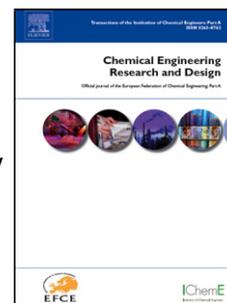


# Journal Pre-proof

A model-based experimental design approach to assess the identifiability of kinetic models of hydroxymethylfurfural hydrogenation in batch reaction systems

Philipp Deussen, Federico Galvanin



PII: S0263-8762(21)00529-3  
DOI: <https://doi.org/10.1016/j.cherd.2021.12.028>  
Reference: CHERD 4753

To appear in: *Chemical Engineering Research and Design*

Received Date: 11 August 2021  
Revised Date: 14 December 2021  
Accepted Date: 16 December 2021

Please cite this article as: Deussen P, Galvanin F, A model-based experimental design approach to assess the identifiability of kinetic models of hydroxymethylfurfural hydrogenation in batch reaction systems, *Chemical Engineering Research and Design* (2021), doi: <https://doi.org/10.1016/j.cherd.2021.12.028>

This is a PDF file of an article that has undergone enhancements after acceptance, such as the addition of a cover page and metadata, and formatting for readability, but it is not yet the definitive version of record. This version will undergo additional copyediting, typesetting and review before it is published in its final form, but we are providing this version to give early visibility of the article. Please note that, during the production process, errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

© 2020 Published by Elsevier.

# A model-based experimental design approach to assess the identifiability of kinetic models of hydroxymethylfurfural hydrogenation in batch reaction systems

Philipp Deussen,<sup>a</sup> Federico Galvanin<sup>a,\*</sup>

<sup>a</sup>*Department of Chemical Engineering, University College London, Torrington Place, London WC1E 7JE, UK  
f.galvanin@ucl.ac.uk*

---

## Highlights

- A review of kinetic models for describing hydroxymethylfurfural (HMF) hydrogenation is proposed
  - A crucial aspect is the model practical identifiability, i.e. the estimability of kinetic parameters from experimental data
  - A three-step approach based on optimal experimental design is proposed to assess the identifiability of candidate models
  - The new approach is tested on proposed HMF hydrogenation models including temperature dependency
  - The most informative conditions for temperature, experiment duration and initial HMF and DMF concentrations are identified
- 

## Abstract

Hydroxymethylfurfural (HMF) is an organic compound that occurs naturally in many foods and is used as feedstock in numerous chemical processes. HMF can be hydrogenated to form DMF, which is an important component in biofuel production. To date, several kinetic models have been proposed and studied in literature for this hydrogenation reaction, including power law models based on reaction species and Langmuir-Hinshelwood-Hougen-Watson (LHHW) models. For these models a critical aspect that has not been addressed in literature is related to their practical identifiability, i.e. the estimability of kinetic parameters from experimental data. Also, none of the existing models propose a temperature dependence of the kinetic parameters.

A three-step approach is presented in this paper where model-based design of experiments (MBD<sub>oE</sub>) techniques are exploited to assess the identifiability of candidate kinetic models of HMF hydrogenation in a batch reaction system. The objective is twofold: 1) to propose new kinetic models of HMF hydrogenation where the temperature is explicitly introduced as experimental design variable

and test the practical estimability of kinetic parameters from concentration data only; 2) to identify the most informative regions of the experimental design space, defined by temperature, experiment duration and initial HMF and DMF concentrations, for achieving a precise estimation of model parameters. Together with a-posteriori statistics obtained from parameter estimation from in-silico data, a MBD<sub>oE</sub> analysis gives a clear representation of the most informative experimental conditions to be explored in the future experimentation underlining distinct areas of practical parametric identifiability.

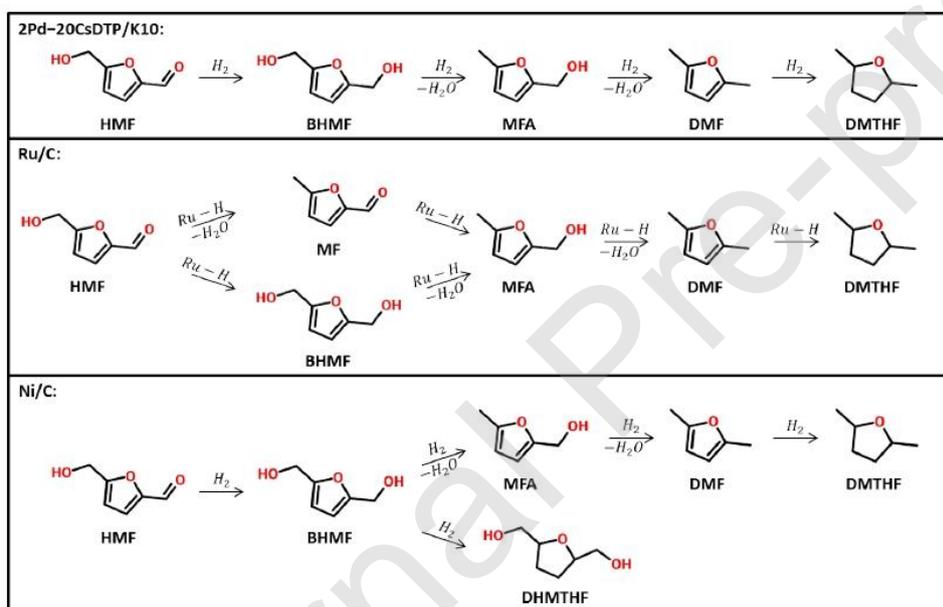
**Keywords:** identifiability analysis, model-based design of experiments, kinetics of HMF hydrogenation

---

## 1. Introduction

Hydroxymethylfurfural (HMF) is an organic compound that can be found naturally in many foods and has various uses in industry as a platform chemical. For example, it can be used as a feedstock in the production of 2,5-Furandicarboxylic Acid (FDCA), which in turn can be used in the production of PET plastic. Furthermore, it finds various applications in the fine chemicals sector, such as pharmaceuticals, flavouring agents or agrochemicals. Most interestingly, it can be used as a feedstock to produce 2,5-Dimethylfuran (DMF), which has potential use as a viable biofuel alternative or addition to diesel or jet fuel. In both cases, tests have shown that the emissions are significantly less polluting and harmful than from their fossil equivalents (van Putten et al., 2013). The hydrogenation reaction that converts HMF to DMF has been studied in literature at variable operating conditions using different catalysts (Grilic et al., 2014, Jain & Vaidya, 2016, Gawade et al., 2016, Gyngazova et al., 2017). Figure 1a shows the main proposed mechanisms for HMF hydrogenation (Bortoli, 2018). As illustrated in the figure, the mechanism that governs the reaction changes when a different catalyst is used. In most cases, the reaction shows a two-step mechanism in which HMF is converted first to bishydroxymethyl furan (BHMF) and then it is further hydrogenated to the highly reactive 5-methyl furfuryl alcohol (MFA). This intermediate rapidly undergoes hydrogenolysis to finally give DMF and small amounts of 2,5-dimethyltetrahydro furan (DMTHF) as main by-product. Most of the studies

present the formation of DMF as the rate determining step. Several kinetic models have been proposed in literature to conduct this reaction effectively and safely in scale and to characterise the hydrogenation reaction in a quantitative way at different experimental conditions. The proposed kinetic models are based on power-law and Langmuir-Hinshelwood-Hougen-Watson (LHHW) mechanisms, as reported in Figure 1b. None of the proposed models includes an explicit temperature dependency, even though the reaction kinetics practically depend on it. Instead, kinetic parameter values have been calculated at discrete temperature points. Furthermore, it is unclear from literature how precisely the kinetic parameters have actually been estimated from data. The impact of parameter uncertainty on the predicted model outputs is a crucial aspect that must be determined to assess how much predictive power these proposed models actually have.



(a)

Authors (Year)	Proposed Mechanism
Grilic, Likozar & Levec (2014)	Power law (1 <sup>st</sup> order)
Luo, Arroyo-Ramirez, Wei & Yun (2015)	Power law (1 <sup>st</sup> order)
Jain & Vaidya (2016)	Single-site LHHW
Gawade, Tiwari & Yadav (2016)	Dual-site LHHW
Gyngazova, Negahdar, Blumenthal & Palkovits (2017)	Power law (1 <sup>st</sup> order)

(b)

*Figure 1 – (a) Summary of the main proposed reaction mechanisms for HMF hydrogenation on selected catalysts (Bortoli, 2018); (b) summary of mechanisms of kinetic models for HMF hydrogenation proposed by various authors to date.*

Once a kinetic model structure is established, a key aspect that must be considered independently of the mechanistic derivations or experimental data that have been used to derive these kinetic models is parameter identifiability (Walter and Pronzato, 1996). A model is considered structurally identifiable if its set of characteristic kinetic parameters can be uniquely determined from data, and if it produces unique and distinct input-output behaviours when different parameter values are used. Many techniques have been proposed to test the structural identifiability of both linear and nonlinear systems. Among them, Laplace transform and Lie derivatives (Walter and Pronzato, 1996), power series extension (Pohjanpalo, 1978) and differential algebra (Margaria et al., 2001) are suitable methods to test identifiability under measurement error-free conditions. A comprehensive review of methods to test model identifiability is reported in Miao et al. (2011).

The aim of kinetic modelling is, of course, not to reproduce experimental data but rather to predict patterns in a larger space of operating conditions based on models calibrated on a limited set of experimental observations. Therefore, the design space must be screened to select the most informative set of experimental conditions in which the input-output structure of models is such that parameters may be practically identified from limited kinetic data (Chis et al, 2016). Model-based Design of Experiments (MBD<sub>oE</sub>) (Franceschini and Macchietto, 2008) techniques can be used to quantify and optimise the expected information to be acquired from a set of experiments given one or more candidate kinetic models (Galvanin et al., 2016). Information metrics based on Fisher Information matrix (FIM) depends on model structure (i.e. set of rate equations), standard deviation of measurement errors and actual value of estimated kinetic parameters. These metrics are optimised to determine the experimental conditions that are mostly informative with respect to the estimation of model parameters, i.e. to minimise parametric uncertainty, and are a valid tool to bracket the most promising regions of experimental conditions to be explored when developing a kinetic model even under model uncertainty.

