

Identifying physico-chemical laws from the robotically collected data

Pascal Neumann,^a Liwei Cao,^{b,c} Danilo Russo,^b Vassilios S. Vassiliadis,^b Alexei A. Lapkin^{b,c*}

^aAachener Verfahrenstechnik – Process Systems Engineering, RWTH Aachen University, Aachen, Germany

^bDepartment of Chemical Engineering and Biotechnology, University of Cambridge, Cambridge, CB3 0AS

^cCambridge Centre for Advanced Research and Education in Singapore, CARES Ltd. 1 CREATE Way, CREATE Tower #05-05, 138602 Singapore

Abstract

A mixed-integer nonlinear programming (MINLP) formulation for symbolic regression was proposed to identify physical models from noisy experimental data. The formulation was tested using numerical models and was found to be more efficient than the previous literature example with respect to the number of predictor variables and training data points. The globally optimal search was extended to identify physical models and to cope with noise in the experimental data predictor variable. The methodology was coupled with the collection of experimental data in an automated fashion, and was proven to be successful in identifying the correct physical models describing the relationship between the shear stress and shear rate for both Newtonian and non-Newtonian fluids, and simple kinetic laws of reactions. Future work will focus on addressing the limitations of the formulation presented in this work, by extending it to be able to address larger complex physical models.

Keywords: model identification; chemical process development; symbolic regression; automated model construction; Mixed-Integer Nonlinear Programming; global optimization

1. Introduction

Today we experience rapid development of a new field of chemical science – digital molecular technology (DMT). It is evident by the increasing number of publications in which synthetic and computational chemistry, or materials development are mixed with machine learning (ML), robotics and artificial intelligence (AI), for example in Refs. [1–6]. DMT is promising

* Corresponding author. A. Lapkin email: aal35@cam.ac.uk

to significantly expand the accessible chemical space and to reduce the price of access to new functional molecules and materials. These are to be achieved through enhanced capabilities in molecules and reactions discovery, and in process development and optimization. The key component of the new DMT methods is the increased volume and quality of chemical data obtained both through data mining, computational chemistry tools and robotic experiments, as lack of data renders ML and AI methods inaccurate and not very useful [7]. Here we ask a question, is it possible to make use of the increased availability of experimental data to enhance our capability in inferring physical knowledge from data by means of algorithmic research? This is driven by the desire to develop predictive models of complex chemical processes, which could be used in optimal control. The approach that we seek to develop should not be based on selecting functional forms from a pre-defined set. This has already been demonstrated within DMT, for example, in selecting suitable kinetic expressions within an automated self-optimization system [8]. Our own interest is in the methods that are inherently not restricted to only the known functional forms and are, therefore, potentially capable of discovering new phenomena.

Automation of chemical research has recently emerged as a highly promising broad methodology [9]. Significant benefits of automated systems include the precise control of the operating conditions, and the ability of generating large, reliable and reproducible datasets [10]. Thanks to the advances in computational power, automated systems have been combined with machine learning (ML) techniques and applied in discovery of novel materials [11] and reactions [12], process development [13] and optimization [8]. Automated systems mainly consist of three components: a decision-making algorithm, a physical platform capable of undertaking the desired class of experiments autonomously and an automated inline/online analytical setup [14]. As the ‘brain’ of automated platforms, the decision algorithm plays a vital role. On this basis, ML algorithms based on statistical surrogate modelling approaches have become increasingly popular, having been shown that they provide accurate models for predictions [15].

On the contrary, the data-driven development of interpretable physical models resisted automation for long [16]. The possibility of building data-driven interpretable and generalizable models for complex and not well understood physical systems is important as these models share the similar structure to those based on first principles and can be transferred to analogous systems, whereas surrogate models cannot be easily generalized [17]. This, in the

longer term, can improve the time and resource efficiency for the product discovery and process optimization, especially for the manufacturability and the scale-up, for which a mechanistic understanding is often crucial [18–20].

The field of algorithmic search for physical models is relatively new, but has seen a number of important advances. Recently, Brunton *et al.* introduced a sparse regression approach to discover equations governing the physics in a chaotic Lorenz system, and in a fluid vortex dynamic system [21]. However, this technique is restricted to a pre-defined algebraic model structure, as selected basis functions are linearly combined. Allowing free-form analytical equations, Bongard and Lipson [22] developed a criterion to find meaningful and complex mathematically invariant models by means of the ML method of symbolic regression (SR). Recent application of SR for physical models can be found in civil engineering [23] and material science [24]. Although successfully proven, the proposed SR was based on a heuristic search that could terminate the optimization in local minima solutions, producing potentially less suitable models than possible. Additionally, as the structure of a model reflects the actual mechanistic interactions within the system studied, these approximate solutions cannot be used reliably to infer any mechanistic information about the system, i.e. to use it to identify chemical reactions mechanisms with certainty. Acknowledging this disadvantage, SR is formulated as a mixed-integer nonlinear programming (MINLP) in Ref. [25,26] and solved to global optimum. However, the method remained in the mathematical domain and has yet to be applied to physical models and noisy experimental data.

This paper aims to advance the method of globally optimal symbolic regression towards automated, data-driven identification of physical models, and its applications to chemical engineering case studies. Compared to additive models in conventional regression and heuristic searches in SR, the globally optimal data-driven modelling technique, without any previously imposed model structure, is expected to discover true underlying relationships more reliably. To accomplish this, a modified optimization formulation of SR is developed and implemented in combination with a framework for physical model selection. As a proof of concept, several case studies were investigated in the areas of rheology and reaction engineering. The purpose is to illustrate an automated research pipeline deriving interpretable and generalizable models and thereby providing access to physical knowledge from data generated in robotic experiments. Within this big picture, closing the loop of utilizing the obtained physical models

in further experimentation and generation of physical knowledge by (automated) interpretation remains for future research, as shown in Figure 1.

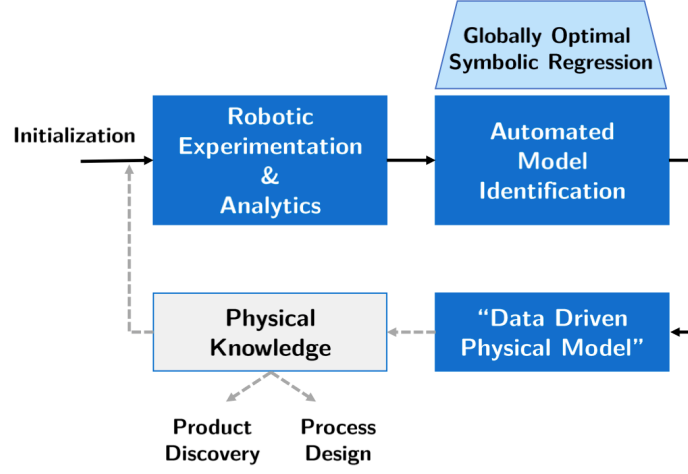


Figure 1. Schematic diagram for automated physical model identification.

2. Materials and Methods

2.1. MINLP formulation

Use will be made of the Directed Acyclic Graph (DAG) description of algebraic functions throughout this work [27]. The MINLP formulation is based on a balanced binary tree superstructure for the representation of the equations describing a physical model. The overall goal is to enable the assembly of free-form algebraic functions by connecting predictor variables and operators, such that the resulting function predicts the dependent variable values accurately. As an example, the structure with nodes in a four-layer balanced binary tree is shown in Figure 2a.

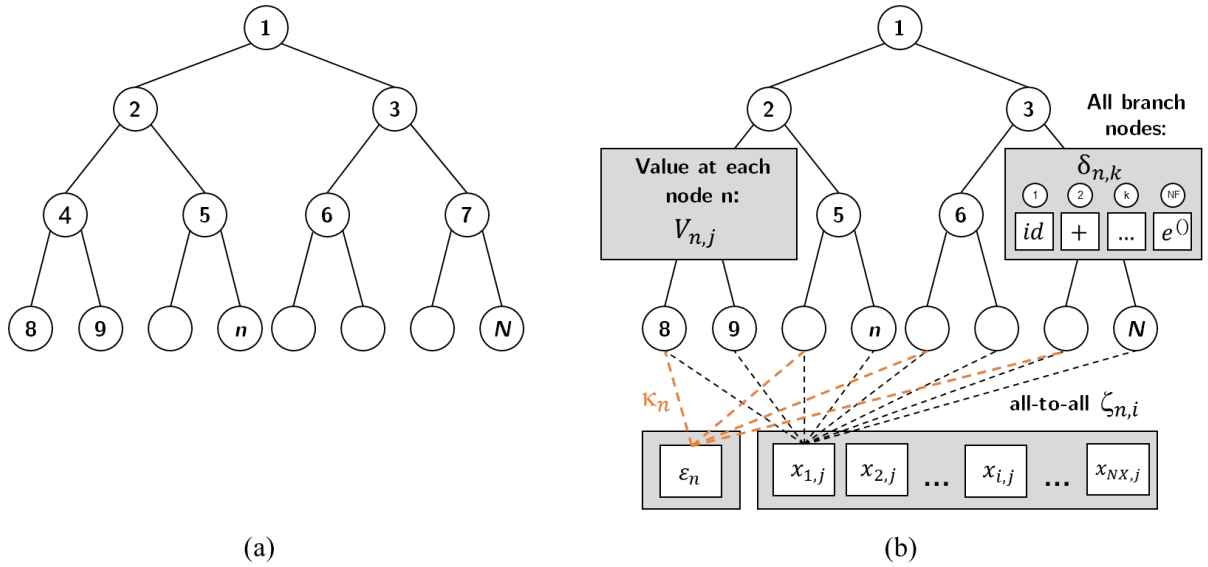


Figure 2. An example of a directed acyclic graph. (a) Binary rooted tree, (b) MINLP set-up in connection with the expression tree.

