrium isotherms from the individual chromatographic band profiles, Chem Eng Sci 54 (1999) 1677–1696, doi:10.1016/s0009-2509(98)00539-9.

[13] F. Gritti, W. Piatkowski, G. Guiochon, Study of the mass transfer kinetics in a monolithic column, J. Chromatogr. A 983 (2003) 51–71, doi:10.1016/s0021-9673(02)01648-5.

[14] A. Felinger, A. Cavazzini, G. Guiochon, Numerical determination of the competitive isotherm of enantiomers, J. Chromatogr. A 986 (2003) 207–225, doi:10.1016/s0021-9673(02)01919-2.

[15] D. Karlsson, N. Jakobsson, K.-J. Brink, A. Axelsson, B. Nilsson, Methodologies for model calibration to assist the design of a preparative ion-exchange step for antibody purification, J. Chromatogr. A 1033 (2004) 71–82, doi:10.1016/j.chroma.2003.12.072.

[16] A. Ladiwala, K. Rege, C.M. Breneman, S.M. Cramer, A priori prediction of adsorption isotherm parameters and chromatographic behavior in ion-exchange systems, Proc. Natl. Acad. Sci. U.S.A. 102 (2005) 11710–11715, doi:10.1073/pnas.0408769102.

[17] G. Carta, A.R. Ubiera, T.M. Pabst, Protein mass transfer kinetics in ion exchange media: measurements and interpretations, Chem. Eng. Technol. 28 (2005) 1252–1264, doi:10.1002/ceat.200500122.

[18] P. Forssén, R. Arnell, T. Fornstedt, An improved algorithm for solving inverse problems in liquid chromatography, Comput. Chem. Eng. 30 (2006) 1381–1391, doi:10.1016/j.compchemeng.2006.03.004.

[19] M. Schröder, E. von Lieres, J. Hubbuch, Direct quantification of intraparticle protein diffusion in chromatographic media, J. Phys. Chem. B 110 (2005) 1429–1436, doi:10.1021/jp0542726.

[20] T. Müller-Späth, G. Ströhlein, L. Aumann, H. Kornmann, P. Valax, L. Delegrange, E. Charbaut, G. Baer, A. Lamproye, M. Jöhnck, M. Schulte, M. Morbidelli, Model simulation and experimental verification of a cation-exchange IgG capture step in batch and continuous chromatography, J. Chromatogr. A 1218 (2011) 5195–5204, doi:10.1016/j.chroma.2011.05.103.

[21] A. Osberghaus, S. Hepbildikler, S. Nath, M. Haindl, E. von Lieres, J. Hubbuch, Determination of parameters for the steric mass action model—a comparison between two approaches, J. Chromatogr. A 1233 (2012) 54–65, doi:10.1016/j.chroma.2012.02.004.

[22] Z. Liu, J. Roininen, I. Pulkkinen, T. Sainio, V. Alopaeus, Moment based weighted residual method—New numerical tool for a nonlinear multicomponent chromatographic general rate model, Comput. Chem. Eng. 53 (2013) 153–163, doi:10.1016/j.compchemeng.2013.02.008.

[23] V. Kumar, S. Leweke, E. von Lieres, A.S. Rathore, Mechanistic modeling of ion-exchange process chromatography of charge variants of monoclonal antibody products, J. Chromatogr. A 1426 (2015) 140–153, doi:10.1016/j.chroma.2015.11.062.

[24] T. Gu, G. Iyer, K.-S.C. Cheng, Parameter estimation and rate model simulation of partial breakthrough of bovine serum albumin on a column packed with large Q Sepharose anion-exchange particles, Sep. Purif. Technol. 116 (2013) 319–326, doi:10.1016/j.seppur.2013.06.004.

[25] M. Rüdt, F. Gillet, S. Heege, J. Hitzler, B. Kalbfuss, B. Guélat, Combined Yamamoto approach for simultaneous estimation of adsorption isotherm and kinetic parameters in ion-exchange chromatography, J. Chromatogr. A 1413 (2015) 68–76, doi:10.1016/j.chroma.2015.08.025.

[26] T. Hahn, P. Baumann, T. Huuk, V. Heuveline, J. Hubbuch, UV absorption-based inverse modeling of protein chromatography, Eng. Life Sci. 16 (2015) 99–106, doi:10.1002/elsc.201400247.

[27] T.C. Huuk, T. Hahn, K. Doninger, J. Griesbach, S. Hepbildikler, J. Hubbuch, Modeling of complex antibody elution behavior under high protein load densities in ion exchange chromatography using an asymmetric activity coefficient, Biotechnol. J. 12 (2017) 1600336, doi:10.1002/biot.201600336.

[28] G. Wang, T. Briskot, T. Hahn, P. Baumann, J. Hubbuch, Estimation of adsorption isotherm and mass transfer parameters in protein chromatography using artificial neural networks, J. Chromatogr. A 1487 (2017) 211–217, doi:10.1016/j.chroma.2017.01.068.