The objective of this study is twofold: 1) to propose a set of new kinetic models for HMF hydrogenation that includes temperature as an explicit experimental design variable; 2) to determine if and at which experimental conditions the relevant model parameters can be determined in a precise and accurate way. A three-step identifiability approach is proposed at the purpose where MBD<sub>oE</sub> techniques are exploited to assess the identifiability of candidate kinetic models of HMF hydrogenation in a batch reaction system. In the first step, candidate kinetic models are reparametrised and temperature dependency is introduced. The design space, defined by temperature, experiment duration and initial HMF and DMF concentrations as experimental decision variables, is screened using Latin Hypercube Sampling (LHS). In the second step, models are assessed in terms of model adequacy and parameter estimability. In the third, final step, a FIM-based MBD<sub>oE</sub> approach is used to identify the most informative regions for kinetic parameter estimation in the design space underlining the regions of parametric identifiability to be used in the experimentation to guarantee a precise estimation of model parameters.

This paper is structured as follow: in Section 2, an overview of the experimental studies performed to obtain the kinetic data and an overview of the kinetic models proposed in literature is given. In Section 3, the suggested MBD<sub>oE</sub> procedure is outlined and the key concepts for the identification of the kinetic models explained. In Section 4 results from the application of the proposed procedure are presented in terms of model fidelity and the experimental conditions at which the kinetic parameters for the new models can be precisely estimated. Lastly, conclusions about the study are summarised in Section 5.

## **2. Materials and Methods**

### **2.1 Review of kinetic models of HMF hydrogenation**

Grilic et al. (2014). This model was investigated by the original authors using a 300mL stainless steel cylindrical autoclave with an impeller. Fourier-Transform Infrared Spectroscopy (FTIR) was used to take online measurements. High-pressure, continuous supply of Hydrogen (8 MPa) as well as a high impeller speed (1000 rpm) was used to exclude mass transfer limitations affecting the measurement.

The kinetic constants include the internal mass transfer limitations. The experiments were conducted at 300 °C. The catalysts used were based on palladium, nickel, nickel-molybdenum and molybdenum-disulphide. The catalysts were added to a support that had previously been dried using wet impregnation. The solvents used were 2-propanol, phenol, m-cresol, anthracene, cyclohexanol, xylene, tetralin and pyridine. This power law model (1-11) lumps measured species together by oxygenated groups (carbonyl and hydroxyl) and their products (CO, CO<sub>2</sub>, H<sub>2</sub>O):

$$r_n = k_n y_i \quad (1)$$

$$y_i = \frac{c_i}{c_{OH}} (t = 0) \quad (2)$$

$$\frac{dy_{OH}}{dt} = -k_1 y_{OH} - k_2 y_{OH} - k_4 y_{CHO} \quad (3)$$

$$\frac{dy_{C=O,tot}}{dt} = k_2 y_{OH} - k_3 y_{CHO} - k_4 y_{CHO} \quad (4)$$

$$\frac{dy_{CHO}}{dt} = k_2 y_{OH} - k_3 y_{CHO} - k_4 y_{CHO} \quad (5)$$

$$\frac{dy_{CO}}{dt} = k_3 y_{CHO} \quad (6)$$

$$\frac{dy_{CO_2}}{dt} = k_4 y_{CHO} \quad (7)$$

$$\frac{dy_{H_2O}}{dt} = k_1 y_{OH} \quad (8)$$

$$\frac{dy_{CH_4}}{dt} = k_5 \quad (9)$$

$$\frac{dy_{C_3H_8}}{dt} = k_6 \quad (10)$$

$$y_{C=O,tot} = y_{CHO} + y_{C=O,ester} \quad (11)$$

In this model the relative concentrations (2) of species  $i$  is  $y_j$  [unitless], the apparent reaction rate of reaction  $n$  is  $r_n$  [min<sup>-1</sup>], the corresponding kinetic constant is  $k_n$  [min<sup>-1</sup>]. The concentrations  $c_i$  are the measured responses and kinetic constants  $k_i$  are the parameters to be estimated. Initial parameter values were estimated numerically from literature values at intervals using the Arrhenius equation, despite the possible effect of mass transfer limitations in the system.

Luo et al. (2015). This model was investigated by the original authors using a stainless steel tubular reactor, 20cm in length and 4.6mm in internal diameter. Measurements were taken using gas chromatography-mass spectrometry offline. The experiment was conducted under 33 bar and at 180

°C. Catalysts tested were platinum palladium, iridium, ruthenium, nickel and cobalt, which were added to a dried support using wet impregnation. Catalyst deactivation rates were experimentally observed and factored into the parameter calculation. Despite high pressure, the original authors acknowledge that mass transfer resistances may be present and should be factored into the kinetic constants. This model is a power law model (12-15) arising from an attempt to quantify the differences between the use of different metal catalysts. The HMF reaction was modeled as a series of first-order, sequential reactions and the rate determining step is the conversion of HMF to BHMF:

$$\frac{dC_{HMF}}{dt} = k_1 C_{HMF} \quad (12)$$

$$\frac{dC_{BMHF}}{dt} = k_1 C_{HMF} - k_2 C_{BMHF} \quad (13)$$

$$\frac{dC_{DMF}}{dt} = k_2 C_{BMHF} - k_3 C_{DMF} \quad (14)$$

$$\frac{dC_{DMTHF}}{dt} = k_3 C_{DMF} \quad (15)$$

In the model  $C_i$  [ $\text{g L}^{-1}$ ] is the concentrations of species  $i$  (measured outputs), and  $k_n$  [ $\text{min}^{-1}$ ], the kinetic constants of reaction  $n$  (to be estimated from data). The species  $i$  are: HMF, BHMF, DMF, DMTHF. Bortoli (2018) points out the structural similarity of this model to the model proposed by Gyngazova et al. (2017).

Jain and Vaidya (2016). This model was investigated using a 100mL reactor with a four-pitched blade at 45°. Measurements were taken offline using high-performance liquid chromatography, mass spectroscopy and a refractive index detector. The experiment was conducted in a temperature range of 40-70 °C and pressures of 0.69-2.07 MPa. Various initial concentrations were investigated on different loadings of ruthenium catalyst on a carbon support. Impeller speeds of 1200 rpm were used to prevent external mass transfer limitations and internal mass transfer limitations were prevented by sizing catalyst pellets according to the Weisz-Prater Criterion (Bindwal and Vaidya, 2014) as cited in Jain and Vaidya (2016). The following rate expression (16) is proposed among a number of potential expressions by the authors to derive the differential balances of each species:

$$r_n = \frac{k_3 K_{H_2} K_B C_{H_2} C_B}{(1 + K_{H_2} C_{H_2} + K_B C_B)^2} \quad (16)$$

Here,  $r_n$  [ $\text{kmol kg}_{\text{cat}}^{-1} \text{min}^{-1}$ ] is the reaction rate and  $k_n$  [ $\text{kmol kg}_{\text{cat}}^{-1} \text{min}^{-1}$ ] is a kinetic constant for the  $n$ -th surface reaction.  $C_i$  [ $\text{kmol m}^{-3}$ ] is the concentration and  $K_i$  [ $\text{m}^3 \text{kmol}^{-1}$ ] is the adsorption rate constant for species  $i$ . Species  $\text{H}_2$  is hydrogen while species B is the substrate in reaction  $n$ .  $K_i$  and  $k_n$  are parameters to be estimated from concentration data. An initial kinetic estimability analysis carried out by Bortoli (2018) reveals that the model is structurally unidentifiable from concentration data only.

Gawade et al. (2016). It was proposed based on an experiment conducted in a 100mL autoclave with pitched turbine impeller. Concentration measurements were taken using mass spectrometry and gas chromatography. The experiment was carried out at an impeller speed of 1000 rpm. The absence of external mass transfer limitations was demonstrated by lowering the impeller speed to 800rpm under constant conditions with no observable effect of the reaction. The solvent used was tetrahydrofuran. The catalyst used was a bifunctional metal-acid palladium-caesium dodeca-tungsto-phosphoric acid catalyst supported on K-10 acidic clay (2Pd-20CsDTP/K-10). The support was dried before catalyst was added to it using wet impregnation. This model (17-20) is based on the observable species participating in the reaction (subscripts: HMF, DMF, BHMF, DMTHF and  $w$  for water) and includes catalyst loading,  $w$  [ $\text{g L}^{-1}$ ], and hydrogen partial pressure,  $P_{\text{H}_2}$  [atm] explicitly.  $C_i$  [ $\text{mol L}^{-1}$ ] and  $K_i$  [ $\text{L mol}^{-1}$ ] represent the concentrations and adsorption coefficients of species  $i$  respectively, while  $k_n$  represents the kinetic constant in reaction  $n$ :

$$-\frac{dC_{\text{HMF}}}{dt} = \frac{k_1 K_{\text{HMF}} C_{\text{HMF}} \sqrt{K_{\text{H}_2} P_{\text{H}_2}} w}{[1 + \sqrt{K_{\text{H}_2} P_{\text{H}_2}} + K_w C_w][1 + K_{\text{HMF}} C_{\text{HMF}} + K_{\text{BHMF}} C_{\text{BHMF}} + K_{\text{DMF}} C_{\text{DMF}} + K_{\text{DMTHF}} C_{\text{DMTHF}}]} \quad (17)$$

$$\frac{dC_{\text{BHMF}}}{dt} = \frac{[k_1 K_{\text{HMF}} C_{\text{HMF}} - k_2 K_{\text{BHMF}} C_{\text{BHMF}}] \sqrt{K_{\text{H}_2} P_{\text{H}_2}} w}{[1 + \sqrt{K_{\text{H}_2} P_{\text{H}_2}} + K_w C_w][1 + K_{\text{HMF}} C_{\text{HMF}} + K_{\text{BHMF}} C_{\text{BHMF}} + K_{\text{DMF}} C_{\text{DMF}} + K_{\text{DMTHF}} C_{\text{DMTHF}}]} \quad (18)$$

$$\frac{dC_{\text{DMF}}}{dt} = \frac{[k_2 K_{\text{BHMF}} C_{\text{BHMF}} - k_3 K_{\text{DMF}} C_{\text{DMF}}] \sqrt{K_{\text{H}_2} P_{\text{H}_2}} w}{[1 + \sqrt{K_{\text{H}_2} P_{\text{H}_2}} + K_w C_w][1 + K_{\text{HMF}} C_{\text{HMF}} + K_{\text{BHMF}} C_{\text{BHMF}} + K_{\text{DMF}} C_{\text{DMF}} + K_{\text{DMTHF}} C_{\text{DMTHF}}]} \quad (19)$$

$$\frac{dC_{\text{DMTHF}}}{dt} = \frac{k_3 K_{\text{DMF}} C_{\text{DMF}} \sqrt{K_{\text{H}_2} P_{\text{H}_2}} w}{[1 + \sqrt{K_{\text{H}_2} P_{\text{H}_2}} + K_w C_w][1 + K_{\text{HMF}} C_{\text{HMF}} + K_{\text{BHMF}} C_{\text{BHMF}} + K_{\text{DMF}} C_{\text{DMF}} + K_{\text{DMTHF}} C_{\text{DMTHF}}]} \quad (20)$$

Kinetic constants and adsorption coefficients are parameters to be estimated from concentration measurements, while hydrogen partial pressure and catalyst loading are explicit inputs. Bortoli (2018)

has found unit inconsistencies in the definition of  $k_n$  for this model and Deussen (2019) could not reproduce the cited concentration profiles using this model.