The formulation presented here is based on Ref [26], but follows a different concept in the set-up of the binary tree in order to reduce the number of binary variables in the global optimization. These modifications are addressed in Section 2.2.

An expression tree consists of $N = 2^{NL} - 1$ nodes, where NL defines the number of layers. All nodes that have a connection with nodes on a lower level, their children nodes, are called branch-nodes (\mathcal{N}_b) or non-leaf nodes and house a mathematical operator. The nodes on the lowest layer in the tree, referred to as leaf-nodes (\mathcal{N}_l), are assigned to a predictor variable $x_{i,j}$ or a constant ϵ_n . In the following sections, we will refer to the total number of activated nodes in the tree as “complexity” of the model (except the ones with an identity operator). Each given data point deployed in the SR is described by two parameters: the value of the predictor variables $x_{i,j}$ and the dependent variable value y_j , which is predicted by the model for each training data point j . As shown in Figure 2b, the input variables are assigned for selection only at the leaf-nodes, while the dependent variable values are used at the root node ($n = 1$) for comparison with the model prediction.

Each node has a value for each data point $V_{n,j}$, which is computed to be used as operator arguments on the layer above. The nodal values at the bottom of the tree are determined by the selection of an input variable or a constant. On branch node layers, nodal values are specified by the selected operator in combination with the node values of their children. The allocation of the input predictor variables $x_{i,j}$ is implemented by the binary variables $\zeta_{n,i}$. Continuous decision variables ϵ_n with bounds ϵ_n^{lo} and ϵ_n^{up} are designated for constants at every even leaf-node. To decide between a variable input and the selection of a constant at the even leaf-nodes, further binaries variables κ_n are assigned. With regard to the branch-nodes, there are binary variables $\delta_{n,k}$ assigned for operator selection, where an operator is active at node n if $\delta_{n,k} = 1$ and inactive if $\delta_{n,k} = 0$. If active, each binary operator is applied using both children nodes while a unary operator uses only the value of the left node $V_{2n,j}$. Five binary operators (addition, subtraction, multiplication, division and power law) and three unary operators (identity, exponential and square root) were implemented. It should be noted that the list of

operators can be extended further, such as cubic, square or logarithm operations as proven in Ref. [26].

Table 1. MINLP Notation: Set

Description	Index	Set	Value
Nodes	N	\mathcal{N}	$\{1, \dots, N\}$
Branch-nodes		\mathcal{N}_b	$\{1, \dots, 2^{NL-1} - 1\}$
Leaf-nodes		\mathcal{N}_l	$\{2^{NL-1}, \dots, N\}$
Even leaf-nodes		\mathcal{N}_l^*	$\{2^{NL-1}, 2^{NL-1} + 2, \dots, N\}$
Algebraic	K	\mathcal{F}	$\{+, -, \dots\}$
Operators			
Predictor	I	\mathcal{X}	$\{1, \dots, NX\}$
Variables			
Input Data Points	J	\mathcal{J}	$\{1, \dots, NE\}$

Table 2. MINLP notation: Parameters.

Description	Parameter
Input variable values	$x_{i,j}$ $i \in \mathcal{X}, j \in \mathcal{J}$
Dependent variable values	y_j $j \in \mathcal{J}$

Table 3. MINLP notation: Decision variables.

Applicability	Description	Variable	Bounds
General	Nodal values	$V_{n,j}$ $n \in \mathcal{N}_b, j \in \mathcal{J}$	$[V_{n,j}^{lo}, V_{n,j}^{up}] \in \mathbb{R}$
Leaf-nodes	Variable selection	$\zeta_{n,i}$ $n \in \mathcal{N}_l, i \in \mathcal{X}$	$\{0,1\}$
	Constant selection	κ_n $n \in \mathcal{N}_l^*$	$\{0,1\}$
	Value of constants	ϵ_n $n \in \mathcal{N}_l^*$	$[\epsilon_{n,j}^{lo}, \epsilon_{n,j}^{up}] \in \mathbb{R}$
	Operator selection	$\delta_{n,k}$ $n \in \mathcal{N}_b, k \in \mathcal{F}$	$\{0,1\}$

Consequently, by using the tree structure and the index assignment given (Tables 1-3), the optimization problem was formulated with the objective to minimize the sum of squared errors

(SSE) between the values computed by the model and the experimental data, Eq. 1, according to Ref. [26].

$$\min \sum_{j=1}^{NE} (y_i - V_{1,j})^2 \quad (1)$$

Eqs. 2-4 enable the operator selection at branch-nodes with a big-M approach, as proposed in Refs. [26] and [28]. The idea of the big-M facilitates the transformation of the disjunctive choice between the operators into linear constraints [28]. If no operator is selected, its nodal value is set to zero by constraints 5-6. Additionally, either none or one operator can be selected at the branch-nodes. Hence, the sum of operator binaries must be less or equal to one, which is constrained by Eq. 7.

In contrast to the branch-node values, the values at the leaf-nodes are determined by equality constraints including the binary selection of predictor variables or constants, Eqs. 8-9. Also, Eqs. 10-11 make sure that either no operand, one variable or one constant can be assigned. Overall, the model should include at least one predictor variable, which is ensured by Eq. 12. For the purpose of completeness, Eqs. 13-17 depict the bounds on decision variables of the MINLP [26].

$$f_k(V_{n,j}, V_{2n,j}, V_{2n+1,j}) \leq M_{n,j,k}^{up} (1 - \delta_{n,k}), n \in \mathcal{N}_b, k \in \mathcal{F}, j \in \mathcal{T} \quad (2)$$

$$f_k(V_{n,j}, V_{2n,j}, V_{2n+1,j}) \geq M_{n,j,k}^{lo} (1 - \delta_{n,k}), n \in \mathcal{N}_b, k \in \mathcal{F}, j \in \mathcal{T} \quad (3)$$

$$g_k(V_{n,j}, V_{2n,j}, V_{2n+1,j}) \leq G_{n,j,k}^{up} (1 - \delta_{n,k}), n \in \mathcal{N}_b, k \in \mathcal{F}, j \in \mathcal{T} \quad (4)$$

$$V_{n,j} \leq V_{n,j}^{up} \sum_{k \in \mathcal{F}} \delta_{n,k}, n \in \mathcal{N}_b, j \in \mathcal{T} \quad (5)$$

$$V_{n,j} \geq V_{n,j}^{lo} \sum_{k \in \mathcal{F}} \delta_{n,k}, n \in \mathcal{N}_b, j \in \mathcal{T} \quad (6)$$

$$\sum_{k \in \mathcal{F}} \delta_{n,k} \leq 1, n \in \mathcal{N}_b \quad (7)$$

$$V_{n,j} = \sum_{i \in \mathcal{X}} \zeta_{n,i} x_{i,j} + \kappa_n \epsilon_n, n \in \mathcal{N}_l^*, j \in \mathcal{T} \quad (8)$$

$$V_{n,j} = \sum_{i \in \mathcal{X}} \zeta_{n,i} x_{i,j}, n \in \mathcal{N}_l \setminus \mathcal{N}_l^*, j \in \mathcal{T} \quad (9)$$

$$\sum_{i \in \mathcal{X}} \zeta_{n,i} + \kappa_n \leq 1, n \in \mathcal{N}_l^* \quad (10)$$

$$\sum_{i \in \mathcal{X}} \zeta_{n,i} \leq 1, n \in \mathcal{N}_l \setminus \mathcal{N}_l^* \quad (11)$$

$$\sum_{n \in \mathcal{N}_l} \sum_{i \in \mathcal{X}} \zeta_{n,i} \geq 1 \quad (12)$$

$$\delta_{n,k} \in \{0,1\}, n \in \mathcal{N}_b, k \in \mathcal{F} \quad (13)$$

$$\zeta_{n,i} \in \{0,1\}, n \in \mathcal{N}_l, i \in \mathcal{X} \quad (14)$$

$$\kappa_n \in \{0,1\}, n \in \mathcal{N}_l^* \quad (15)$$

$$V_{n,j} \in [V_{n,j}^{lo}, V_{n,j}^{up}], n \in \mathcal{N} \quad (16)$$

$$\epsilon_{n,j} \in [\epsilon_{n,j}^{lo}, \epsilon_{n,j}^{up}], n \in \mathcal{N}_l^* \quad (17)$$

Due to the binary architecture and the commutative nature of addition and multiplication, the expression tree contains many mathematically invariant models (symmetries). The design of the formulation should therefore impede redundancies. Eqs. 18-19 resemble cuts in the tree such that, if a unary operator is selected, the unused part towards the right children node is set to zero [26]. The Eqs. 20-23 assure that if an operator is selected on a lower layer of the expression tree, there is an operator attached to the parental node [26]. Likewise, it ensures that the children of a node with value zero also have no operator or variables attached.

