

# DEVELOPING GREY-BOX DYNAMIC PROCESS MODELS

C. de Prada \*, D. Hose\*\*, G. Gutierrez\*, J.L. Pitarch\*

\* *Department of Systems Engineering and Automatic Control, University of Valladolid  
c/ Real de Burgos s/n. Sede Mergelina EII, 47011, Valladolid, Spain  
(e-mails: prada@autom.uva.es , gloria@autom.uva.es , jose.pitarch@autom.uva.es)*

\*\* *Institute of Engineering and Computational Mechanic, Universität Stuttgart  
Pfaffenwaldring 9, 70569 Stuttgart, Germany (e-mail: dominik.hose@itm.uni-stuttgart.de)*

**Abstract:** This paper presents a methodology for developing grey models of process systems, that is, models that, being based on fundamental principles and laws of nature, combine them with sub-models obtained from experimental data. The method follows two steps: the first one takes advantage of what is known, while the second uses the data and mixed-integer optimization algorithms to identify the structure and parameters of the remaining parts of the model. The method is illustrated in a challenging biotechnological process: the Acetone-Butanol-Ethanol (ABE) fermentation process.

© 2018, IFAC (International Federation of Automatic Control) Hosting by Elsevier Ltd. All rights reserved.

**Keywords:** Modelling methodology, grey-box models, structure, identification, fermentation process.

## 1. INTRODUCTION

Quite often people say that modelling is an art, which surely is true. Nevertheless, every art requires an associated technique and developing dynamic models of industrial processes is not an exception. Traditionally, two main approaches have been used: First-principles and Data-based methodologies.

The former tries to develop models formulating equations according to physical laws that are pertinent to the process considered. Typical formulations in the process industry use mass and energy balances, phase equilibrium, etc. This requires knowledge of the process and the concerned laws and good judgement in establishing hypothesis that support the validity of the model. Deciding which phenomena and equations should be incorporated into the model is not easy, as they represent the compromise between simplicity of use and fidelity in the representation of reality. The time needed to develop these models should not be underestimated, but the use of a modern simulation environment and a certain experience should facilitate the task. As physic-chemical laws are usually valid under a wide range of conditions, one important advantage when using these models refers to the confidence they provide and the associated extrapolation capabilities.

On the opposite hand, data based methods try to discover the model that relates several process variables by analyzing and correlating sets of experimental values. Then, a model structure with adjustable parameters is proposed and the parameters can be estimated so that the models output adjusts as well as possible to the experimental data. Here, a selection of the candidate model structure and the parameter estimation algorithm are the key components. These types of models are, in principle, easier to formulate and understand, and have the advantage of being closely related to reality as they are obtained directly from data, their main drawback being the limited extrapolation capabilities they have outside the region covered by the experimental data used in their development.

Notice that, in practice, the first principle models always incorporate unknown parameters, hence, when a model of reality is proposed, developing it always involves a stage of data collection and adjustment of the model parameters so that the model outputs fit the experimental data, as in (1):

$$\min_{p, x_0} \sum_{t=1}^N (y_m(t) - y(u, x, p, t))^2 \quad \text{s.t. :} \quad (1)$$

$$\dot{x} = f(x, u, p), \quad x(0) = x_0$$

$$y = h(x, u, p)$$

where  $u$  is the vector of known manipulated variables,  $x$  are the state variables with initial value  $x_0$ ,  $y_m$  are the process measurements of the output variables  $y$  of the model and  $p$  represents the model parameters. The dynamic model is described in terms of the vector functions  $f(\cdot), h(\cdot)$ .

In the same way, a proper selection of the model structure in data based models implies certain knowledge of the interactions and phenomena taking place in the process, so that both methodologies have some points in common. The main differences are related to which element, knowledge or data analysis, play the central role.

Choosing one or the other approach is dictated mainly by the final aim and use. The model requirements are different, for instance, for operators' training than for controller design. If the purpose is using the model to take decisions about process operation, then probably a dynamic optimization problem similar to (2) is to be solved at regular time intervals:

$$\min_u J(u) \quad \text{s.t. :} \quad \dot{x} = f(x, u, p), \quad x(0) = x_0 \quad (2)$$

$$y = h(x, u, p), \quad g(x, u, p) \leq 0$$

where  $u$  is now the vector of decision variables,  $x$  are the state variables,  $y$  the model outputs and  $p$  represents the model parameters. The dynamic model is again described in terms of the functions  $f(\cdot), h(\cdot)$ , whereas  $g(\cdot)$  denotes the

problem specific constraints and  $J(u)$  is the cost function to be minimized. Solving (2) may involve a lot of computation depending on the size and structure of the model and the degrees of freedom.