Gyngazova et al. (2017). The experimental setup upon which the model (21-26) was developed used a 50mL stainless steel autoclave with magnetic stirrer. Samples were analysed offline using a gas chromatograph. The experiment was carried out in THF solvent on a nickel catalyst that was added in to a previously dried carbon support using wet impregnation. Temperature and pressure ranged from 150-190°C and 100-140 bar, respectively. External mass transfer limitations were avoided using stirrer speeds of 900 rpm, while internal mass transfer limitations were avoided by using small (< 50µm) catalyst pellets such that the sizing satisfied the Weisz-Prater Criterion. Furthermore, the absence of heat transfer limitations was ensured by taking samples from inside the solid catalyst and the liquid phase of the reaction to ensure that the temperature was controlled in both phases. Finally, the authors repeated the experiment with fresh and recycled catalyst and due to the low difference in these measurements concluded that the deactivation was negligible. This is, therefore, a model developed at conditions ensuring a kinetically-controlled regime. The authors derived a simple system of ordinary differential equations where the non-measurable parameters are the apparent kinetic constants of each reaction:

$$-\frac{dC_{HMF}}{dt} = k_1 C_{HMF} \quad (21)$$

$$\frac{dC_{BHMF}}{dt} = k_1 C_{HMF} - k_2 C_{BHMF} - k_5 C_{BHMF} \quad (22)$$

$$\frac{dC_{MFA}}{dt} = k_2 C_{BHMF} - k_3 C_{MFA} \quad (23)$$

$$\frac{dC_{DMF}}{dt} = k_3 C_{MFA} - k_4 C_{DMF} \quad (24)$$

$$\frac{dC_{DMTHF}}{dt} = k_4 C_{DMF} \quad (25)$$

$$\frac{dC_{DHMTHF}}{dt} = k_5 C_{BHMF} \quad (26)$$

This model describe the concentration in time of HMF, MFA, BHMF, DMF, DMTHF.  $C_i$  [mol m<sup>-3</sup>] is the concentration of species  $i$  (to be measured) while  $k_n$  [s<sup>-1</sup>] is the kinetic constant of reaction  $n$  (to be estimated from data). The model by Luo et al. (2015) is structurally similar to this model but neglects

MFA (i.e. equation (24)). Bortoli (2018) analysed this model, finding no structural identifiability issues. Deussen (2019) designed optimal experiments at the temperatures that had been experimentally investigated. Optimal designs tended towards higher initial reagent concentrations and temperatures. Due to these successful preliminary studies and the large scope of species the model encompasses, this model was chosen as a valid benchmark model upon which to build new models.

## 2.2 Main factors investigated in kinetic studies

The following factors affect the results obtained in the aforementioned kinetic studies: *i*) catalyst type and loading (resulting in potential mass transfer or kinetic limitations); *ii*) choice of solvent; *iii*) temperature; *iv*) pressure; *v*) initial reagent concentrations.

- Solvent: The choice of solvent will affect the conversion and selectivity of this three-phase reaction system. Both polar and non-polar, alcoholic and non-alcoholic, solvents may be used. Side reactions may occur with certain catalysts. Three main metrics to assess solvent impacts are suggested: the Hildebrand solubility parameter ( $\delta$ ), polarity and dielectric constant. Table S4.1 in Appendix 4 provides a visual summary of the catalysts and solvents which have been studied (Deussen, 2019). Any values of kinetic parameters will be specific to the solvents the respective experiments were conducted with, making the model comparison more challenging.
- Catalyst type and loading: Different catalysts will result in different reaction pathways. Since the studies were carried out on different catalysts, direct comparisons between kinetic parameter values are not possible. However, in all setups, a higher catalyst loading implies higher reaction speeds until a point, where higher loading results in unwanted polymerisation reactions. The maximum catalyst loading is limited by the experimental setup.
- Mass transfer limitations: This three-phase reaction can be externally mass transfer limited at the phase interface or internally in the catalyst. The external limitation can be overcome by sufficient mixing. Internal mass transfer resistances can be overcome by catalyst pellet sizing, using the Weisz-Prater Criterion (Gyngazova et al., 2017).

- Temperature: A higher temperature speeds up the main and the side reactions. For the model to account for this effect accurately, all relevant reactions must explicitly be included and the setup must indeed be kinetically controlled.
- Hydrogen pressure: The hydrogen supply pressure impacts conversion and selectivity. A pressure so high that constant hydrogen concentrations can be assumed is good for model simplification and pressure control (Bortoli, 2018). However, pressures above 30 bar may crack open Furan rings of HMF and lead to undesirable side-reactions (Hu et al., 2014).
- Initial HMF and DMF concentrations: The initial concentrations of HMF and DMF can be varied as long as certain ratios of substrate to product are maintained (Bortoli, 2018).

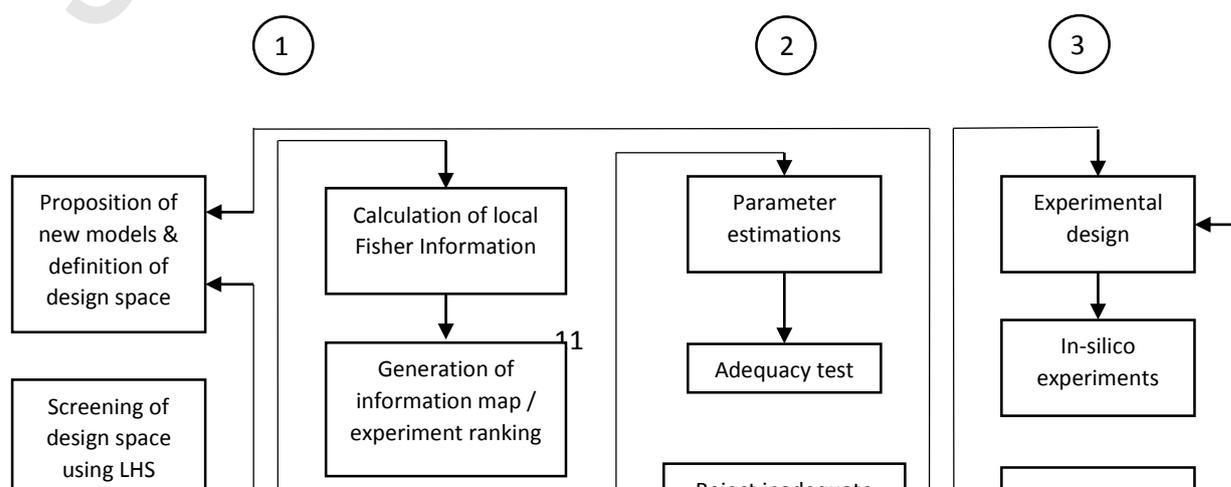
### 3. Proposed Model Identification Procedure

A sketch of the proposed model identification procedure is shown in Figure 2. The three distinct steps of the procedure are:

- 1) Initial proposition and screening of models in a preliminary design space;
- 2) Checking models for parameter estimability and adequacy;
- 3) Model-based experimental design for precise parameter estimation.

In step 1) potential kinetic models are proposed and a preliminary experimental design space is defined (i.e. space of potential variability of experimental decision variables). The design space is explored using Latin Hypercube Sampling (LHS) and in-silico experiments are carried out to assess a preliminary identifiability of model parameters based on sensitivity and correlation analysis.

Experiments are then ranked based on local Fisher information analysis (Galvanin et al., 2016). In step 2) a nonlinear parameter estimation is carried out and a-posteriori statistics, including lack-of-fit tests and assessment of kinetic parameter precision, are used to reject inadequate models.



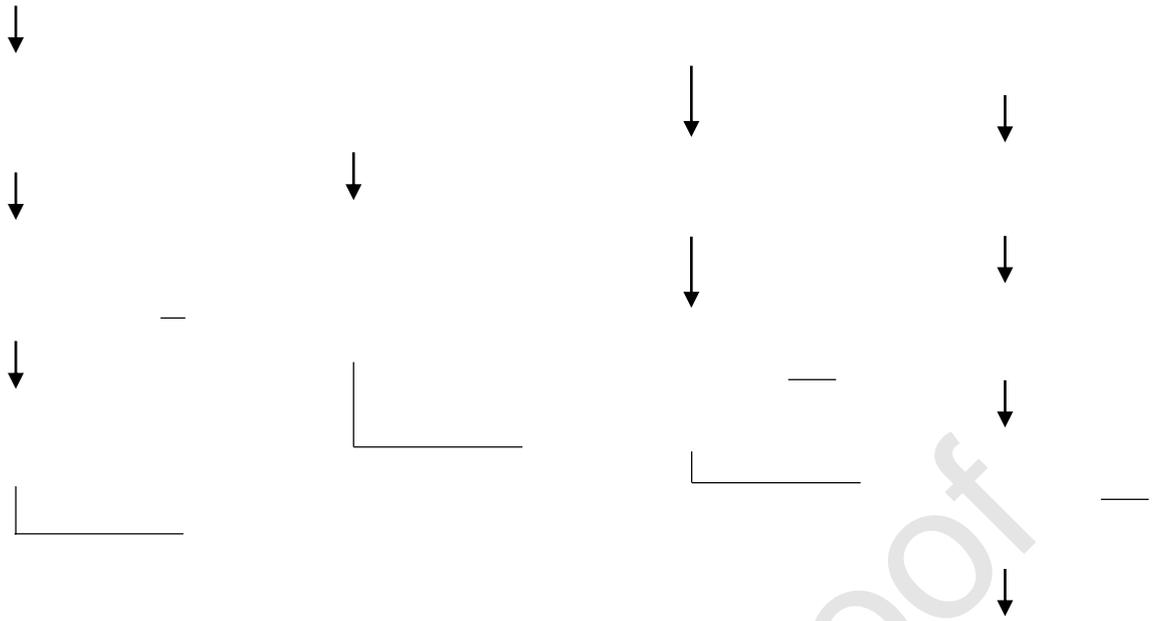


Figure 2 – Structure of the proposed model identification procedure.