Additionally, symmetry breaking cuts (SC) to remove redundant solutions, which are caused by the commutative nature of addition and multiplication, were implemented. Eq. 24 is sufficient for one data point $j = j'$ to impose an order on the values of the children nodes [26]. The symmetry breaking cuts also posed as big-M constraints where M_n^* is set using interval arithmetic on the bounds of the two children node values [28].

$$\sum_{k \in \mathcal{F}_{unary}} \delta_{n,k} \leq 1 - \sum_{k \in \mathcal{F}} \delta_{2n+1,k}, n \in \{1, \dots, 2^{NL-2} - 1\} \quad (18)$$

$$\sum_{k \in \mathcal{F}_{unary}} \delta_{n,k} \leq 1 - \sum_{i \in \mathcal{X}} \zeta_{2n+1,i}, n \in \{2^{NL-2}, \dots, 2^{NL-1} - 1\} \quad (19)$$

$$\sum_{k \in \mathcal{F}} \delta_{n,k} \geq \sum_{k \in \mathcal{F}} \delta_{2n,k}, n \in \{1, \dots, 2^{NL-2} - 1\} \quad (20)$$

$$\sum_{k \in \mathcal{F}} \delta_{n,k} \geq \sum_{k \in \mathcal{F}} \delta_{2n+1,k}, n \in \{1, \dots, 2^{NL-2} - 1\} \quad (21)$$

$$\sum_{k \in \mathcal{F}} \delta_{n,k} \geq \sum_{i \in \mathcal{X}} \zeta_{2n,i} + \kappa_n, n \in \{2^{NL-2}, \dots, 2^{NL-1} - 1\} \quad (22)$$

$$\sum_{k \in \mathcal{F}} \delta_{n,k} \geq \sum_{i \in \mathcal{X}} \zeta_{2n+1,i}, n \in \{2^{NL-2}, \dots, 2^{NL-1} - 1\} \quad (23)$$

$$V_{2n,j'} - V_{2n+1,j'} \geq M_{n,j'}^{SC} (1 - \sum_{k=\{+,*\}} \delta_{n,k}), n \in \mathcal{N}_b, \exists j' \in \mathcal{J} \quad (24)$$

Without a doubt and from experience with big-M formulations in the Mathematical Programming community, this will lead to rather loose lower bounding in the associated Branch-and-Bound (B&B) traditionally used to solve Mixed-Integer Linear (MILP) and Mixed-Integer Nonlinear Programming (MINLP) problems. Indeed, such were the observations reported later in the computational results of this work, and hence the serious limitations that is a challenge for future development of this rigorous methodology.

2.2. Details of Modifications of the Formulation

The previously reported MINLP formulation [26] was modified in order to reduce the number of required binary variables. This is expected to advance the overall efficiency in solving the MINLP as less decision variables have to be determined in the global optimization. The main difference is that in the presented formulation the tree is always fully constructed for a given number of pre-defined tree layers NL . The predictor variables and constants can only be assigned to the lowest layer. The inclusion of an identity function allows to pass up the values without any change to a higher layer in the expression trees.

In comparison to that, in Ref. [26] predictor variables and constants are available for selection at every node in the expression tree superseding an identity function. If those leaf-node operands are selected on a higher layer, their children nodes as well as the other subsequent lower levels are discarded. Hence, the tree is not set up necessarily to the maximum of allowed layers. By introducing an identity function in the modified version of the MINLP, as was proposed with the theoretical formulation of problem in the first recorded publication in the open literature on the topic in chemical engineering [25], the binaries on branch nodes for $x_{i,j}$ and ε_n can be spared. The significance in the reduction is increasing with the number of layers and the number of overall considered input variables. Furthermore, the full construction of the tree allowed to replace big-M constraints at the leaf-nodes by equality constraints due to linearity in the binary variables. Another main difference is the asymmetric supply of continuous variables as constants. This design further reduces the number of required binaries as well as the symmetries in the superstructure set-up. In case of a four-layer tree this would reduce their number by four.

2.3. Solver

For the aims of this work, the choice of solvers is limited to those that can deterministically solve MINLPs to global optimum. According to Ref. [29], the general list of feasible non-convex global MINLP solvers contains ANTIGONE [30], BARON [31], Couenne [32], LindoGlobal [33], and SCIP [34]. According to the results of the comparative solver study [26], BARON solves more SR problems and converges faster than all other solvers. Hence, BARON (version 18.5.9) as commercial general-purpose solver in deterministic global optimization was selected in combination with IBM CPLEX as sub-solver. Upon solver completion, the optimization results are analyzed within Pyomo allowing to translate the optimal decision variable values into the corresponding algebraic equation. This model can

then be evaluated at different inputs for prediction as well. For further simplification, Python's library for symbolic mathematics SymPy was used.

2.4. Physical Model Selection

The SR is to be performed with experimental systems data and it is to be acknowledged that all measurements have an error. Following a globally optimal approach targeting exclusively model accuracy (Eq. 1), errors are represented in the final model what is also known as overfitting. Hence, methodological measures have to be included to restrict the influence of errors on the final model and assure generalisation capabilities with low errors to unseen data. In case of SR, the errors in the training data are propagated through the expression tree and the selected operators apply to the data including errors. The limited robustness to noise is especially prevalent among SR due to its maximal flexibility in constructing free-form models [35].

To only extract the relevant terms describing the main signal and to preferably exclude the errors, the complexity of the final model is penalized. Model complexity is restricted to a threshold C . The identity function does not add to the complexity of a model. Consequently, the true complexity must be discounted by nodes with an identity function assigned. This complexity criterion is included as additional constraint in the MINLP (Eq. 25).

$$\sum_{n \in \mathcal{N}_b} \sum_{k \in \mathcal{F}\{\text{id}\}} \delta_{n,k} + \sum_{n \in \mathcal{N}_l} \sum_{i \in \mathcal{X}} \zeta_{n,i} + \sum_{n \in \mathcal{N}_l^*} \kappa_n \leq C \quad (25)$$

By limiting the flexibility allowed, overfitting can be reduced and sparse models found. This also increases interpretability. Furthermore, this constraint filters mathematical invariants including more terms from the search space.

Next, it is proposed to identify a portfolio of the most accurate models with varying complexity C by solving multiple MINLPs in parallel to global optimality. Among the results, one model is selected that is as sparse as possible to allow interpretation and knowledge extraction but is also as complex as necessary to describe the underlying physical system without overfitting. Hence, the portfolio models are to be compared with regard to validation error. Due to the growing flexibility, the training error is assumed to be the lowest for the model with the highest allowed complexity. Without requiring assumptions about the underlying true model, these can

be compared quantitatively by a data set for validation to check for overfitting. With the purpose of also comparing extrapolation capabilities of the models, the validation data set is created by extracting the data points at extrema of the predictor variables. Finally, the model selection can be based on lowest validation error which also determines the required model complexity. The framework is illustrated in Figure 3.

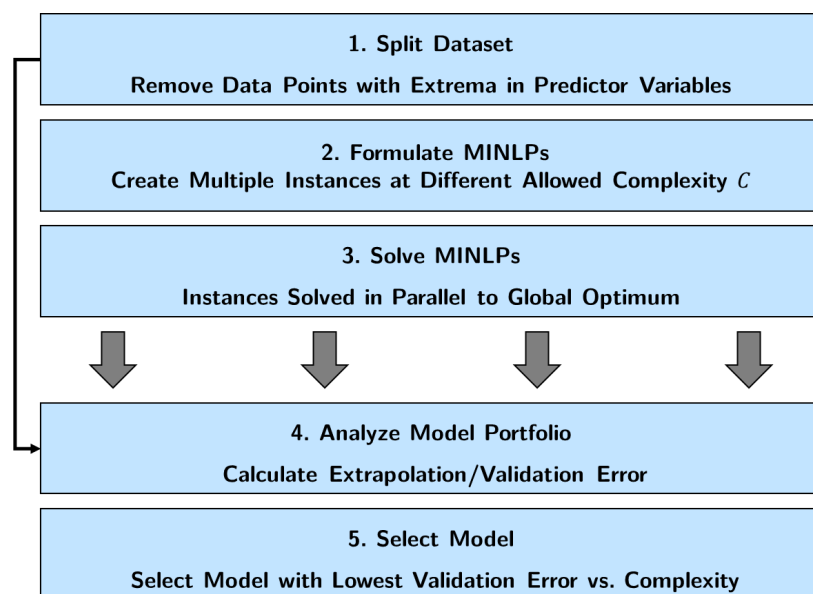


Figure 3. Framework for the Automated Identification and Selection of Physical Models via Symbolic Regression.