This formulation assumes that the structure of the model is correct and the parameters are being estimated in the right way. Normally, a first principles model is preferable in decision making, as it may provide more confidence and a wider range of validity. Nevertheless, these assumptions and expectations may fail due to:

- The difficulty to model certain relationships among variables, due to unknown phenomena, complex relationships, impossibility to measure certain variables required to adjust the model parameters, etc.
- The computational load associated to very detailed models that makes solving (2) impractical

In these cases, it may be sensible to combine what is known with certainty about the model, such as equations representing mass or energy balances, with other types of models obtained from measurements, representing the more difficult or complex parts of the process. This results in a hybrid model that, being a mixture of “white box” first principles and “black box” sub-models is known as a “grey box” one. There are many good reviews and publications on modelling, both covering first principles (Cellier, 1991) and data based approaches (Zou, Li and Zhang, 2017), but there is a lack of literature for the systematic development of grey-box ones.

Therefore, this paper presents a methodology for developing grey-box models, dealing in particular with the problem of structure identification of the black-box elements of the models and discussing the software tools available. The method is illustrated in a challenging application: the Acetone-Butanol-Ethanol (ABE) fermentation process.

The paper is organized as follows: The modelling methodology is explained in Section 2 besides software tools that help in these tasks. The ABE process and a preliminary model are described in Section 3. Then, the grey-box model for the ABE process is developed in Section 4 together with some validation results. Finally, the paper ends with some remarks.

## 2. GREY BOX MODELLING

When developing first-principle models in the process industry, one may face situations in which a partial set of equations  $f(\cdot), h(\cdot)$  is known, as in (2), but a subset of variables  $z(x, u)$  is not (or too complex to model). However, these relations between  $z, x$  and  $u$  are required to complete the model.

$$\dot{x} = f(x, u, z(x, u), p), \quad y = h(x, u, z(x, u), p) \quad (3)$$

### 2.1 Problem formulation

Given the partial model (3), and assuming that one experiment has been performed so that values for  $u$  and  $y_m$  are available from the process via any data-collection system, the problem can be formulated as finding the sub-model  $z(x, u)$  and parameters  $p$  such that the response of (3) fits the experimental values in the best possible way.

A typical approach to the problem is proposing a functional structure for  $z(x, u)$ , such as  $z = r(x, u, \theta)$  and, as measurements for  $z$  rarely are accessible, adjust the parameters  $\theta$  by solving the following optimization problem:

$$\min_{\theta, p, x_0} \sum_{t=1}^N (y_m(t) - y(u, x, p, t))^2 \quad \text{s.t.}:$$

$$\dot{x} = f(x, u, r(x, u, \theta), p), \quad x(0) = x_0 \quad (4)$$

$$y = h(x, u, r(x, u, \theta), p)$$

As the initially proposed structure for  $r(\cdot)$  will likely not be correct, the fit is to be repeated with a modified candidate structure following a trial and error procedure. This is a time consuming procedure with no guarantee of success.

### 2.2 Proposed modelling methodology

Instead, the following two-stage approach can be used.

1. **Estimation.** Variables  $z$  are considered as independent and included in the data-fitting problem (5) as new decision variables with an appropriate parameterization  $z_i$ :

$$\min_{z_i, p, x_0} \sum_{t=1}^N (y_m(t) - y(u, x, p, t))^2 \quad \text{s.t.}:$$

$$\dot{x} = f(x, u, z_i, p), \quad x(0) = x_0 \quad (5)$$

$$y = h(x, u, z_i, p), \quad c(x, u, z_i, p) \leq 0$$

2. **Regression.** Simulate the model  $\dot{x} = f(x, u, z_i, p)$  in (5) to generate values for  $x$ , using  $u$  and the estimated  $z_i$  as inputs. Find correlations for  $z_i$  with any  $x$  and/or  $u$ , and formulate regression constraints  $z(x, u)$ . In this way, all values are consistent with the remaining model.