In the final step 3) model-based design of experiments (MBDoE) is carried out and applied to in-silico experiments. Statistics obtained from parameter estimation are analysed to verify the impact of design on parameter estimation. The conditions that appear optimal from the information map / experiment ranking in step 1) are used as initial guess conditions to facilitate the MBDoE optimisation.

### 3.1 Initial proposition and screening of models in a preliminary design space

In the first phase of the proposed model identification procedure, kinetic models are proposed. These models are based on the kinetic models from literature (see Section 2.1 for further details). Chosen models will be modified to include temperature dependency using Arrhenius expression (27):

$$k_i = A_i e^{-\frac{E_{A_i}}{RT}} \quad (27)$$

As in multiple kinetic studies reaction rate constant values are given, these will be expressed in the form

$$\ln(k_i) = -\frac{E_{A_i}}{R} \left(\frac{1}{T}\right) + \ln(A_i) . \quad (28)$$

in order to obtain initial parameter guesses for pre-exponential factors ( $A_i$ ) and activation energies ( $E_{A_i}$ ) to initialise parameter estimation algorithms. The reparametrisation proposed by Buzzi-Ferraris

and Manenti (2009) was adopted to decrease the structural correlation associated with Arrhenius expressions by re-balancing the parameters to be estimated:

$$k_i = B_i e^{C_i \left( \frac{1}{T} - \frac{1}{T_m} \right)} \quad (29)$$

In (29)  $B_i$  [ $\text{min}^{-1}$ ] and  $C_i$  [K] are new parameters to be estimated replacing  $E_{Ai}$  [ $\text{J mol}^{-1}$ ] and  $A_i$  [ $\text{min}^{-1}$ ] and  $T_m$  [K] is an average reference temperature. Initial guesses for these new parameters may be determined from the values of the initial guesses for parameters  $A_i$  and  $E_{Ai}$  by re-arranging (29) as follows:

$$B_i = \exp \left( \ln(A_i) + C_i / T_m \right) \quad (30)$$

$$C_i = -E_{Ai} / R \quad (31)$$

The proposed kinetic models are represented by a system of differential and algebraic expressions that takes a general form (32) of state variables  $\mathbf{x}(t)$ , control variables (inputs)  $\mathbf{w}$ , and parameters to be estimated  $\boldsymbol{\theta}$ . Measured responses,  $\hat{\mathbf{y}}(t)$ , and  $\mathbf{x}(t)$  are related by function  $\mathbf{h}(\mathbf{x}(t))$  (30). The design vector  $\boldsymbol{\phi}$  (34) is located in the design space  $\Phi$  and comprises initial conditions  $\mathbf{y}^0$ , duration  $\tau$  and a set of sampling times  $\mathbf{t}_{\text{sp}}$ .

$$\mathbf{f}(\mathbf{x}(t), \mathbf{w}, \boldsymbol{\theta}, \tau) = 0 \quad (32)$$

$$\hat{\mathbf{y}}(t) = \mathbf{h}(\mathbf{x}(t)) \quad (33)$$

$$\boldsymbol{\phi} = (\mathbf{y}^0, \tau, \mathbf{w}, \mathbf{t}_{\text{sp}}) \quad (34)$$

The design space,  $\Phi$  should be defined at this stage based on the features of the experimental setup and the practical limitations offered by the equipment. It is desirable, for ease of comparison, that all models will be within the same design space, that should be contained in the range of model validity. Furthermore,  $\Phi$  should be as large as possible so that any experimental region in which the models are more informative can be captured by the model. A set of experiments must be designed to screen the entire design space (Franceschini and Macchietto, 2008). Latin Hypercube Sampling (LHS) is an exploratory design used to screen the design space by considering previously generated design points in the generation of the subsequent points. This significantly reduced the number of points needed, especially in multi-dimensional problems, to cover the entire design space. In this way the risk of

incurring in data clusters and local information optima is reduced (Montgomery, 2012). In this study the LHS points were generated using a freely available MATLAB function (Khaled, 2020). The conditions generated using the LHS function are used to run in-silico experiments. These experiments will be analysed using information metrics (32-34), which include local sensitivity (35), variance-covariance (36) and correlation (37) metrics.

The local dynamic sensitivities of the measurable outputs  $\hat{y}_i$  ( $1, \dots, n_m$ ) (i.e. set of concentrations in the specific case study) to changes in the parameters  $\boldsymbol{\theta}$  ( $1, \dots, n_\theta$ ) can be expressed by the sensitivity matrix  $\mathbf{Q}$  for all the experiments generated by LHS ( $1, \dots, n_{exp}$ ):

$$\mathbf{Q} = \begin{bmatrix} \frac{\partial \hat{y}_1}{\partial \theta_1} & \dots & \frac{\partial \hat{y}_1}{\partial \theta_{n_\theta}} \\ \vdots & \ddots & \vdots \\ \frac{\partial \hat{y}_{n_m}}{\partial \theta_1} & \dots & \frac{\partial \hat{y}_{n_m}}{\partial \theta_{n_\theta}} \end{bmatrix} \quad (35)$$

In dynamic models  $\mathbf{Q}$  theoretically exists at every point in time and can be computed to gain information at every sampling point (Franceschini and Macchietto, 2008). The  $n_\theta \times n_\theta$ -dimensional variance-covariance matrix  $\mathbf{V}^\theta$  (36) is computed from  $\mathbf{Q}$  and the standard deviation of measurement errors ( $\sigma$ ):

$$\mathbf{V}^\theta = [\mathbf{Q}^T \sigma^{-2} \mathbf{Q}]^{-1} \quad (36)$$

The correlation matrix  $\mathbf{C}$  is obtained by computing each  $kl$ -th element from the elements of  $\mathbf{V}^\theta$ :

$$C_{kl} = \frac{V_{kl}^\theta}{\sqrt{V_{kk}^\theta V_{ll}^\theta}} \quad (37)$$

Correlation coefficients  $C_{kl}$  may range from -1 to 1 (total anticorrelation to total correlation), with critical values  $> 0.95$ . Correlations close to 1/-1 indicate models where the practical estimability of model parameters is difficult.

The Fisher Information Matrix (FIM) (38), which is an approximate estimate of the Hessian of the log-likelihood function of parameters  $\boldsymbol{\theta}$ , is obtained from the elements of  $\mathbf{Q}$  and  $\sigma$ :

$$\mathbf{H} = \sum_{i=1}^{n_{exp}} \sum_{j=1}^{n_m} \left[ \frac{1}{\sigma_{ij}^2} Q_{ij}^T Q_{ij} \right] \quad (38)$$

The FIM (38) can be used as a scalar measure of information, e.g. by taking suitable metrics  $\psi$  (trace, determinant, maximum eigenvalues). This FIM metric,  $\psi(\mathbf{H}^T)$ , may then be plotted on a

multidimensional graph visualising the information of design vectors sampled from  $\phi$  using LHS. If more than two design variables have been considered, multiple graphs will be drawn to visualise the distribution of information in the full design space. The individual points may either represent experimental conditions, in which case the graph will show an experiment ranking, or represent a continuous information surface, as in Deussen and Galvanin (2021).

When there are high correlations between model parameters ( $C_{kl}$  terms in (37)) and low sensitivities (low values of sensitivity elements in (35)) the FIM approaches singularity and the model becomes non-identifiable. However, it may be the case that a model contains “sloppy” parameters rather than that it is globally structurally unidentifiable (Chis, 2016). Models are affected by sloppiness if parameters can change by order of magnitude without significantly altering the predictive behaviour of the model. From a quantitative point of view, a model is said to contain sloppy parameters when there are differences between the eigenvalues of a FIM that exceed three orders of magnitude. However, the model may still be practically identifiable if there are experimental conditions under which the sloppy parameters can be identified. Sloppiness is a measure of the ratio of the axes of the confidence area, while the identifiability is better described as the area (i.e. the determinant of variance-covariance). Sloppy models are common in biological systems or systems involving organic compounds and is hence very relevant to this investigation (Chis et al., 2016). Should the model be structurally unidentifiable in the space investigated, for example due to high parameter correlations, it must be reparametrised, reformulated or rejected at this point in the procedure.

### 3.2 Checking models for parameter estimability and adequacy

The parameters are determined using the parameter estimation function in gPROMS, which conducts nonlinear parameter estimation using maximum likelihood methods. Parameter estimation results are given in terms of statistical model adequacy and parameter precision. Relevant a-posteriori statistics are *i*) the Chi-Square ( $\chi^2$ ) test for model adequacy; *ii*) Student’s t-test for parameter precision. The  $\chi^2$  test is passed when the sum of weighted residuals (39) is smaller than a reference  $\chi^2$  value taken from a tabulated  $\chi^2$  distribution with  $(n_m \cdot n_{sp} - n_{\theta})$  degrees of freedom:

$$\chi^2 = \sum_{i=1}^{n_m} \sum_{j=1}^{n_{sp}} \frac{(y_{ij} - \hat{y}_{ij})^2}{\sigma_i^2} \quad (39)$$

The Student's  $t$ -test is passed when the  $t$ -values (40) obtained from parameters are larger than a tabulated reference  $t$ -value, i.e.

$$t_i = \frac{\hat{\vartheta}_i}{\sqrt{V_{ii}^{\vartheta}}} > t(1 - \alpha/2) \quad i = 1, \dots, n_{\vartheta} \quad (40)$$

In (40)  $t(\cdot)$  is the  $t$ -value obtained from a Student's distribution with degree of freedom equal to  $n_{exp} n_m - n_{\vartheta}$  and level of confidence given by the argument in brackets (Asprey and Naka, 1999).