2.5. Materials

All chemicals (glycerol $\geq 99.5\%$, isopropyl alcohol $\geq 99.5\%$, carboxymethyl cellulose, 4-nitrophenyl acetate (esterase substrate), 4-nitrophenol $\geq 99\%$, potassium chloride $\geq 99.0\%$, sodium hydroxide $\geq 98.0\%$, methanol $\geq 99.9\%$) were purchased from Sigma Aldrich and used as received. Water was obtained using a Maxima (USF) Milli-Q system. Viscosity experiments on a Newtonian fluid were carried out using a commercially available emulsion.

2.6. Experimental Procedures

Automated viscosity measurement

Samples of aqueous solutions were prepared as a batch of 14 samples in an automated fashion by means of two R-Series pumping modules (Vapourtec Ltd.) coupled with a Gilson fraction collector, and controlled by a custom-written algorithm. Carboxymethyl cellulose was pre-diluted in water.

The viscosities of aqueous glycerol and carboxymethyl cellulose solutions were measured by means of the custom-built automated capillary viscometer shown schematically in Figure 4.

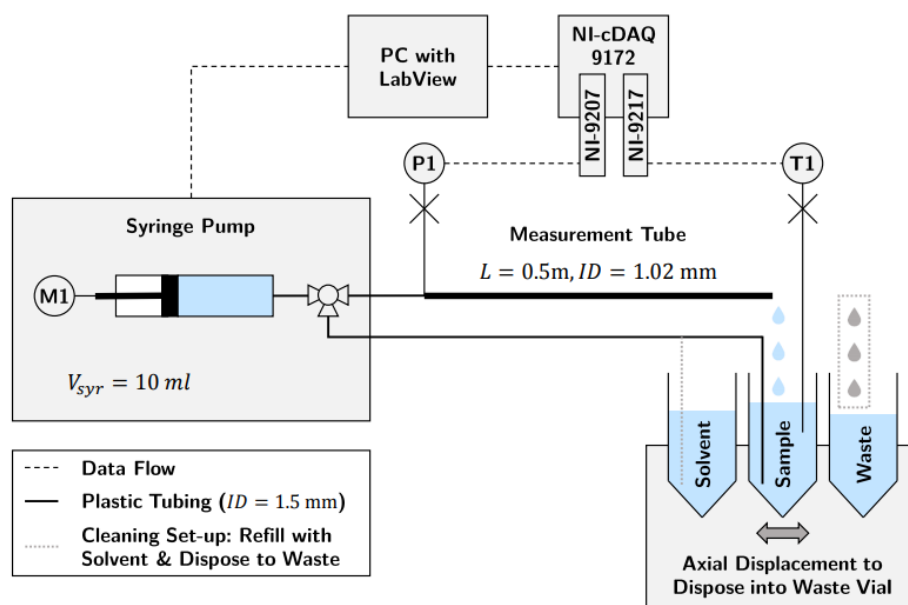


Figure. 4. Schematic diagram of the custom-built capillary viscometer. Solvent withdrawal and waste disposal for cleaning are indicated by dashed line. Vials are displaced axially for waste disposal.

To automatically control the system, a user interface was implemented in LabView 2014 (SP1), that allowed to visualize, analyse, and save the acquired signals. For the non-Newtonian solutions, a special mode was established to test multiple predefined flow rates in an automated sequence. A syringe pump (CETONI neMESYS 290N) was selected to ensure accurate and constant flows free from pulsation, and equipped with a 10 mL glass syringe (CETONI, ID 14.57 mm). The pump was connected to an automated three-ports valve (CETONI valve, max. pressure 3 bar) to empty and re-fill the syringe using different routes of fluid flow. When the syringe was emptied, the valve directed the fluid through the stainless-steel tube (L = 0.5 m, ID = 1.02 mm, OD = 1.59 mm), placed horizontally at the same height as the syringe. The pressure drop relative to ambient pressure was measured using a microfluidic sensor (P1, Elveflow MPS3) with a range of 2 bar and accuracy of ± 4 mbar. The pressure sensor signal was acquired via a voltage and current input module (National Instruments NI-9207). The temperature of the samples was constantly recorded by a resistance temperature detector (T1, Bearing Sensor Platinum, 100 Ω) and was found to be always in the range 21 ± 0.8 °C. The inlet line was automatically switched from the samples to the solvent, 2-propanol, and waste reservoirs for the cleaning and drying routines to prevent cross-contamination with the next

sample. Each measurement was carried out in triplicate. The acquired time series data were processed to derive time-independent data series. This was done using an algorithmic scheme in Python to average only the steady-state pressure-drop signal values. The obtained flow rate (Q) – pressure drop (Δp) data can be used to calculate viscosity according to Hagen-Poiseuille law for Newtonian fluids (Eq. 25), and the Weissenberg-Rabinowitsch equation for non-Newtonian solutions (Eq. 26).

$$\eta = \frac{\pi R^4 \Delta p}{8 Q L} \quad (25)$$

$$\eta = \frac{\pi R^4 \Delta p}{2 Q L} \left(3 + \frac{d \ln Q}{d \ln \Delta p} \right)^{-1} \quad (26)$$

where R and L are the radius and the length of the adopted tubing, respectively.

Before taking advantage of the automated set-up a calibration procedure was conducted to reduce the systematic measurements errors. This consisted of: (i) a pressure sensor calibration carried out using a digital pressure indicator (Druck, DPI 600/IS) applying a static pressure with air, and (ii) a full set-up calibration with aqueous glycerol solutions at known concentrations. In the latter case viscosity of the same samples was also measured by a commercial rotational viscometer (Rheometric Scientific, ARES G2) in order to rule out the systematic error introduced by the automated set-up. Additional information is provided in the Electronic Supplementary Information (ESI).

Reaction kinetic experiments

Reaction kinetic data was collected for hydrolysis of 4-nitrophenyl acetate (PNPA) under basic conditions as a case study. For each experimental run, three stock solutions were prepared consisting of PNPA (at the desired concentration) in 3.0 % (v/v) aqueous methanol, 3 mol·L⁻¹ KCl, and an aqueous NaOH solution at a fixed pH. The adopted conditions for each kinetic experiment are provided in the ESI. The 1 mL of each solution was directly mixed in a spectrophotometric agitated disposable cuvette. Absorption spectra (300-500 nm) were collected at fixed time intervals (Agilent, Cary 60). Absorption data at 400 nm were converted to PNP concentration. Calibration was carried out using different aqueous solutions at a known concentration of PNP at the same methanol and KCl concentrations and pH of the tested solutions. PNPA concentrations were calculated as its initial concentration minus the

concentration of the formed product according to the literature [36,37], since no by-products formation was reported under the adopted conditions.

3. Results and Discussion

3.1. Proof of concept

In the first instance, the methodology described in Section 2.1. was applied to data without errors to assess global optimization in SR, gaining a deeper understanding of its performance. For reasons of simplicity a function with the same structure of Arrhenius law was considered, but without units and physically relevant parameters (Eq. 27). Arrhenius law is applicable in rheology as well as in reaction kinetics.

$$y = 3 \exp\left(\frac{8}{x}\right) \quad (27)$$

Ten data points were randomly sampled in the x interval (10, 40). The calculations were performed on an Intel® Core™ i5-3337U CPU @ 1.80 GHz processor. A schematic representation of a four-layers tree for Arrhenius law identification is shown in Figure 5.

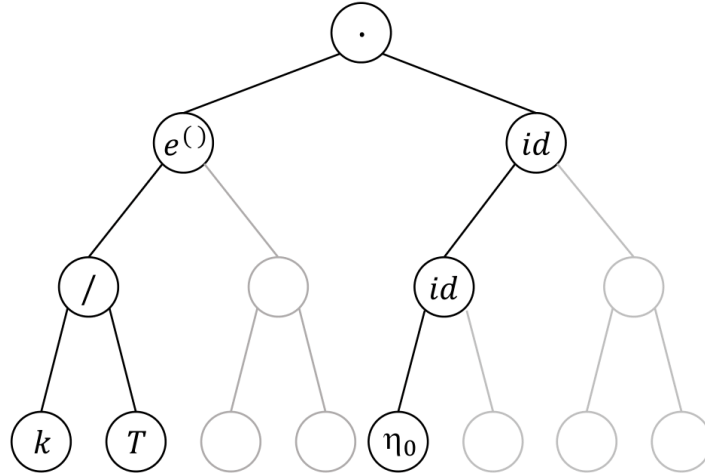


Figure 5. Binary expression tree for Arrhenius equation.

The resulting MINLP consisted of 54 binary and 154 continuous variables as well as 1174 constraints. The included operators were $\mathcal{F} = \{id, +, -, \cdot, /, \exp\}$. The globally optimal model (Eq. 28) was found within 29 s. It is a mathematically invariant model of the true one as there is an unconstrained functional search space for symbolic regression.

$$y = \frac{\exp(1.099)}{\exp(-\frac{8.0}{x})} = 3.0 \exp\left(\frac{8.0}{x}\right) \quad (28)$$

In the second instance, the function shown in Eq. 29 was adopted to evaluate the impact of the modifications in the MINLP formulation (Section 2.2) in terms of CPU time until convergence.

$$y = \frac{2x_1}{5 - x_2} \quad (29)$$

In comparison to the previously reported formulation,[26] the influence of the number of included operators (NF), data points (NE), predictor variables (NX), and tree layers (NL) is studied. Figure 6 summarises the obtained results.