Finally, the functional  $z(x, u)$  is added to the model in order to get the final expression (3).

Note that additional constraints  $c(x, u, z_i, p) \leq 0$  have been included in Stage 1 to guarantee that the values of  $z$  over time conform to physically admissible ones, such as being positive, larger than other variables, etc.

The parameterization  $z_i, i = 1, \dots, n$  can be simple (constant values over a set of discrete-time intervals), or more complex ones based on collocation points, according to the problem nature and the expected time evolution of these variables. Additionally, the consideration of  $z$  as independent variables adds extra degrees of freedom that facilitate the model fit to the experimental values. Note that (5) could be formulated alternatively as a dynamic data reconciliation problem with robust estimators (Huber, 2014) instead of the quadratic cost, if enough measurements were available to achieve enough redundancy.

Solving (5) provides a set of points  $z_i$  coherent with both the experimental data and the model. The resolution can be done either via sequential or simultaneous approaches: Depending on the problem structure, a combination of a dynamic simulator and an optimization algorithm (rSQP like SNOPT or an evolutionary one) can be a good choice, but modern optimization environments like CasADi (Andersson et al., 2012) or Pyomo, (Hart et al., 2012) offer excellent features, including automatic discretization by orthogonal collocation and auto-

matic differentiation, that facilitate the use of efficient interior point codes such as IPOPT in a simultaneous approach.

There are different ways to approach the problem in Stage 2. Among them, an option is to postulate a flexible general structure with, for instance, a neural network adjusting its parameters later on to fit the  $z_i$ ,  $x$  and  $u$  values. Nevertheless, this approach has some drawbacks: on the one hand, it does not take advantage of the partial knowledge that one may have about  $z$  and, on the other hand, it does not guarantee a feasible extrapolation when  $z$  takes values outside the estimated range, unless extra conditions are imposed. A good alternative is to use mixed-integer optimization and global methods to select among a combination of user-provided potential basis functions, those that provide the best fit taking into account possible extra constraints to guarantee physical coherence. Algebraic modelling environments such as ALAMO (Cozad, Sahinidis and Miller 2014; 2015) offer very good support to the fitting task using global MINLP solvers like BARON and adaptive-sampling procedures. In the next section, this methodology is applied to a challenging biotechnological process.

### 3. THE ABE PROCESS

The Acetone-Butanol-Ethanol (ABE) fermentation process has experienced a rise in popularity due to its possibilities in the production of bio-butanol which is being used in a lot of products such as bio fuels (Mayank et al. 2013). In our case, the ABE installation is a batch process (Hose et al. 2016), i.e. it is carried out in a closed fermenter without any input or output flow, the cells consume the substrate and generate the useful products. Only the temperature and the pH are controlled around a certain value. However, other forms such as the continuous fermentation are possible. The actual fermenter used in this study can be seen in Figure 1.



Figure 1. ABE fermenter with control equipment.

Initially, the substrate, in this case glucose, is given into the reactor along with the so-called inoculum, the microbiological starting culture, in this case *Clostridium acetobutylicum*. After a short lag phase in which the bacteria adjust themselves to the new environment, the acid production phase or acidogenesis starts in which mainly acetate, butyrate and lactate are produced. This phase is typically indicated by a rapid growth of cells until the pH has dropped from about 7 to 4.5. In a second step, the solventogenesis, the cell number stalls or even decreases and the production of the solvents

starts, that is acetone, butanol and ethanol, which is the main aim of the process. The temperature is usually maintained constant at its optimum which lies between 30 and 40°C throughout the whole process. For further insight into the biochemistry of ABE fermentation one can refer to Hubert, Andersch and Gottschalk (1982).

#### 3.1 Process model

The important features to be modelled have been identified as: cell growth and death dynamics, substrate utilization for the acid, solvent production and cell maintenance, as well as inhibition mechanisms due to an excess of both substrate and solvents in the broth. A macroscopic model is employed to capture the quantitative process dynamics with the purpose of determining the best operating conditions later on. The use of microscopic models based on metabolic pathways has not been considered as they are not adequate for the final use of the models in economic optimization of the process operation.