The chosen confidence interval in this study is 95% ( $\alpha = 0.05$ ). If conditions (40) are satisfied for all parameters this can be interpreted as an index of satisfactory parameter precision. Statistical tests (39,40) are applied to all experiments designed using the LHS-generated conditions. Conclusions about the model can then be drawn based on if and where the model fulfils both these criteria, i.e. if the tests fail considering all experiments throughout the design space, the model must be rejected.

### 3.3 Experimental design for precise parameter estimation

If the model is not rejected (i.e. is identifiable), optimal experiments may be designed. In designs aiming at precise parameter estimation, the objective function should be a suitable measure of the information contained in each experiment. The FIM (38) is used as a basis for developing the objective function. The design vector  $\phi$  contains the variables that can be adjusted freely. However, the matrix  $\mathbf{H}$  must be converted to a scalar before optimisation. Many optimisation criteria have been proposed to this end, including criteria for parallel and sequential designs (Galvanin et al., 2006).

Popular criteria are:

- A-optimality maximises the trace of the FIM:

$$Obj = \max (\text{tr} (\mathbf{H} (\phi, \vartheta))) \quad (41)$$

- E-optimality maximises the smallest Eigenvalue of the FIM:

$$Obj = \max (\min (\lambda (\mathbf{H} (\phi, \vartheta)))) \quad (42)$$

- D-optimality maximises the determinant of the FIM:

$$Obj = \max |\mathbf{H} (\phi, \vartheta)| \quad (43)$$

- SV-optimality maximises a selected eigenvalue of the FIM:

$$Obj = \max (\lambda_i (\mathbf{H} (\boldsymbol{\phi}, \boldsymbol{\theta})) \quad . \quad (44)$$

The computer tools used to carry out MBDoe was gPROMS Process Builder (Process Systems Enterprise, 2019) version 1.4.0 for model validation, parameter estimation and optimal experimental design. Matlab version R2018b was used to perturb the parameters used for in-silico experiments, to introduce error to measurements taken from in-silico and to generate LHS points. The optimisation criteria that are supported in gPROMS Process Builder and which will be used in this investigation are A-, E-, and D-optimal designs (Franceschini and Macchietto, 2008). Once the experiments have been designed, in-silico experiments are run at the determined conditions to generate new data points. In these in-silico experiments, unlike in the screening experiments from step 1, the parameter values are perturbed from literature values to check whether the model can produce statistically meaningful results when checked against a real system, i.e. to evaluate the impact of parametric mismatch on the acquired information. This perturbation is a uniform error distribution added using the ‘rand’ function in MATLAB. The model generates in-silico concentration profiles that are further perturbed with gaussian noise to reflect the measurement error.

As an additional part of the third step of the procedure, it is beneficial to rank the experiments that were designed from the LHS using one of the optimality metrics. One can compute the trace of the FIM from each experiment in the screening phase and build “information maps”, i.e. visualisation of the information in the space of experimental decision variables. Conditions corresponding to the highest amount of information will then be used to initialise MBDoe algorithms that will operate on the variance-covariance matrix of model parameters. If the function is convex, a FIM-based experiment ranking can easily indicate where global and local information maxima are located and guide the optimal experimental design (Galvanin et al., 2016). Furthermore, this approach provides a simple graphical representation to guide future experiments by visualising where high information experiments and potential discontinuities are expected.

#### 4. Results

As a result of the previous studies in literature (Bortoli, 2018; Deussen, 2019), the model (45-50) by Gyngazova et al. (2017) is chosen as a basis to propose new model configurations where temperature

dependency is made explicit. The model identification procedure will then be applied to each distinct model configuration: M1, M2 and M3. Common to all configurations:

$$-\frac{dC_{HMF}}{dt} = k_1 C_{HMF} \quad (45)$$

$$\frac{dC_{BHMF}}{dt} = k_1 C_{HMF} - k_2 C_{BHMF} - k_5 C_{BHMF} \quad (46)$$

$$\frac{dC_{MFA}}{dt} = k_2 C_{BHMF} - k_3 C_{MFA} \quad (47)$$

$$\frac{dC_{DMF}}{dt} = k_3 C_{MFA} - k_4 C_{DMF} \quad (48)$$

$$\frac{dC_{DMTHF}}{dt} = k_4 C_{DMF} \quad (49)$$

$$\frac{dC_{DHMTHF}}{dt} = k_5 C_{BHMF} \quad (50)$$

The resulting model configurations differ from the original Gyngazova et al. (2017) model by the following:

- M1: inclusion of temperature dependency through the Arrhenius equation

$$k_i = A_i e^{-\frac{E_{Ai}}{RT}} \quad (51)$$

Set of parameters to be estimated:  $A_i, E_{Ai}$

- M2: inclusion of temperature dependency through reparametrised Arrhenius equation

$$k_i = B_i e^{C_i \left( \frac{1}{T} - \frac{1}{T_m} \right)} \quad (52)$$

Set of parameters to be estimated:  $B_i, C_i$

- M3: M1 model reduction considering the estimation of activation energies only

$$k_i = A_i e^{-\frac{E_{Ai}}{RT}} \quad (53)$$

Set of parameters to be estimated:  $E_{Ai}$

Since the new model configurations will only affect the parameters to be estimated and not the measured outputs, a common design space  $\phi$  can be defined for all of them. The design vector  $\phi$  consisted of initial HMF and DMF concentrations, temperature and experiment duration, i.e.

$$\phi = [C_{HMF}(0) \quad C_{DMF}(0) \quad T \quad \tau]^T \quad (54)$$

The design space is defined by the following lower/upper bounds on design variables

- Initial concentration of HMF ( $C_{HMF}(0)$ ):  $0.01 - 0.2 \text{ g L}^{-1}$ ;

- Initial concentration of DMF ( $C_{\text{DMF}}(0)$ ): 0.01 – 0.2 g L<sup>-1</sup>;
- Reaction temperature ( $T$ ): 300 – 600 K;
- Experiment duration ( $\tau$ ): 60 – 360 min.

Ranges for temperature ( $T$ ) and experiment duration ( $\tau$ ) are defined by upper and lower bounds considering conditions that are achievable in the actual experimental setup, as well as to limit the formation of undesired by-products towards both ends of the range and to avoid slow reactions paces that are impractical for experimentation at the lower bound. The same range of initial concentrations previously proposed in literature (Gawade et al., 2016 & Gyngazova et al., 2017) was adopted in this study. Under these conditions the catalyst to substrate ratio remains similar (as there is no loading term in the model) and the hydrogen concentration can safely be assumed to be well in excess. It is assumed that 30 sampling points can be taken and that they are evenly spaced in time.

#### 4.1 Model configuration 1 (M1)

##### 4.1.1 Step 1: Initial proposition and screening of models in a preliminary design space

The original parameters  $k_i$  are now replaced by Arrhenius parameters  $A_i$  and  $E_{A_i}$  (278). This doubles the number of kinetic parameters to be estimated. For the initial guesses, the Arrhenius equations were re-written (289) and a straight line fitted to the experimental data points from Gyngazova et al. (2017). Fitting results and initial guess values for kinetic parameters are reported, respectively, in Figure S1 and Table S1.1 of Appendix 1. M1 adequately represents the physical system as its behaviour reflected the experimental results at the temperatures that were investigated originally in Gyngazova et al. (2017). In-silico data were generated using LHS conditions from Table S1.2 and the initial guesses from Table S1.1 to determine whether the models are practically identifiable. The standard deviation used to generate the concentration data was  $\sigma = 0.003$  [g L<sup>-1</sup>] based on experimental evidence from other hydrogenation reactions in similar setups (Bindwal and Vaidya, 2013 & 2014).

Table 1 – Correlation matrix for M1: results show the correlations between the pre-exponential factors and their respective activation energies; critical correlations are indicated in red.

$\vartheta$	#	1	2	3	4	5	6	7	8	9	10
$A_1$	1	1									

A <sub>2</sub>	2	-0.18	1								
A <sub>3</sub>	3	-0.23	-0.35	1							
A <sub>4</sub>	4	-0.03	0	-0.06	1						
A <sub>5</sub>	5	-0.17	<b>0.95*</b>	-0.34	0.04	1					
E <sub>A1</sub>	6	<b>1*</b>	-0.18	-0.22	-0.03	-0.17	1				
E <sub>A2</sub>	7	-0.19	<b>1*</b>	-0.35	0	<b>0.95*</b>	-0.19	1			
E <sub>A3</sub>	8	-0.22	-0.35	<b>1*</b>	-0.06	-0.34	-0.22	-0.35	1		
E <sub>A4</sub>	9	-0.03	0	-0.05	<b>1*</b>	0.04	-0.03	0	-0.05	1	
E <sub>A5</sub>	10	-0.17	0.94	-0.33	0.04	<b>1*</b>	-0.18	0.94	-0.33	0.04	1

The correlation analysis carried out on M1 (Table 1) shows several critical correlations (in red) as effect of the standard Arrhenius formulation. This model should, therefore, be rejected or reformulated. In addition, the information content obtained from pre-exponential factors is very low compared to the activation energies, see Figure S1.1 in Appendix 2. This has a critical impact on the information content of experiments. The most informative experiments have both a narrow information gap between pre-exponential constants and activation energies and generate a significant amount of information related to activation energies. This is realised, for example, in experiments 1 and 9-12. From the ranking of experiments (Figure 3) the shape of the objective function can be deduced in the space of operating conditions. There is an evident information peak at the conditions of experiment 1, the conditions of which are:  $C_{\text{HMF}}(0) = 0.105 \text{ g L}^{-1}$ ,  $C_{\text{DMF}}(0) = 0.193 \text{ g L}^{-1}$ ,  $T = 511 \text{ K}$  and  $\tau = 229 \text{ min}$ . The parameter estimation tests (see Section 4.1.2) will be carried out despite model rejection to investigate the estimability behaviour of M1 and the effect of parameter correlations.