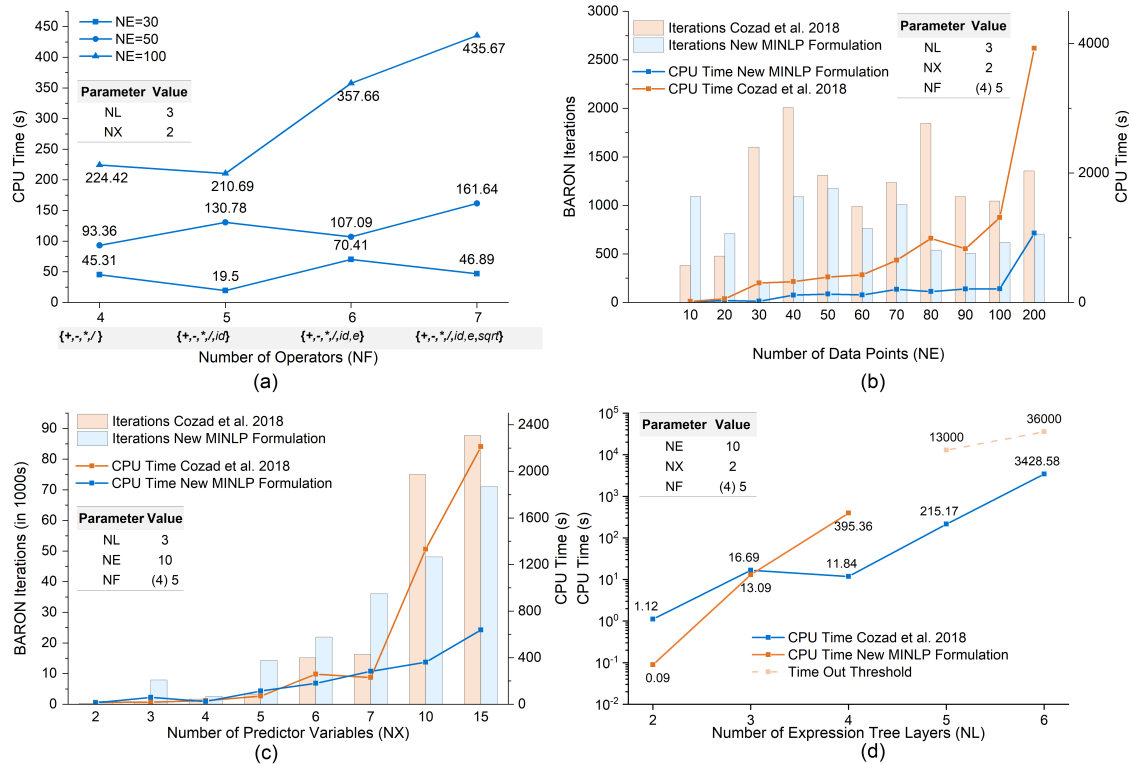


Figure 6. Scalability of the proposed methodology: (a) Scalability in the number of included operators (b) Scalability in the number of data points (c) Scalability in the number of prediction variables (d) Scalability in the number of expression tree layers.

The scale-up with regard to an increasing number of potential operators was studied first (Figure 6a). A comparison with the published format [26] is not provided as, for algebraic operators, there is neither a conceptual nor a mathematical difference in the two formulations.

Since the number of binary variables and big-M constraints increase linearly with the number of operators, a growth in the computational time is expected. The results obtained for $NE = 30$ and $NE = 50$ did not prove a consistent growth in CPU time. An overall linear trend became apparent when using an increasing number of data points.

The number of data points in the training set affects the number of nodal values and constraints in a linear fashion. The results obtained at different NE are summarised in Figure 6b and compared with the previously reported formulation [26]. For both formulations an increase in CPU time was observed, and in all cases, our adapted formulation showed improved scalability in the number of data points.

As described in Section 2.2, the conceptual difference in the modified formulation allows to reduce the number of binary variables which are required to allocate the predictor variables in the tree. As a result, a difference in performance should be observed when increasing the number of predictor variables. The expected improvements became evident at higher quantities of potential predictors and are shown in Figure 6c.

As the last parameter, the influence of the number of allowed layers in the expression tree was considered. By growing the tree in terms of the number of layers, an exponentially increasing number of nodes is added. Accordingly, the number of variables and constraints increases exponentially. As a result, the increase in CPU time is also exponential, as shown in Figure 6d. In the case of layer-scalability, the previous MINLP formulation is superior. The new proposed formulation timed out after a few hours for a number of layers greater than four. For the function under study a three-layered tree was sufficient and the ability to discard sub-layers is in favour over constructing the whole tree with identity functions. This advantage might diminish if more complex functions are sought within higher layered trees. Applicable to both formulations, this result confirms the high computational expense of SR due to the combinatorial search space. The exponential scale-up behaviour in tree layers could strongly limit the method's ability to identify more complex models.

3.2. Newton's Law of Viscosity

The data collected from a commercially available emulsion sample by means of the automated capillary viscometer were used to identify the simple linear relationship between shear stress (τ_w) and shear rate ($\dot{\gamma}_w$) at the wall of the tubing (Eq. 30).

$$\tau_w = \eta \dot{\gamma}_w \quad (30)$$

where

$$\tau_w = \frac{\Delta p R}{2 L} \quad (31)$$

$$\dot{\gamma}_w = \frac{4 Q}{\pi R^3} \quad (32)$$

For the parameter identification an expression tree with three layers, including the shear rate ($\dot{\gamma}_w$) as the only predictor variable, and ten experimental data points were used. Two data points were not included in the training set and used for the calculation of the extrapolation error. The set of operators included the basic operators and a power law $\mathcal{F} = \{id, +, -, \cdot, /, ^\wedge\}$. The resulting MINLP consisted of 28 binary and 58 continuous variables and 434 constraints. Five different instances were solved in parallel of different complexities $C = \{3,4,5,6,7\}$.

The obtained portfolio of models, Table 4, initially consisted of five models. As the complexity is constrained by an upper bound (inequality), similar models with the same complexity are identified. These, together with invariant models at higher complexities, were neglected.

Table 4. Physical model identification: Newton's law of viscosity.

Model Complexity	Identified model	Training error	Extrapolation error	Computational Time (s)
3	$\tau = 0.106\dot{\gamma}$	0.622	0.060	112
5	$\tau = 0.099\dot{\gamma}^{1.023}$	0.478	0.252	2926
7	$\tau = \dot{\gamma}^{0.213}\dot{\gamma}^{0.188}$	0.339	2.540	1139

To choose the best model among the identified ones, the prediction of the models was plotted together with the experimental data, see Figure 7, and the training and extrapolation errors were compared. As expected, the training error decreases with complexity of the model as there is more flexibility allowed to SR. However, the comparison of the extrapolation errors shows the superiority of Newton's law model ($C = 3$), whereas the other identified models suffer from overfitting. Overall, the Newton's law model can be selected as the sparsest model with the

highest generalisation capability, and can be easily interpreted to generate knowledge about the physics of the system under investigation.

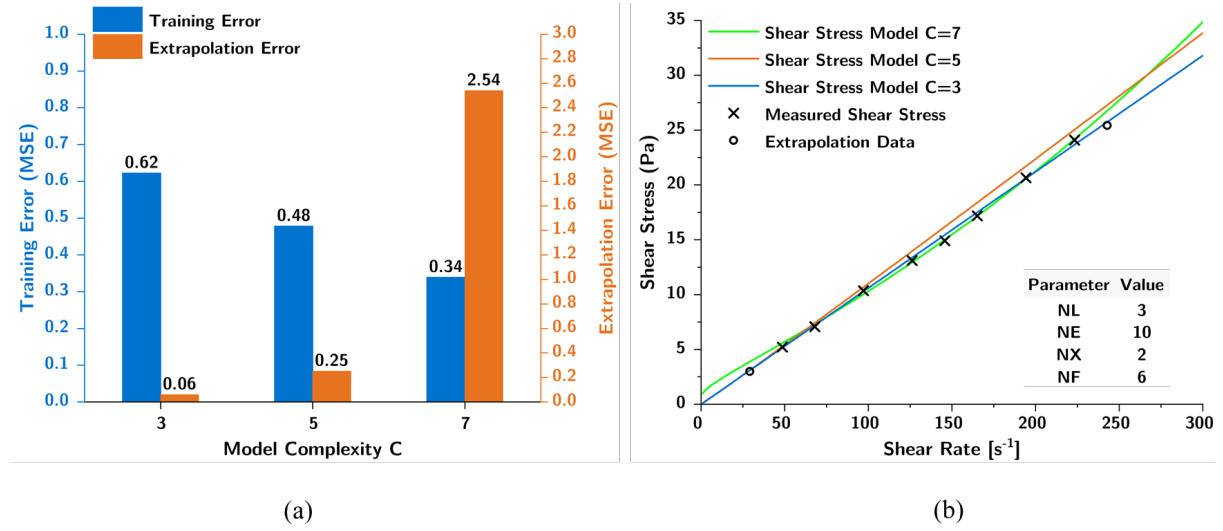


Figure 7. Physical model selection for Newtonian power law: (a) measured shear stress and shear rate: data used for model training (b) errors in training and extrapolation data set for model identification

The model identification was conducted with only ten data points, highlighting the sparsity of required data points of the presented method compared to other data-driven methods. This is especially beneficial in chemistry, where data points can be expensive to generate.