The nomenclature of the concentrations (in g/L) is given by  $X$ : Biomass (*C. acetobutylicum*),  $S$ : Substrate (Glucose),  $P_a$ : Solvent (Acetone),  $P_b$ : Solvent (Butanol) and  $P_e$ : Solvent (Ethanol). Considering that the volume is sensibly constant, one possible model that captures the features explained above and should, therefore, be able to reproduce the general trajectories of experimental data is (6):

$$\begin{aligned} \dot{X} &= \mu X - \lambda X, & \dot{S} &= -Y_{xs} \mu X - mX, \\ \dot{P}_a &= Y_{xa} \mu X, & \dot{P}_b &= Y_{xb} \mu X, & \dot{P}_e &= Y_{xe} \mu X \end{aligned} \quad (6)$$

Basically the model is composed of mass balances for cells, substrate and products, that we know must be satisfied. The accumulation of cells per time unit is equal to the difference between inflow and outflow of cells (which are zero as the process operates in batch mode), the rate of growth and the rate of cell death (both assumed proportional to the number of cells). Similar arguments are used to derive the other equations. The model employs a growth term  $\mu$ , death and maintenance coefficients,  $\lambda$  and  $m$ , as well as the rates  $Y_{xs}$ ,  $Y_{xa}$ ,  $Y_{xb}$  and  $Y_{xe}$  indicating how substrate is converted into cells, and products are generated as a result of the cell activity. Nevertheless, it is also well known that the growth term  $\mu$  is not constant, but depends on several factors and the relation among them,  $\mu = f(X, S, P_a, P_b, P_e)$ , is not well known.

Several models for  $\mu$  have been proposed in the literature, some of which are shown in Table 1 (Heijnen and Romain, 1995; Yang and Tsao, 1994). Note, that the original model by Monod is by far the most popular due to its simplicity and the fact that it often suffices to reproduce the general behavior of cell growth according to the Michaelis-Menten kinetics. The other models are often extensions of the one by Monod.

#### 3.2 Direct parameter estimation

For identification and validation purposes two data sets were provided. The corresponding experiments were carried out in the batch fermenter of Figure 1 with glucose as the substrate and the bacteria *C. acetobutylicum* as biomass. Measurements of concentrations of biomass, substrate and the three

Table 1. Several models for  $\mu$ . In the model by Yang,  $P_{aa}$  denotes acetate,  $P_{ba}$  denotes butyrate and  $P_l$  denotes lactate.

Model	$\mu$
Monod	$\bar{\mu} \frac{S}{S+K}$
Teissier	$\bar{\mu} \left(1 - e^{-\frac{S}{K}}\right)$
Haldane	$\bar{\mu} \frac{S}{K_1 S^2 + S + K}$
Hinshelwood	$\bar{\mu} \frac{S}{S+K} \prod_{i \in (a,b,c,aa,ba)} (1 - K_p P_i)$
Yang	$\frac{\bar{\mu} S}{S+K} \left(1 - \left(\frac{P_{aa}}{C_{maa}}\right)^{maa} - \left(\frac{P_{ba}}{C_{mba}}\right)^{mba} - \left(\frac{P_b}{C_{mb}}\right)^{mb} - m_1 \left(\frac{P_{aa}}{C_{maa}}\right)^{maa} \cdot \left(\frac{C_b}{C_{mba}}\right)^{mb} - m_2 \left(\frac{P_{ba}}{C_{mba}}\right)^{mba} \cdot \left(\frac{P_b}{C_{mba}}\right)^{mb} - m_3 \frac{5.6 - pH}{1.6}\right)$

products were taken over the batch cycle. Unfortunately, no information about the temperature or pH during the experiments were provided making it impossible to include these variables in the models.