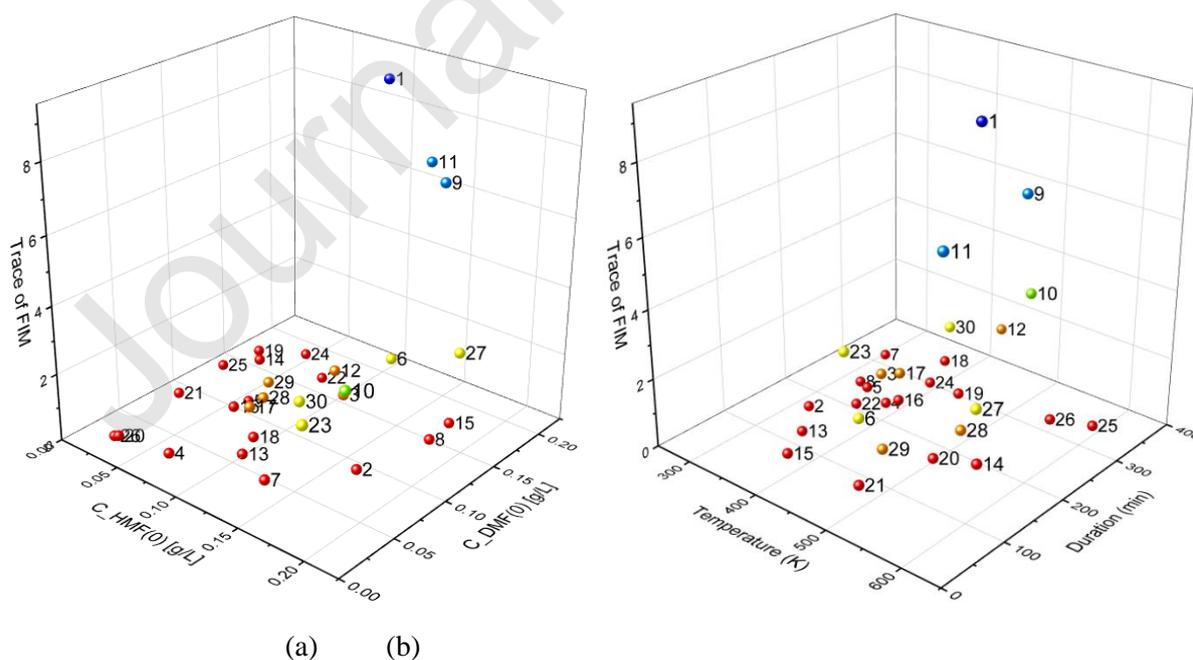


Figure 3 – Experiments ranked by trace of FIM for M1 as a function of (a) initial reagents concentration; (b) experiment duration and temperature. Blue indicates high information level; red indicates low information level.

#### 4.1.2 Checking models for parameter estimability

M1 showed high parameter correlation as part of step 1 tests. However, parameter estimation results (Table 2) show that this model passed both the  $\chi^2$  test ( $\chi^2 = 4944.70$ ;  $\chi^2_{\text{critical}} = 5561.91$ ) and Student's t-test. These can be considered surprising results and worth closer scrutiny, since a problem of high correlation should usually result in poor parameter estimation. As shown in Table 1, there are critical correlations between the pre-exponential factors and the activation energies as well as  $A_2$  and  $A_5$  and  $A_5$  and  $E_{A2}$ .  $E_{A2}$  and  $E_{A5}$  have an almost critical correlation of 0.94. Since a wide design space has been investigated with many experiments and a high number of sampling points, information from additional experiments would most not help on reducing parameter correlation. A structural identifiability issue is present for this model where the correlations persist for any potential combination of design variables in the design space (Chis et al., 2016). An additional issue is the low sensitivity of the measured concentrations to the pre-exponential factors (several orders of magnitude lower contribution of all  $A_i$  to Fisher information, as illustrated in Figure S1.1 in Appendix2). This is an often-encountered problem in Arrhenius expressions (Schwaab et al., 2008).

Table 2 – Parameter estimation results including estimated values, initial guesses and 95% confidence t-values for M1.

$\vartheta$	Final Value	Initial Guess	95% t-value	$\vartheta$	Final Value	Initial Guess	95% t-value
$A_1$	$2.75 \times 10^6$	$1.62 \times 10^6$	14.69	$E_{A1}$	68.87	58.79	321.01
$A_2$	$2.34 \times 10^7$	$5.34 \times 10^6$	4.36	$E_{A2}$	71.23	64.48	100.51
$A_3$	$1.36 \times 10^5$	$5.16 \times 10^5$	2.08	$E_{A3}$	48.91	42.00	33.01
$A_4$	$8.00 \times 10^5$	$6.35 \times 10^5$	14.17	$E_{A4}$	77.34	83.61	250.58
$A_5$	$5.06 \times 10^6$	$3.13 \times 10^5$	4.15	$E_{A5}$	74.28	62.17	99.02
Reference value for 95%:			1.65	Reference value for 95%:			1.65

Due to the low sensitivities and extremely high parameter correlations, the estimated values for  $A_i$  can become unrealistically large (Table 2) without particularly affecting model predictions, and thus the standard deviations of kinetic model parameters appearing at the denominator of Student's t-values

(404) become very small by comparison. Thus, the t-test is not reliable in this case and the model must be rejected for further scrutiny due to structurally high correlations.

## 4.2 Model configuration 2 (M2)

### 4.2.1 Step 1: Initial proposition and screening of models in a preliminary design space

It is apparent that the introduction of  $T_m$  in (30) introduces a new degree of freedom for M2. Initially,  $T_m$  was fixed in the middle of the design range at 450K. However, when applied to represent the system behaviour in the preliminary design space, the resulting discrepancy between the concentration profile predicted by M2 and the original Gyngazova et al. (2017) model used to generate in-silico experimental measurements was so large that the effect of  $T_m$  was further investigated. It was found that fixing  $T_m$  at the experimental temperature provided the best fit to predictions from the original Gyngazova et al. model across the remaining temperature range (except near the discontinuity where  $T = T_m$ ). Even still, there was a discrepancy of about an order of magnitude between these optimal initial guesses and the experimentally determined kinetic constants for all reactions in the investigated set of experimental temperatures. From this we can conclude that M2 is an inaccurate description of the physical system in the preliminary design space thus needs to be rejected at the model proposition step. However, since M2 adequately represents the behaviour of the system around a reference temperature point, M2 might represent the system better in a reduced temperature range. This is supported by the fact that M2 kinetic parameter estimates approach the experimentally determined values as  $T_m$  approaches the temperature at which the experimental points were taken. As a result, step 1 and 2 of the proposed procedure will be carried out to verify whether M2 succeeds at reducing the correlation between parameters, even though in the preliminary design space M2 does not represent the system well. Results may be useful for future work to verify the model adequacy in a smaller temperature range around optimal design points. As in M1, to carry out the preliminary analysis, a design space has been defined and LHS points generated. The design space and therefore the points generated using the LHS technique for the initial screening are the same as for M1 and can be found in Appendix 1 (Table S1.2). As was done in M1, noise was added to the experimental points before

the parameter estimation. A  $T_m$  value of 450K was chosen as the average reference temperature, a value located in the middle of the investigated temperature range. The simulations for experiments 11, 14, 25, 27 & 28 failed to converge numerically during simulation, due to the particularly high residuals obtained at the upper bound of the temperature range, with all the failed calculations clustering around 550-600 K.

*Table 3 - Correlation matrix for M2: results show the correlations between the pre-exponential factors and their respective activation energies; critical correlations are indicated in red.*

9	#	1	2	3	4	5	6	7	8	9	10
B <sub>1</sub>	1	1									
B <sub>2</sub>	2	-0.23	1								
B <sub>3</sub>	3	-0.24	-0.38	1							
B <sub>4</sub>	4	0	0.01	0.02	1						
B <sub>5</sub>	5	-0.12	0.84	-0.29	-0.11	1					
C <sub>1</sub>	6	-0.39	0.14	0.13	0.01	0.12	1				
C <sub>2</sub>	7	0.15	-0.39	0.16	-0.03	-0.33	-0.32	1			
C <sub>3</sub>	8	0.1	0.11	-0.22	-0.01	0.11	-0.35	-0.4	1		
C <sub>4</sub>	9	0.03	0.03	0.04	0.54	-0.06	-0.07	-0.05	-0.15	1	
C <sub>5</sub>	10	0.15	-0.39	0.16	-0.03	-0.33	-0.32	<b>1</b>	-0.4	-0.05	1

The preliminary analysis shows that M2 succeeds at greatly reducing the correlation between parameters (Table 3) despite there is one critical correlation remaining between  $C_2$  and  $C_5$ . In addition, the contributions to the overall information of the parameter groups  $B_i$  and  $C_i$  are more regularly distributed as illustrated in Figure S2.2 of Appendix 2. The FIM traces for these groups of parameters are about 10 orders of magnitude apart in Figure S2.2 (M2) as opposed to ca. 20 orders of magnitude in Figure S2.1 (M1). However, the distinctive information behaviour of the two groups of parameters can still be seen clearly (Figure S2.2).