[29] V. Kumar, A.M. Lenhoff, Mechanistic modeling of preparative column chromatography for biotherapeutics, Annu. Rev. Chem. Biomol. Eng. 11 (2020) 235–255, doi:10.1146/annurev-chembioeng-102419-125430.

[30] J. Diedrich, W. Heymann, S. Leweke, S. Hunt, R. Todd, C. Kunert, W. Johnson, E. von Lieres, Multi-state steric mass action model and case study on complex high loading behavior of mAb on ion exchange tentacle resin, J. Chromatogr. A 1525 (2017) 60–70, doi:10.1016/j.chroma.2017.09.039.

[31] J. Samuelsson, P. Sajonz, T. Fornstedt, Impact of an error in the column hold-up time for correct adsorption isotherm determination in chromatography, J. Chromatogr. A 1189 (2008) 19–31, doi:10.1016/j.chroma.2007.10.032.

[32] J.W. Cooley, J.W. Tukey, An algorithm for the machine calculation of complex Fourier series, Math. Comput. 19 (1965) 297, doi:10.1090/s0025-5718-1965-0178586-1.

[33] Abraham. Savitzky, M.J.E. Golay, Smoothing and Differentiation of Data by Simplified Least Squares Procedures, Anal. Chem. 36 (1964) 1627–1639, doi:10.1021/ac60214a047.

[34] L.L. Schumaker, P. Dierckx, Curve and surface fitting with splines, Math. Comput. 63 (1994) 427, doi:10.2307/2153590.

[35] S. Butterworth, On the theory of filter amplifiers - experimental wireless - 1930, Experiment. Wireless Wireless Eng. 7 (1930) 536–541.

[36] P. Virtanen, R. Gommers, T.E. Oliphant, M. Haberland, T. Reddy, D. Cournapeau, E. Burovski, P. Peterson, W. Weckesser, J. Bright, S.J. van der Walt, M. Brett, J. Wilson, K.J. Millman, N. Mayorov, A.R.J. Nelson, E. Jones, R. Kern, E. Larson, C.J. Carey, İ. Polat, Y. Feng, E.W. Moore, J. VanderPlas, D. Laxalde, J. Perktold, R. Cimrman, I. Henriksen, E.A. Quintero, C.R. Harris, A.M. Archibald, A.H. Ribeiro, F. Pedregosa, P. van Mulbregt, S. 1.0 Contributors, Author Correction: sciPy 1.0: fundamental algorithms for scientific computing in Python, Nat. Methods 17 (2020) 352, doi:10.1038/s41592-020-0772-5.

[37] K. Pearson VII, Note on regression and inheritance in the case of two parents, Proc. R. Soc. Lond. 58 (1895) 240–242, doi:10.1098/rspl.1895.0041.

[38] H.B. Curry, The method of steepest descent for non-linear minimization problems, Q. Appl. Math. 2 (1944) 258–261, doi:10.1090/qam/10667.

[39] F. Hayes-Roth, Review of, Adaptation in natural and artificial systems by John H. Holland", The U. of Michigan Press, 1975, ACM SIGART Bulletin 15 (1975), doi:10.1145/1216504.1216510.

[40] M.D. McKay, R.J. Beckman, W.J. Conover, Comparison of three methods for selecting values of input variables in the analysis of output from a computer code, Technometrics 21 (1979) 239–245, doi:10.1080/00401706.1979.10489755.

[41] H. Niederreiter, Low-discrepancy and low-dispersion sequences, J. Number Theory 30 (1988) 51–70, doi:10.1016/0022-314x(88)90025-x.

[42] K. Deb, A. Pratap, S. Agarwal, T. Meyarivan, A fast and elitist multiobjective genetic algorithm: NSGA-II, IEEE Trans. Evol. Comput. 6 (2002) 182–197, doi:10.1109/4235.996017.

[43] K. Deb, H. Jain, An Evolutionary Many-Objective Optimization Algorithm Using Reference-Point-Based Nondominated Sorting Approach, Part I: solving Problems With Box Constraints, IEEE Trans. Evol. Comput. 18 (2014) 577–601, doi:10.1109/tevc.2013.2281535.

[44] V. Chankong, Y.Y. Haimes, Multiobjective decision making: theory and methodology, Courier Dover Publicat. (2008).

[45] M.A. Branch, T.F. Coleman, Y. Li, A Subspace, Interior, and conjugate gradient method for large-scale bound-constrained minimization problems, SIAM J. Sci. Comput. 21 (1999) 1–23, doi:10.1137/s1064827595289108.

[46] F.-M. De Rainville, F.-A. Fortin, M.-A. Gardner, M. Parizeau, C. Gagné, DEAP, in: Proceedings of the Fourteenth International Conference on Genetic and Evolutionary Computation Conference Companion - GECCO Companion '12, 2012, doi:10.1145/2330784.2330799.