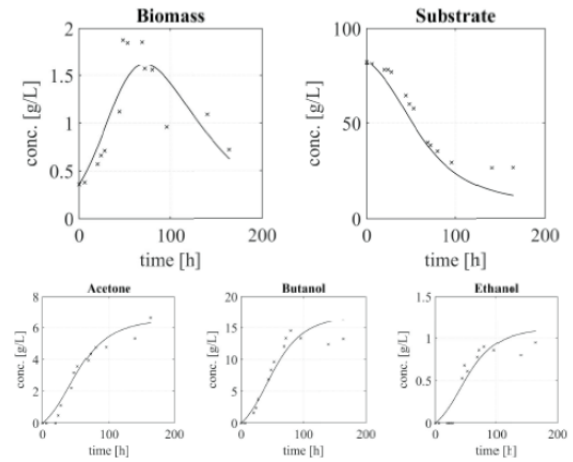
The traditional approach to parameter estimation chooses a model structure for  $\mu$  and then solves (4) to estimate the model parameters. In our case, as an example, we have chosen the Hinshelwood model from Table 1, including an inhibitory effect by butanol, as in (7):

$$\mu = \mu_0 \frac{S}{S+K_\mu} (1 - K_b P_b) \quad (7)$$

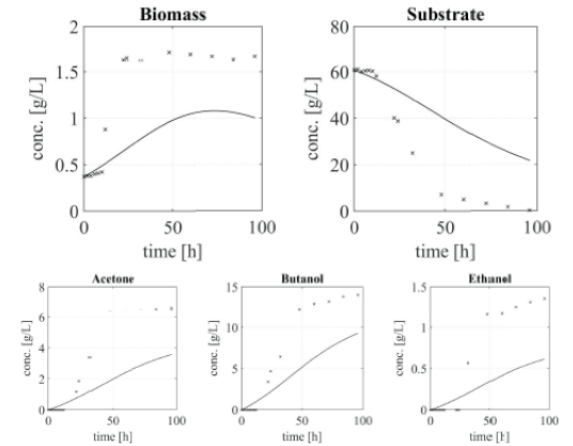
The initial value of the concentrations is assumed to be known from measurements. The coefficients  $\lambda$  and  $m$ , as well as rates  $Y_{xs}$ ,  $Y_{xa}$ ,  $Y_{xb}$  and  $Y_{xe}$  and parameters  $\mu_0$ ,  $K_\mu$ ,  $K_b$ , are unknown decision variables of (4). As the concentrations and parameters should be positive, the constraints  $y = [X, S, P_a, P_b, P_e] \geq 0$ ,  $p = [\lambda, m, Y_{xs}, Y_{xa}, Y_{xb}, Y_{xe}, \mu_0, K_\mu, K_b] \geq 0$  are imposed to the dynamic optimization, as well as normalization factors in the objective function.

A sequential approach, connecting the dynamic simulation of (6)-(7) with a nonlinear optimization (NLP) algorithm is often used for this type of parameter estimation (Boada et al., 2016). Here, the genetic algorithm provided in MATLAB was used to solve the optimal fitting to the experimental data. Note that the problem is highly non-linear and non-convex, which justifies the use of the evolutionary algorithm to avoid local minima in spite of the higher computation times. The results comparing the model response (line) and the experimental data (dots) can be seen in Figure 2a. The cell dynamics are approximated reasonably well and, although the substrate utilization yields some error, this fitting could be considered adequate. Note that the training set on the left presents some wrong data, because the concentrations of butanol and ethanol cannot decrease over time.

Unfortunately, the model validation depicted in Figure 2b demonstrates how difficult finding universal models that apply in the general case are, as they do not fit at all in different conditions.



a) Response against the identification dataset.



b) Response against the validation dataset.

Figure 2. Model response (line) and experimental data (dots).

The underlying problem with this technique is that a model for the cellular growth  $\mu$  has to be chosen in advance which may not apply in this particular case. However, because all models provided in Table 1 are heuristic, they cannot generally be considered applicable and have to be chosen carefully.

#### 4 GREY MODEL OF THE ABE PROCESS

Next, the methodology presented in Section 2 will be used to derive a grey model for the ABE process. Other successful applications are reported in, for instance, Pitarch et al. (2017).

##### 4.1 Estimation with free growth rate

In the first stage, the growth rate  $\mu$  is parameterized over time and its values are considered as independent coefficients to be estimated besides the other model parameters using a formulation similar to (5). Defining  $N + 1$  discretization points  $t_i$ ,  $i = 1, \dots, N$ , a common parameterization is  $\mu(t) = \mu_i, t \in [t_i, t_{i+1}]$  as in Figure 3.