3.3. Non-Newtonian Power Law

Identification of a non-Newtonian power law was used to prove that the model selection framework favours higher complexity models where required. Eleven experimental data points were collected using 1% w/w aqueous carboxymethyl cellulose at different flow rates. As for the Newton's law identification, an expression three with three layers was used, including the shear rate ($\dot{\gamma}_w$) as the only predictor variable. Five MINLP instances were solved $C = \{3, 4, 5, 6, 7\}$ in parallel. The resulting MINLPs have 24 binary and 65 continuous variables and 486 constraints. The data were pre-processed scaling down the apparent shear rate (Eq. 32) by a factor of 10 before training. Especially when including power law operations, this allowed to keep the variable bounds and big-M values calculated by interval arithmetic low, reducing the overall search space of the solver.

With the afore-mentioned settings, the portfolio of six models, Table 5, was obtained within 13 min. The mathematically invariant and similar models were discarded. It is worth mentioning that the same power law was found for $C = \{5,6,7\}$. The resulting portfolio consists of three different models of different complexity.

Table 5. Physical model identification: non-Newtonian power law.

Model Complexity	Identified model
3	$\tau = 0.1\dot{\gamma}$
4	$\tau = 4.448 + 0.1\dot{\gamma}$
5	$\tau = 0.712\dot{\gamma}^{0.671}$

It is noteworthy that all the computed models can be physically interpreted as the ones describing Newtonian fluids, Bingham fluids and non-Newtonian (power law) fluids. The model selection was carried out comparing their training and validation errors. Once again, two experimental data points at the edges of the investigated range of shear rates were taken aside and used to evaluate the extrapolation performance of the obtained models. The three candidate models together with their performance on the training and validation data are shown in Figure 8. In this case, both the training and the extrapolation error decrease with the complexity indicating that the most complex power law is the most appropriate for the description of the experimental data.

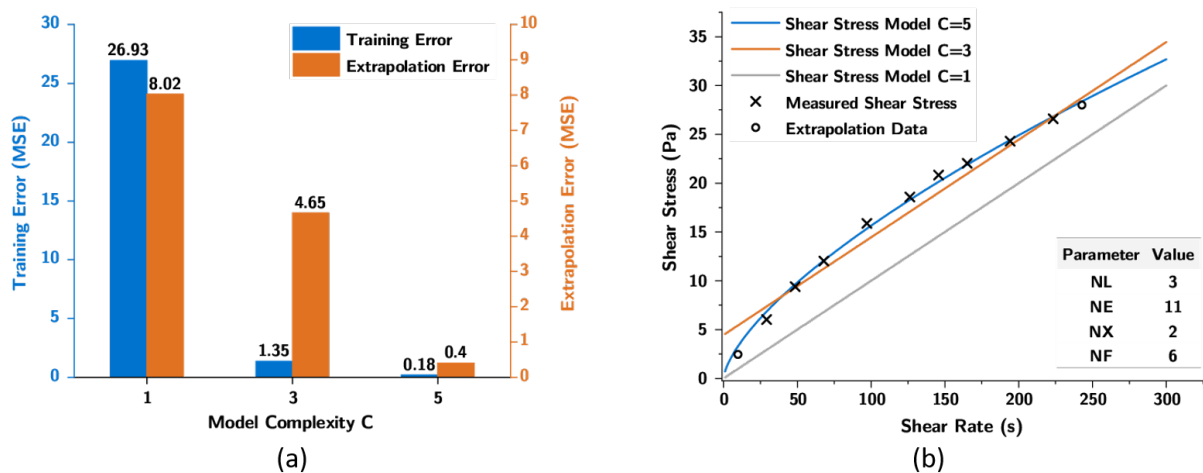


Figure 8. Physical model selection for non-Newtonian power law: (a) measured shear stress and shear rate: data used for model training (b) errors in training and extrapolation data set for model identification.

3.4. First-order Kinetic Law

Previous examples show the potential of the adopted methodology to discover sparse and interpretable models to describe the viscous behaviour of different fluids with a limited amount of experimental data. Moreover, a simple procedure was proven to be effective to select the most appropriate model within the obtained portfolio. In the following the same procedure was applied to learn a kinetic model of a simple test reaction for which a large amount of experimental data was available.

According to literature [38], hydrolysis of carboxylic acid esters can be described by first order kinetic law, Eq. 33.

$$r = k_h[PNPA] \quad (33)$$

where $[PNPA]$ is the molar concentration of the investigated compound (4-nitrophenyl acetate) and the kinetic constant k_h can be expressed as shown in Eq. 34.

$$k_h = k_N + k_A[H^+] + k_B[OH^-] \quad (34)$$

Under the adopted experimental conditions ($\text{pH} > 10.52$) the terms k_N and $k_A[H^+]$ are negligible and the overall kinetic law is given in Eq. 35.

$$r = k_B[OH^-] [PNPA] \quad (35)$$

In the first attempt, experimental data were collected at a fixed pH of 10.52, at different PNPA concentrations. Concentrations vs time data series were pre-processed to obtain an approximation of the reaction rate over time using the centered difference approximation.

For this example, a three-layer tree structure was allowed, including $[PNPA]$ as the only predictor variable. Due to the relatively low values of the measured reaction rates ($10^{-7} - 10^{-9} \text{ mol} \cdot \text{L}^{-1} \cdot \text{s}^{-1}$), they were expressed as $\text{mmol} \cdot \text{L}^{-1} \cdot \text{h}^{-1}$.

Five MINLP instances were solved $C = \{3,4,5,6,7\}$ in parallel. The resulting MINLPs have 21 binary and 219 continuous variables and 1354 constraints. 31 experimental data points were split into 25 training examples and 6 validation data points in the range $[PNPA] \in (1.11 \cdot 10^{-3} - 4.63 \cdot 10^{-2} \text{ mmol} \cdot \text{L}^{-1})$. The validation points were chosen at the lower/upper end of the dataset in order to test the extrapolation ability of the model.

The model portfolio without doubling is summarized in Table 6. The identified model with the lower extrapolation error is the true underlying first-order kinetic law governing the physics of the chemical system.

Table 6. Physical model identification: First order kinetic law.

Model Complexity	Identified model	Training error	Extrapolation error	Computational Time (s)
3	$r = 8.49[PNPA]$	$3.10 \cdot 10^{-4}$	$2.4 \cdot 10^{-3}$	62.6
5	$r = 8.42[PNPA] + 0.00133$	$2.80 \cdot 10^{-4}$	$2.7 \cdot 10^{-3}$	84.3

3.5. First-order Kinetic Law: dependence on pH

In the second attempt experimental data collected at different pH were included in the training algorithm to identify the dependence of the kinetic constant on the $[OH^-]$ concentration. The training data set consisted of 80 experimental data obtained varying $[PNPA]$ and $[OH^-]$ in the ranges $5.04 \cdot 10^{-4} - 4.55 \cdot 10^{-2} \text{ mmol} \cdot \text{L}^{-1}$ and $3.33 \cdot 10^{-1} - 4.33 \text{ mmol} \cdot \text{L}^{-1}$, respectively.

An expression tree with three layers was used again. The set of operators included the basic operators and a power law $\mathcal{F} = \{id, +, -, \cdot, /, ^\wedge\}$. The resulting MINLP consisted of 25 binary and 450 continuous variables, and 2973 constraints. Five different instances were solved in parallel of different complexities $C = \{3,4,5,6,7\}$. 80 experimental data points were split into 80 % training examples and 20 % validation data in the ranges $[PNPA] \in (5.00 \cdot 10^{-4} - 4.63 \cdot 10^{-2} \text{ mmol} \cdot \text{L}^{-1})$ and $[OH] \in (3.33 \cdot 10^{-1} - 4.33 \cdot 10^{-2} \text{ mmol} \cdot \text{L}^{-1})$. Reaction rates were expressed as $\text{mmol} \cdot \text{L}^{-1} \cdot \text{h}^{-1}$.

According to the examples reported in Sections 3.4, invariant models were obtained for $C = 4$ and 6, and discarded. The obtained portfolio of models is summarized in Table 7.

Table 7. Physical model identification: kinetics dependence on pH.

Model Complexity	Identified model	Training error	Extrapolation error	Computational Time (s)
3	$r = 33.2[PNPA]$	6.70	4.59	241
5	$r = 26.6[PNPA][OH^-]$	0.03	0.002	226
7	$r = (24.7 + [OH^-])[PNPA][OH^-]$	0.022	0.003	4287

As shown, both errors are of 2 or 3 orders of magnitude higher when $C = 3$, whereas similar training errors were obtained when $C = 5$ and $C = 7$. However, the lowest extrapolation error suggests that the model with complexity $C = 5$ is the most suitable one for the description of the kinetic behaviour of the system.