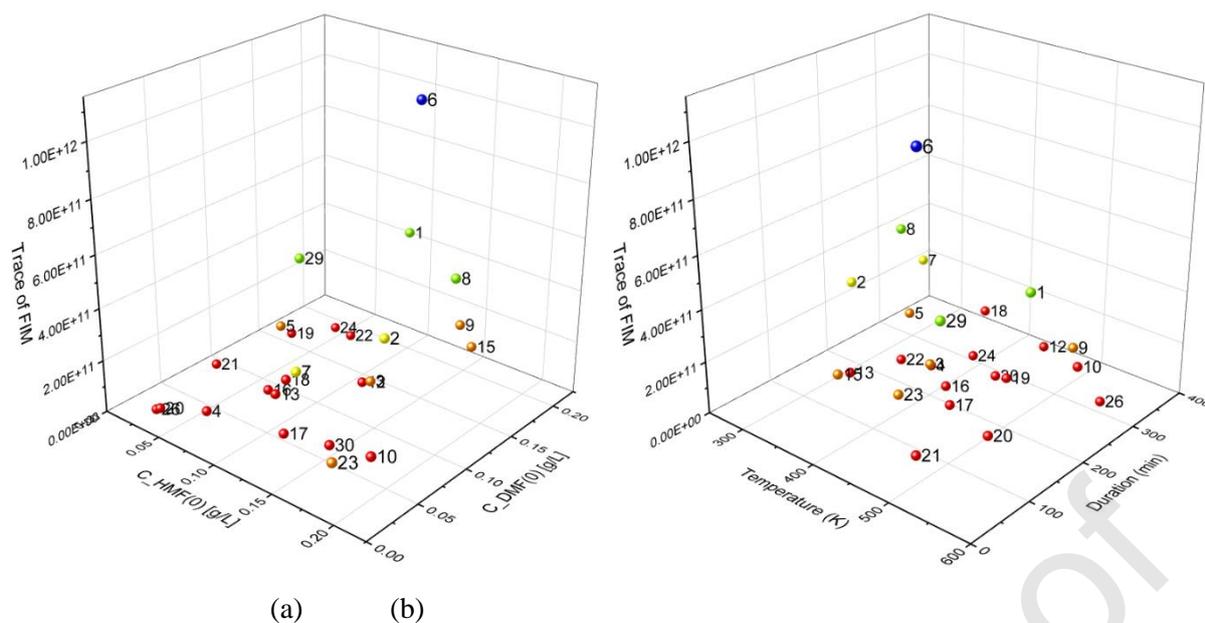


Figure 4 – Experiments ranked by trace of FIM for M2 as a function of (a) initial reagents concentration; (b) duration and temperature. Blue indicates high information levels; red indicates low levels of information.

The most informative region of the design space (Figure 4) is towards higher initial reagent concentrations ( $C_{\text{HMF}}(0) \approx 0.10\text{-}0.16 \text{ g L}^{-1}$ ,  $C_{\text{DMF}}(0) \approx 0.10\text{-}0.20 \text{ g L}^{-1}$ ) (Figure 4a), average temperature and short experimental duration (Figure 4b). The most informative experiment (experiment 6) is characterised by the following conditions:  $C_{\text{HMF}}(0) = 0.152 \text{ g L}^{-1}$ ,  $C_{\text{DMF}}(0) = 0.140 \text{ g L}^{-1}$ ,  $T = 486.15 \text{ K}$  and  $\tau = 67 \text{ min}$ .

#### 4.2.2 Step 2: Checking models for parameter estimability

M2 is adequate to represent the system only in a design space where high temperature experimental points (i.e. points in the temperature range 550-600 K) are removed ( $\chi^2 = 3910.13$ ;  $\chi^2_{\text{critical}} = 4647.00$ ). In this range of experimental conditions a precise estimation of kinetic parameters can be achieved (Table 4). However, the stationarity criterion for the log-likelihood in parameter estimation is not fulfilled in these parameter estimates, indicating very steep slopes of the objective function, i.e. due to high sensitivities, small changes in conditions can result quickly in the non-stationarity of relevant covariance values.

Table 4 – Parameter estimation results including estimated values, initial guesses for M2 and 95% confidence intervals excluding outliers in temperature range 550-600 K.

$\vartheta$	Final Value	Initial Guess	95% t-value	$\vartheta$	Final Value	Initial Guess	95% t-value
B <sub>1</sub>	0.03	0.028	111.05	C <sub>1</sub>	-8147.46	-8078.42	37.13
B <sub>2</sub>	0.11	0.12	30.80	C <sub>2</sub>	-7399.94	-7996.39	10.24
B <sub>3</sub>	0.28	0.32	13.37	C <sub>3</sub>	-5289.37	-5443.95	6.06
B <sub>4</sub>	$0.90 \times 10^{-3}$	$0.90 \times 10^{-3}$	138.29	C <sub>4</sub>	-9222.75	-9225.76	481.57
B <sub>5</sub>	0.01	0.01	25.95	C <sub>5</sub>	-6973.79	-7634.95	8.48
Reference value for 95%:			1.65	Reference value for 95%:			1.65

The experimental ranking (Figure 4) confirms a very steep variation of  $\text{tr}(\mathbf{H})$  between different experimental conditions. The very high values of Fisher information obtained are cause of concern for a reliable quantification of information. A potential explanation is the extremely high sensitivities observed for parameters  $C_i$  (Figure S2.2 in Appendix 2), which lead to high values of information metrics, but also to a more erratic behaviour in the overall information. This is clear by comparing the distribution of information of  $\text{tr}(\mathbf{H})$  in Figure S2.2 (M2) and the one realised in Figure S2.1 (M1). In the latter case there is a limited number of experiments characterised by a very low information level. The model does succeed at balancing the information that can be obtained on model parameters and on reducing parameter correlation and can be promising for further validation studies in a tighter temperature range.

### 4.3 Model configuration 3 (M3)

#### 4.3.1 Step 1: Initial proposition and screening of models in a preliminary design space

Thus far, M1 and M2 have been rejected for further analysis. M1, however, is a better representation of the physical system in a larger experimental design space. The main problem is the high correlations of pre-exponential factors and activation energies. A sensitivity analysis of the two groups of parameters ( $A_i$  and  $E_{ai}$ ) has been carried out to see if either of them or at least some individual parameters can be excluded from parameter estimation. Results are illustrated in Appendix 3. The analysis was carried out at fixed operating conditions ( $T = 423$  K,  $C_{\text{HMF}}(0) = 0.0667$  g L<sup>-1</sup>,  $C_{\text{DMF}}(0) = 0.0222$  g L<sup>-1</sup>) where the pre-exponential factors were subjected to changes of +20% while keeping the activation energies constant (Table S3.1). As shown in Figure S3.1 the pre-exponential factor has a limited effect on the concentration profile of the reaction as compared to the activation energy for the same relative change in the model parameters. As a consequence, a new model (M3) is proposed where the pre-exponential factors are fixed at the values estimated in M1 and the activation

energies are estimated, so that the temperature dependency is still maintained. A preliminary analysis of M3 was carried out using the same 30 experiments generated using LHS. Results show that, in this case, there are no significant correlations between model parameters (Table 5).

Table 5 – M3: correlation matrix obtained after parameter estimation.

9	#	1	2	3	4	5
E <sub>A1</sub>	1	1.00				
E <sub>A2</sub>	2	-0.17	1.00			
E <sub>A3</sub>	3	-0.17	-0.22	1.00		
E <sub>A4</sub>	4	0.05	0.05	-0.29	1.00	
E <sub>A5</sub>	5	-0.05	0.75	-0.14	0.06	1.00

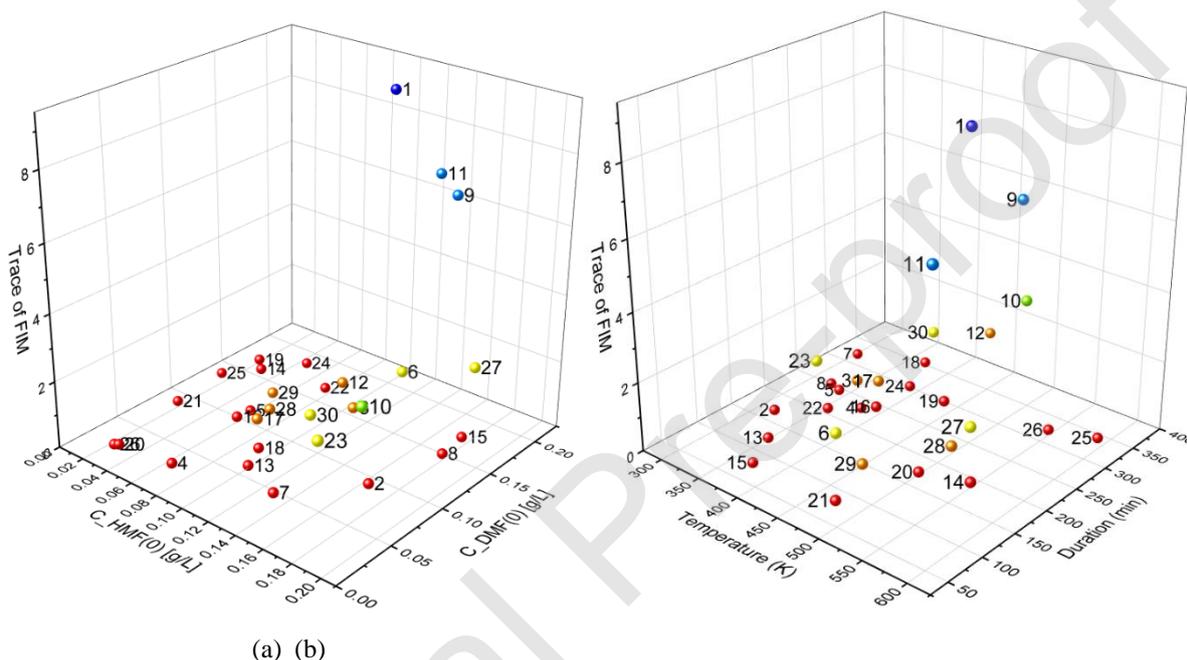


Figure 5 – Experiments ranked by trace of FIM for M3 as a function of (a) initial reagents concentration; (b) experiment duration and temperature. Blue indicates high information; red indicates low information.

#### 4.1.2 Step 2: Checking models for parameter estimability

As can be seen from Table 6, results are satisfactory for the estimation of all the parameters except  $E_{A3}$ , a parameter affected by a large uncertainty. The model passed the  $\chi^2$  test overall ( $\chi^2 = 5220.13$ ;  $\chi^2_{\text{critical}} = 5566.99$ ) but has some limitations on representing experiments with particularly low MFA concentrations. This model is practically identifiable despite the results for  $E_{A3}$ , as experiments with particular low MFA concentrations produced a very low amount of information to estimate this specific parameter. The Fisher Information analysis (Figure S2.3 in Appendix 2) as well as the ranking of experiments (Figure 5) allow an overview of this initial assessment. It is important to note

that there is very little variation in the overall information between these information metrics (M3) and the metrics that included the pre-exponential factors (M1). By comparing Figure 5 (M3) and Figure 3 (M1) the information levels and distribution obtained from experiment rankings clearly show that only negligible information was lost by omitting the pre-exponential factors (value of  $\text{tr}(\mathbf{H})$  is around 9 for experiment 1 in both cases). This confirms that the sensitivities of the pre-exponential factors and the corresponding low information were causal to the problems in M1. Thus, M3 presents a solution to the problem of parametric identifiability and will be considered in further analysis, where optimal experiments can be designed by MBDoE to improve the confidence on the relevant kinetic parameters.