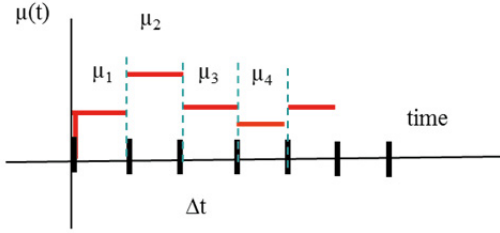


Figure 3. A zero-order parameterization of  $\mu(t)$ .

The estimation problem can then be formulated as follows:

$$\begin{aligned}
 \min_{\mu, p} \quad & \sum_{i=1}^N [y(t_i) - y_m(t_i)]^T W [y(t_i) - y_m(t_i)] + \beta \sum_{i=1}^N [\Delta \mu_i]^2 \\
 \text{s.t.} \quad & \dot{X} = \mu X - \lambda X, \quad \dot{S} = -Y_{xs} \mu X - mX, \\
 & \dot{P}_a = Y_{xa} \mu X, \quad \dot{P}_b = Y_{xb} \mu X, \\
 & \dot{P}_e = Y_{xe} \mu X, \quad \mu(t) = \mu_i, \quad t \in [t_i, t_{i+1}), \\
 & y = [X, S, P_a, P_b, P_e] \geq 0, \\
 & p = [\lambda, m, Y_{xs}, Y_{xa}, Y_{xb}, Y_{xe}] \geq 0
 \end{aligned} \tag{8}$$

With  $\Delta \mu_i$  denoting changes over consecutive time instants. Note that, in addition to the positive constraints on the concentrations and parameters, a Tikhonov  $L_2$  regularization term has been added to the quadratic objective, penalizing the changes over time of  $\mu$ , so that increasing or decreasing the weighting factor  $\beta$  we can favour more smooth evolutions or give more freedom to fit the experimental data  $y_m$ . In the objective function  $W$  is a normalization matrix that can also be used to balance the quality of the adjustment among the components of the vector  $y$ . Note also that, in spite of the larger number of parameters to be estimated, the structure of (8) is now far more simple, and the problem can be recast as a quadratic programming one.

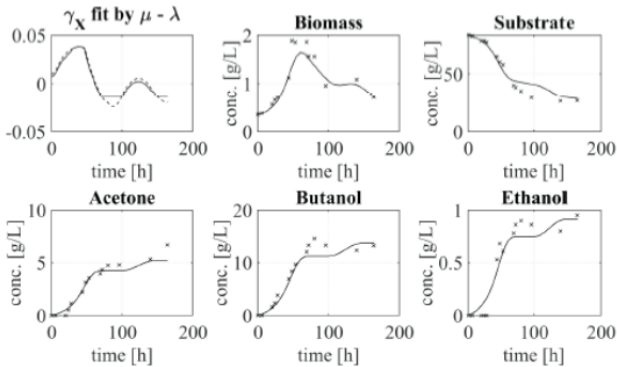


Figure 4. Comparison of the model responses and experimental data (dots) for the biomass, substrate and products. Upper left corner in dots: estimated values of  $\mu_i - \lambda$ .

The solution for a factor  $\beta = 1000$  can be seen in Table 2 and Figure 4, which also shows the time evolution of the model output besides the experimental data for the biomass, substrate and products. The fit is as good as or better than the one displayed in Figure 2a. In addition, the figure shows in dots the time evolution of the estimated values of the term  $\mu(t) - \lambda$ , denoted by  $\gamma_x$ .

Table 2: Values of the estimated model parameters.

Parameter	Value	Parameter	Value
$\lambda$	0.008	$Y_{xa}$	2.625
$m$	0.106	$Y_{xb}$	6.869
$Y_{xs}$	19.224	$Y_{xe}$	0.458

#### 4.1 Estimation of the growth rate with ALAMO

In order to complete the model, a functional relation  $\mu = f(X, S, P_a, P_b, P_e)$  has to be found. This can be done using the estimated values of  $\mu$ ,  $\mu_i$ , together with the ones of the other variables obtained by simulation. For instance, the ones depicted in Figure 4, which are consistent with the white model.