4. Conclusions

Based on a MINLP formulation for global SR reported in the mathematical domain, a different approach in setting up the superstructure was introduced to reduce the number of binary variables involved in globally optimal SR. In addition, this formulation is complemented with a framework to enable the automated identification of physical models from crude data.

The new approach was found to outperform the previously proposed ones in term of computational time when increasing the number of included operators, predictor variable and experimental data. As examples, the developed method allowed to correctly identify the models underlying the rheological behaviour of Newtonian and non-Newtonian fluids, as well as simple kinetic laws, also in the case of sparse data sets, which is a common scenario in chemical process development.

A significant limitation of the methodology was found in the exponential scale-up of the computational time for (a) an increasing number of adopted layers in the tree necessary to represent complex algebraic structures of analytical function type models, and (b) an increasing number of data points.

This serious issue of computational efficiency enhancement cannot be resolved by parallelization. This, at present, limits the identification of more complex models for which a

larger number of algebraic operators is often needed. Work is currently underway on alternative approaches using rigorous mathematical programming, such as the one presented in this work, as well as complementary methodologies to derive globally optimal fitted model structures.

Acknowledgements

PN is grateful for his Erasmus funding received for the exchange between RWTH Aachen and the University of Cambridge and that the exchange programme is co-funded by the Department of Chemical Engineering and Biotechnology, University of Cambridge, and Sustainable Reaction Engineering group of Prof. A. Lapkin. LC is grateful to BASF for co-funding her PhD. This project is also co-funded by the UKRI project “Combining Chemical Robotics and Statistical Methods to Discover Complex Functional Products” (EP/R009902/1), and co-funded by the National Research Foundation (NRF), Prime Minister’s Office, Singapore under its Campus for Research Excellence and Technological Enterprise (CREATE) program as a part of the Cambridge Centre for Advanced Research and Education in Singapore Ltd (CARES).

References

- [1] C.W. Coley, R. Barzilay, T.S. Jaakkola, W.H. Green, K.F. Jensen, Prediction of Organic Reaction Outcomes Using Machine Learning, *ACS Cent. Sci.* 3 (2017) 434–443.
- [2] A. Echtermeyer, Y. Amar, J. Zakrzewski, A. Lapkin, Self-optimisation and model-based design of experiments for developing a C–H activation flow process, *Beilstein J. Org. Chem.* 13 (2017) 150–163.
- [3] H. Gao, T.J. Struble, C.W. Coley, Y. Wang, W.H. Green, K.F. Jensen, Using Machine Learning To Predict Suitable Conditions for Organic Reactions, *ACS Cent. Sci.* 4 (2018) 1465–1476.
- [4] M.I. Jeraal, N. Holmes, G.R. Akien, R.A. Bourne, Enhanced process development using automated continuous reactors by self-optimisation algorithms and statistical empirical modelling, *Tetrahedron.* 74 (2018) 3158–3164.
- [5] C.W. Coley, W. Jin, L. Rogers, T.F. Jamison, T.S. Jaakkola, W.H. Green, R. Barzilay, K.F. Jensen, A graph-convolutional neural network model for the prediction of chemical reactivity, *Chem. Sci.* 10 (2019) 370–377.
- [6] V. Duros, J. Grizou, W. Xuan, Z. Hosni, D.L. Long, H.N. Miras, L. Cronin, Human versus Robots in the Discovery and Crystallization of Gigantic Polyoxometalates,

- Angew. Chemie - Int. Ed. 56 (2017) 10815–10820.
- [7] G. Skoraczynski, P. Dittwald, B. Miasojedow, S. Szymkuć, E.P. Gajewska, B.A. Grzybowski, A. Gambin, Predicting the outcomes of organic reactions via machine learning: are current descriptors sufficient?, *Sci. Rep.* 7 (2017) 3582.
 - [8] B.J. Reizman, K.F. Jensen, Feedback in Flow for Accelerated Reaction Development, *Acc. Chem. Res.* 49 (2016) 1786–1796.
 - [9] C. Houben, A.A. Lapkin, Automatic discovery and optimization of chemical processes, *Curr. Opin. Chem. Eng.* 9 (2015) 1–7.
 - [10] A.A. Lapkin, P.K. Plucinski, Chapter 1. Engineering Factors for Efficient Flow Processes in Chemical Industries, in: *Chem. React. Process. under Flow Cond.*, Royal Society of Chemistry, Cambridge, 2009: pp. 1–43.
 - [11] C.J. Richmond, H.N. Miras, A.R. de la Oliva, H. Zang, V. Sans, L. Paramonov, C. Makatsoris, R. Inglis, E.K. Brechin, D.-L. Long, L. Cronin, A flow-system array for the discovery and scale up of inorganic clusters, *Nat. Chem.* 4 (2012) 1037–1043.
 - [12] D.W. Robbins, J.F. Hartwig, D.W.C. MacMillan, A simple, multidimensional approach to high-throughput discovery of catalytic reactions., *Science.* 333 (2011) 1423–7.
 - [13] A.A. Lapkin, P.K. Heer, P.-M. Jacob, M. Hutchby, W. Cunningham, S.D. Bull, M.G. Davidson, Automation of route identification and optimisation based on data-mining and chemical intuition, *Faraday Discuss.* 202 (2017) 483–496.
 - [14] J.P. McMullen, M.T. Stone, S.L. Buchwald, K.F. Jensen, An Integrated Microreactor System for Self-Optimization of a Heck Reaction: From Micro- to Mesoscale Flow Systems, *Angew. Chemie Int. Ed.* 49 (2010) 7076–7080.
 - [15] D.P. Solomatine, A. Ostfeld, Data-driven modelling: some past experiences and new approaches, *J. Hydroinformatics.* 10 (2008) 3–22.
 - [16] M. Schmidt, H. Lipson, Distilling Free-Form Natural Laws from Experimental Data, n.d. <http://science.sciencemag.org/> (accessed March 5, 2019).
 - [17] O. Wolkenhauer, Why model?, *Front. Physiol.* 5 (2014) 21.
 - [18] B.J. Reizman, K.F. Jensen, Feedback in Flow for Accelerated Reaction Development, *Acc. Chem. Res.* 49 (2016) 1786–1796.
 - [19] C.S. Horbaczewskyj, C.E. Willans, A.A. Lapkin, R.A. Bourne, *Green Chemical Engineering*, Wiley-VCH, Weinheim, 2018.
 - [20] A. A. Lapkin, A. Voutchkova, P. Anastas, A conceptual framework for description of complexity in intensive chemical processes, *Chem. Eng. Process. Process Intensif.* 50

- (2011) 1027–1034.
- [21] S.L. Brunton, J.L. Proctor, J.N. Kutz, Discovering governing equations from data by sparse identification of nonlinear dynamical systems., *Proc. Natl. Acad. Sci. U. S. A.* 113 (2016) 3932–7.
 - [22] J. Bongard, H. Lipson, Automated reverse engineering of nonlinear dynamical systems, 2007. <http://www.pnas.org/content/pnas/104/24/9943.full.pdf> (accessed October 2, 2018).
 - [23] B. Tarawneh, W. AL Bodour, K. Al Ajmi, Intelligent Computing Based Formulas to Predict the Settlement of Shallow Foundations on Cohesionless Soils, *Open Civ. Eng. J.* 13 (2019) 1–9.
 - [24] Y. Wang, N. Wagner, J.M. Rondinelli, Symbolic regression in materials science, n.d. <https://arxiv.org/pdf/1901.04136.pdf> (accessed January 22, 2019).
 - [25] V.S. Vassiliadis, Y. Wang, H. Arellano-Garcia, Y. Yuan, A Novel Rigorous Mathematical Programming Approach to Construct Phenomenological Models, *Comput. Aided Chem. Eng.* 37 (2015) 707–712.
 - [26] A. Cozad, N. V. Sahinidis, A global MINLP approach to symbolic regression, *Math. Program.* (2018) 1–23.
 - [27] E.R. Gansner, S.C. North, K.P. Vo, DAG—a program that draws directed graphs, *Softw. Pract. Exp.* 18 (1988) 1047–1062.
 - [28] I.E. Grossmann, Review of Nonlinear Mixed-Integer and Disjunctive Programming Techniques, *Optim. Eng.* 3 (2002) 227–252.
 - [29] P. Belotti, C. Kirches, S. Leyffer, J. Linderoth, J. Luedtke, A. Mahajan, Mixed-integer nonlinear optimization, *Acta Numer.* 22 (2013) 1–131.
 - [30] R. Misener, C.A. Floudas, ANTIGONE: Algorithms for coNTinuous / Integer Global Optimization of Nonlinear Equations, *J. Glob. Optim.* 59 (2014) 503–526.
 - [31] M.R. Kılınç, N. V. Sahinidis, Exploiting integrality in the global optimization of mixed-integer nonlinear programming problems with BARON, *Optim. Methods Softw.* 33 (2018) 540–562.
 - [32] P. Belotti, J. Lee, L. Liberti, F. Margot, A. Wächter, Branching and bounds tightening techniques for non-convex MINLP, *Optim. Methods Softw.* 24 (2009) 597–634.
 - [33] Y. Lin, L. Schrage, The global solver in the LINDO API, *Optim. Methods Softw.* 24 (2009) 657–668..
 - [34] S. Vigerske, A. Gleixner, SCIP: global optimization of mixed-integer nonlinear

- programs in a branch-and-cut framework, *Optim. Methods Softw.* 33 (2018) 563–593.
- [35] M. Quade, M. Abel, K. Shafi, R.K. Niven, B.R. Noack, Prediction of dynamical systems by symbolic regression, *Phys. Rev. E* 94 (2016) 12214.
- [36] J. Klausen, M.A. Meier, R.P. Schwarzenbach, Assessing the Fate of Organic Contaminants in Aquatic Environments: Mechanism and Kinetics of Hydrolysis of a Carboxylic Ester, 1997. <https://pubs.acs.org/sharingguidelines> (accessed April 9, 2019).
- [37] Peter S. Marrs, Class Projects in Physical Organic Chemistry: The Hydrolysis of Aspirin, *J. Chem. Educ.* 81 (2004) 870–873. www.JCE.DivCHED.org (accessed April 9, 2019).
- [38] H.J. Goren, M. Fridkin, The Hydrolysis of p-Nitrophenylacetate in Water. Mechanism and Method of Measurement, *Eur. J. Biochem.* 41 (1974) 263–272.