For this task, ALAMO (Automatic Learning of Algebraic MOdels) has proven to be a useful tool to get algebraic surrogate models from given data sets, (Cozad, Sahinidis and Miller, 2014). For this purpose, it provides some standard basis functions, such as monomials, logarithms or exponentials, which can be combined. Nevertheless, in this context, its strengths lie in the possibility of including user defined basis functions as well as constraints guaranteeing physical sense of the estimated functional, even outside the range of the experimental data. ALAMO automatically picks the more suitable basis functions through a mixed-integer optimization, refining the solution with adaptive sampling in the regions where the model presents larger errors.

In order to select the basis functions, notice that we can drop the dependency on  $X$ , as the model (6) already includes the term  $\mu X$ . At the same time, it suffices to let  $\mu$  depend on  $S$  and only one product, e.g.  $P_b$ , as the other ones are proportional to it. Then, in addition to the standard basis functions of ALAMO, two user-defined ones are included in the list:

$$\frac{S}{1 + S/K_s}, \quad \frac{S}{1 + P_b/K_b} \tag{9}$$

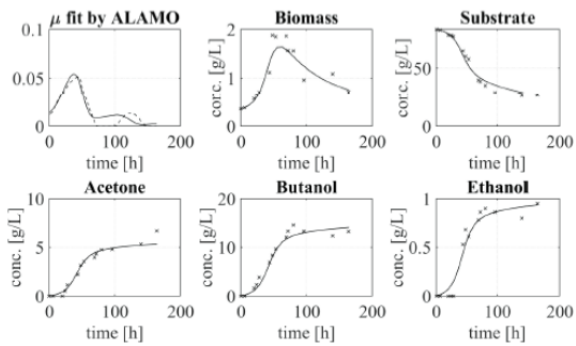
The optimization problem to be solved is given by (10), where the combination of basis function  $f_i$  that best fits the values  $\mu_i$  is selected by the binary variables  $w_j$ .  $\Lambda$  stands for an allowed range of the  $\alpha$  parameters and  $l$  is the maximum number of basis functions allowed in the solution.

$$\begin{aligned}
 \min_{\alpha_j, w_j} \quad & \sum_{i=1}^N \left[ \mu_i - \sum_{j=1}^M \alpha_j f_j(S_i, P_{bi}) \right]^2 \quad \text{s.t.} : \\
 & \Lambda^{low} w_j \leq \alpha_j \leq \Lambda^{up} w_j \quad j = 1, \dots, M \\
 & \sum_{j=1}^M w_j \leq l \quad w_j \in [0, 1]
 \end{aligned} \tag{10}$$

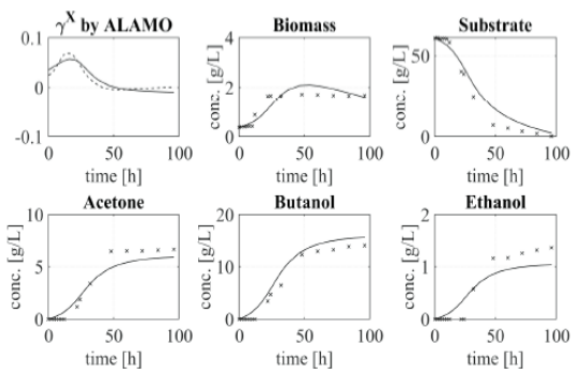
The problem is solved using the BARON code and the solution obtained is given below.

$$\mu = 3.7 \frac{S}{S + 300} - 460 S P_b - 0.0095 \frac{S}{1 + P_b / 14.042} \tag{11}$$

Top left picture in Figure 5a depicts the results of the fit comparing the values of  $\mu_i$  with the approximation (11). This expression is then incorporated into (6) to get the final grey-box model below.



a) Fit to the training dataset. Top left: fit of the growth rate.



b) Response against the validation dataset.

Figure 5. Fit of the grey-box model to the experimental data.

$$\begin{aligned}
 \dot{X} &= \mu X - 0.008 X \\
 \dot{S} &= -19.224 \mu X - 0.106 X \\
 \dot{P}_a &= 2.625 \mu X & \dot{P}_b &= 6.869 \mu X & \dot{P}_e &= 0.458 \mu X \\
 \mu &= 3.7 \frac{S}{S + 300} - 460 S P_b - 0.0095 \frac{S}{1 + P_b / 14.042}
 \end{aligned} \quad (12)$$

The model output is compared to experimental training and validation datasets in Figure 5, showing a good agreement. It becomes evident that the reliability of the simulation is highly dependent on the  $\mu$ -model accuracy, since a small error in the approximation by ALAMO still leads to some error. Comparing Figures 2b with 5b, the proposed approach yields much better (although not perfect) fits. Note that, as was already mentioned, some experimental data are not reliable.

## 5 CONCLUSIONS

A methodology for the development of grey-box models, combining first principles and data driven models, has been presented and tested successfully. Its strengths, compared to traditional approaches, are that fewer assumptions about the model have to be made, leaving additional degrees of freedom. Hence, resulting optimization problems are handier: the computational cost is lower due to the decomposition of the optimization in two steps.

Further advantages in the demonstration example include that the cellular growth term is obtained without a trial-and-error procedure.

## ACKNOWLEDGMENT

The authors wish to thank the Chemical Engineering DPT. of UVa for the experimental data provided as well as to the EU and the Spanish MINECO/FEDER for their support through the projects H2020-SPIRE CoPro (Grant Agreement n° 723575) and INOPTCON (DPI2015-70975).

## REFERENCES

- Andersson J., Akesson J., and Diehl M., (2012), “CasADi: A symbolic package for automatic differentiation and optimal control”. *Recent Advances in Algorithmic Differentiation*, vol. 87, *Lecture Notes in Comp. Science and Eng.*, pag. 297–307. eds. Forth S., Hovland P., Phipps E., Utke J. and Walther A., Springer, Berlin Heidelberg.
- Boada Y., Pitarch J.L., Vignoni A., Reynoso-Meza G., Picó J., (2016), “Optimization Alternatives for Robust Model-based Design of Synthetic Biological Circuits”. 11<sup>th</sup> *DYCOPS-CAB Conf*, pag. 821-826.
- Cellier, F., (1991). *Continuous Systems Modelling*. Springer.
- Cozad A., Sahinidis N.V., and Miller D.C., (2014), “Learning surrogate models for simulation-based optimization. *AIChE Journal*, 60(6), 2211-2227.
- Cozad A., Sahinidis N.V., and Miller D.C., (2015), “A combined first-principles and data-driven approach to model building”. *Computers & Chemical Eng.*, 73:116-127.
- Hart W.E., Laird C., Watson J-P., and Woodruff D.L., (2012) *PYOMO – Optimization Modeling in Python*. Vol. 67. Springer.
- Heijnen J.J. and Romein B., (1995), “Derivation of kinetic equations for growth on single substrates based on general properties of a simple metabolic network”. *Biotechnology progress*, 11(6):712-716.
- Hose D., de Prada C. and González G., (2016), “Modelling and Identification of ABE Fermentation Processes”. In *Actas de las XXXVII Jornadas de Automática*, 45, CEA.
- Huber, P., (2014), Robust statistics, in: M. Lovric (Ed.), *International Encyclopedia of Statistical Science*, Springer, Berlin/Heidelberg, , pp. 1248–1251
- Hubert B., Andersch W., and Gottschalk G., (1982), “Continuous production of acetone and butanol by clostridium acetobutylicum in a two-stage phosphate limited chemostat”. *Eur. Jour. of applied microbiology and biotechnology*, 15(4):201-205.
- Mayank R., Ranjan A., and Moholkar V.S., (2013), “Mathematical models of abe fermentation: review and analysis.” *Critical reviews in biotechnology*, 33(4):419-447.
- Pitarch J.L, Palacín C.G., de Prada C., Voglauer B. and Seyfriedsberger G., (2017) “Optimisation of the resource efficiency in an industrial evaporation system”. *Journal of Process Control*, 56:1–12.
- Xiaoping Y., and Tsao G.T., (1994) “Mathematical modeling of inhibition kinetics in acetone-butanol fermentation by clostridium acetobutylicum”. *Biotechnology progress*, 10(5):532-538.
- Zou W., Li C., and Zhang N., (2017) “A T-S Fuzzy Model Identification Approach based on a Modified Inter Type-2 FRCM Algorithm”, *IEEE Trans. on Fuzzy Systems*, In Press, DOI: 10.1109/TFUZZ.2017.2704542.