Supplementary material

Identifying physico-chemical laws from robotically collected data

Pascal Neumann,^a Liwei Cao,^{b,c} Danilo Russo,^b Vassilios S. Vassiliadis,^b Alexei A. Lapkin^{b,c1}

^aAachener Verfahrenstechnik – Process Systems Engineering, RWTH Aachen University, Aachen, Germany

^bDepartment of Chemical Engineering and Biotechnology, University of Cambridge, Cambridge, CB3 0AS

^cCambridge Centre for Advanced Research and Education in Singapore, CARES Ltd.

1 CREATE Way, CREATE Tower #05-05, 138602 Singapore

Viscometer Automation and Calibration

Figure S1. shows the developed graphical user interphase in LabView. It implements all the main functions including the manual and automated control of the syringe pump, data acquisition and saving into a comma delimited file. The collected data of temperature, pressure and volumetric flow rate are acquired at a frequency of 1500 Hz with the described National Instruments modules and is then averaged to reduce random errors. For calibration purposes, the Hagen-Poiseuille law for Newtonian liquids and the empirical viscosity models for the viscosity of glycerol-water mixtures at varying weight contents were implemented. Moreover, the functionalities to measure the viscosity at different shear rates automatically and the drying procedure after the cleaning with isopropanol were designated in the interface [1,2].

¹ Corresponding author. A. Lapkin email: aal35@cam.ac.uk

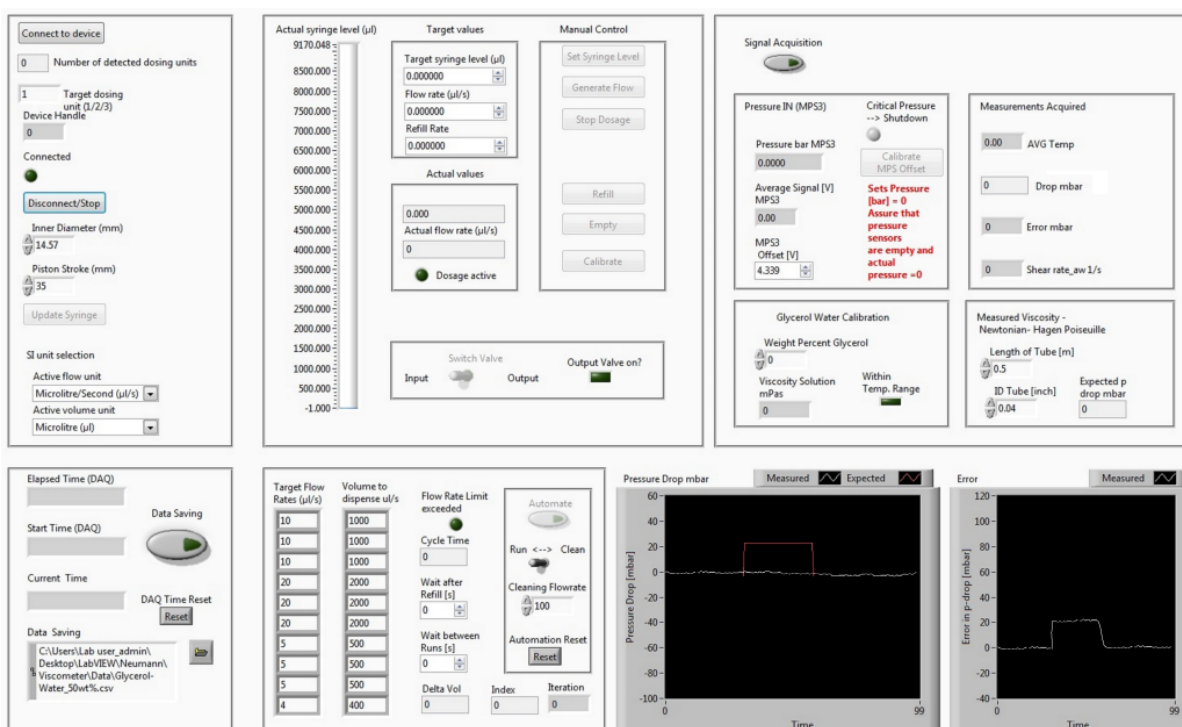


Figure S1. LabView front Panel – Graphical User Interface for the Automated Operation of the capillary viscometer.

Before the whole set-up was calibrated with the glycerol-water mixtures, the pressure sensor (Elveflow MPS3) was calibrated in combination with the NI-9702 module. This was necessary as the supplier calibration is based on their own Elveflow amplification and data acquisition system which was not available for the set-up. The sensor was calibrated with a static air pressure by closing the sensor outlet and applying a known pressure with the DRUCK digital pressure indicator. The voltage signals were then acquired in LabView for one to two minutes for each constant pressure within the interval (0,2] bar with steps of 0.1 bar. Based on the arithmetic mean for each, the calibration curve was plotted in Figure S2. The resulting linear fit was implemented in the LabView program, to convert the voltage signals into measured pressure.

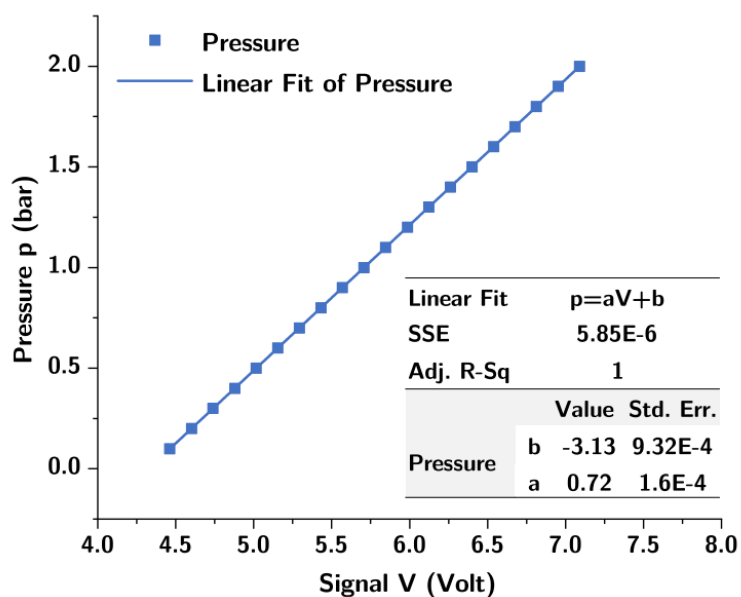


Figure S2. Calibration of the Pressure Sensor Elveflow MPS3 (0-2bar) with DRUCK DPI 600/IS.

Condition adopted for the kinetic experiments

The initial conditions for all the carried-out kinetic experiments are summarized in Table S1. All the reaction mixtures contain 1 M KCl and 1% (v/v) MeOH.

Table S1. Initial conditions for the kinetic study of 4-nitrophenyl acetate (PNPA) hydrolysis reaction.

Run	[PNPA] ₀ (mM)	pH
1	$3.43 \cdot 10^{-2}$	10.52
2	$5.76 \cdot 10^{-2}$	10.52
3	$1.95 \cdot 10^{-2}$	10.52
4	$1.28 \cdot 10^{-2}$	10.52
5	$1.30 \cdot 10^{-2}$	11.07
6	$2.86 \cdot 10^{-2}$	11.30
7	$3.36 \cdot 10^{-2}$	11.37
8	$3.70 \cdot 10^{-2}$	11.07

References

- [1] A. Volk, C.J. Kähler, Density model for aqueous glycerol solutions, *Exp. Fluids*. 59 (2018) 75. doi:10.1007/s00348-018-2527-y.
- [2] N.S. Cheng, Formula for the viscosity of a glycerol-water mixture, *Ind. Eng. Chem. Res.* 47 (2008) 3285–3288. doi:10.1021/ie071349z.