#### 4.3.3 Step 3: Model-based experimental design for precise parameter estimation

A model-based optimal experimental design has been carried out to identify the most informative experimental conditions for the estimation of critical parameters. As initial guesses for the parameter values, the values determined by linear regression from the original experimental points published by Gyngazova et al. (2017) are used as “true” underlying parameters for the experiment against which the model will be checked, a random error of within 15% of the original parameter guess has been added to verify the robustness of the proposed optimal experimental design. The five most informative experiments as obtained from the ranking of LHS experiments are summarised in Table S3.2 in Appendix 3. The ranking showed that experimental design conditions from experiments 1 and 11 are the most informative (as illustrated in Figures 5a and 5b) with experiment 11 providing an improved balance between the information content related to multiple parameters (see Figure S2.3 in Appendix 2). Since these two experiments are similarly informative and close in the experimental design space, they yield very similar results when used to initialise the MBDoE optimisation. Experiments 23 and 30 are characterised by a lower level of information, and lie on a different region of the information space as deduced from the rankings (Figures 5a and 5b, yellow points) and might thus lead to local optimum when used to initialise the MBDoE optimisation. An experiment design was carried out using a D-optimal design criterion as the objective function. The initial guess

conditions were those suggested by experiment 1 from the LHS screening (Table S3.2) and the resulting optimally designed conditions were the following:

- $C_{HMF}(0)$ : 0.200 g L<sup>-1</sup>;
- $C_{DMF}(0)$ : 0.198891 g L<sup>-1</sup>;
- $T = 455.42$  K;
- $\tau = 594$  min ( $\sim 10$  h).

Interestingly, MBDoE tends to drive the conditions towards high initial concentrations for HMF and DMF and temperatures positioned in the middle of the experimental design space (300 – 600 K). Even if not shown for the sake of conciseness, it has been verified that different MBDoE experimental design criteria (E-, A- optimal) yielded very similar results in terms of optimal experimental design conditions.

Using these D-optimally designed conditions, a further in-silico experiment and a subsequent parameter estimation were carried out using only the data generated from this experiment. Parameter estimation results are shown in Table 6. The full set of model parameters could be estimated with satisfactory statistics in terms of parameter precision. After the D-optimal designed experiment, moderate parameter correlations have been observed (Table S3.2 in Appendix 3), and the model is adequate to represent the concentrations for all the chemical species as illustrated by the  $\chi^2$  test results. An information map showing the regions of maximum information is shown in Figure 6. Results are in line with the prediction of information from the LHS as they are in the area of the design space where a high information cluster is observed (dark red area in the figure). A second, local, maximum for information is present at temperatures around 330 K and initial concentrations in the middle of the design space ( $C_{DMF}(0)$  around 0.10), but a significantly lower amount of information is generated at these conditions.

It is important to stress that one, single, optimally designed experiment provided better statistics than all of the experiments combined that were used in the LHS-based parameter estimation. This is likely due to the fact that very low MFA concentrations, which make  $E_{A3}$  particularly difficult to estimate, are realised in several regions of the design space sampled by LHS and tend to provide a significantly

low level of information for the estimation of this critical parameter. This is clearly indicating that, in the specific case study, the explorability of the design space might simply lead to a waste of time and analytical resources, while an MBDoE approach can be more effective for kinetic model identification.

Table 6 – Parameter estimation results for M3 from LHS experiments and after the D-optimally designed in experiment including estimated values, initial guesses, 95% confidence t-values and  $\chi^2$  test statistics (reference  $\chi^2$  values are indicated in brackets, \*\* indicates parameters failing t-test, \* indicates responses failing  $\chi^2$  test).

9	LHS Experiments			D-optimal design		
	Final Value	Initial Guess	95% t-value	Final Value	Initial Guess	95% t-value
E <sub>A1</sub>	66.99	58.79	1939.57	59.31	58.79	99.33
E <sub>A2</sub>	66.49	64.48	676.13	65.16	64.48	174.59
E <sub>A3</sub>	6.57**	42.00**	0.001**	42.35	42.00	39.19
E <sub>A4</sub>	76.36	83.61	3772.45	83.75	83.61	636.16
E <sub>A5</sub>	63.86	62.17	492.01	62.64	62.17	149.84
Reference t-value (95%):			1.65	Reference t-value (95%): 1.654		
<b><math>\chi^2</math> test statistics</b>						
	LHS Experiments			D-optimal design		
$\chi^2$ - Total	5220.13 (5566.99)			$\chi^2$ - Total	135.14 (206.87)	
$\chi^2$ - C <sub>BHMF</sub>	783.94 (1113.39)			$\chi^2$ - C <sub>BHMF</sub>	11.45 (37.65)	
$\chi^2$ - C <sub>MFA</sub>	1260.58* (1113.39)			$\chi^2$ - C <sub>MFA</sub>	19.9 (37.65)	
$\chi^2$ - C <sub>DHMTFH</sub>	914.26 (1113.39)			$\chi^2$ - C <sub>DHMTFH</sub>	34.64 (37.65)	
$\chi^2$ - C <sub>DMF</sub>	825.82 (1113.39)			$\chi^2$ - C <sub>DMF</sub>	33.66 (37.65)	
$\chi^2$ - C <sub>DMTHF</sub>	874.21 (1113.39)			$\chi^2$ - C <sub>DMTHF</sub>	21.88 (37.65)	
$\chi^2$ - C <sub>HMF</sub>	861.32 (1113.39)			$\chi^2$ - C <sub>HMF</sub>	13.61 (37.65)	

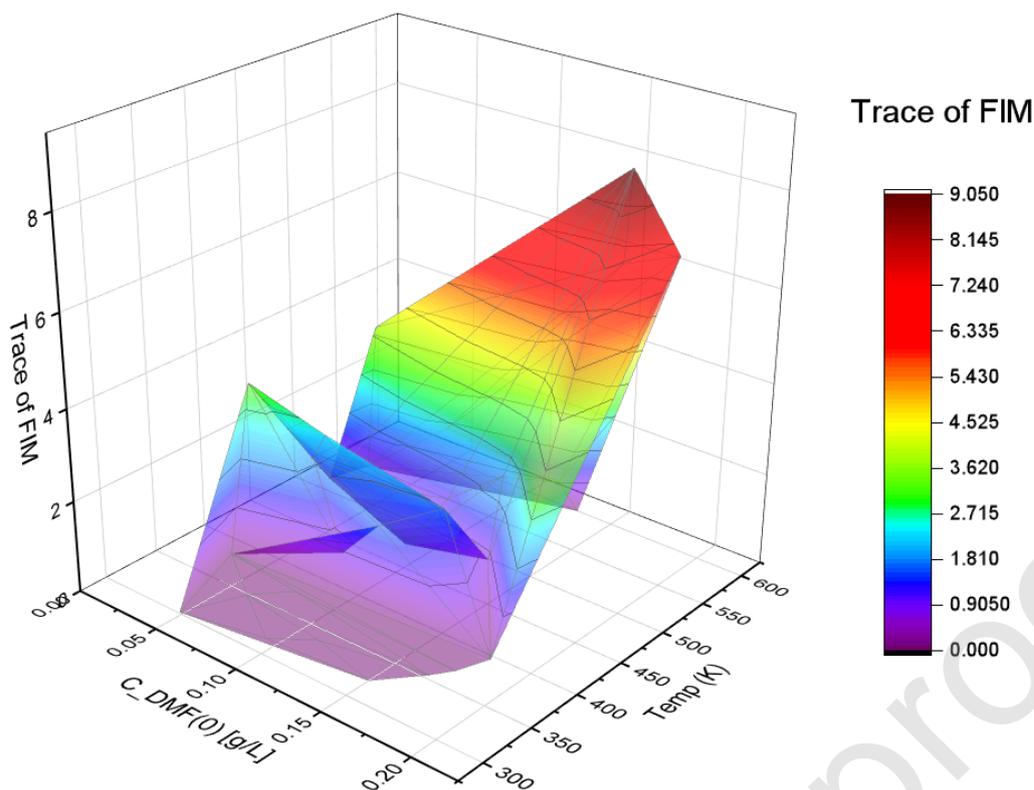


Figure 6 – Information map showing the regions of maximum information.

## 5. Conclusions and Outlook

A new model-based experimental design approach has been proposed in this paper to assess the parametric identifiability of kinetic models of HMF hydrogenation in batch reaction systems. The approach combines: *i*) a preliminary model-based analysis of information carried out by screening the experimental design space for the proposed kinetic models; *ii*) a model assessment using a-posteriori statistics obtained after parameter estimation from in-silico generated data; *iii*) a model-based experimental design, to identify the most promising regions of experimental conditions to explore in future experiments. Three model configurations (M1, M2, M3) based on the model proposed by Gyngazova et al. (2017) were investigated to determine the parameters for the HMF hydrogenation reaction where an explicit temperature dependency was introduced for the kinetic parameters. In M1, this dependency is represented by an Arrhenius expression and the full set of model parameters ( $A_i$  and  $E_{Ai}$ ) is estimated. In M2, a re-parametrisation in the form suggested by Buzzi-Ferraris and Manenti (2009) replaced the Arrhenius expression, with two new parameters ( $B_i$  and  $C_i$ ) being used to

represent the kinetic rate constants. M3 configuration was based on a reduced model approach where only the activation energies ( $E_{Ai}$ ) were estimated. These additional relationships made the model identification task more challenging. In particular, M1 was rejected for further analysis following preliminary assessments due to the high correlations between kinetic model parameters. M2 was rejected as this configuration does not accurately represent the system in a large design space, but only in a limited range of temperatures. However, the use of a model re-parametrisation succeeds on reducing the existing correlation between kinetic model parameters. M3 was chosen for further analysis as it was discovered that a limited correlation between parameters could be realised whilst maintaining a temperature dependency and a statistically reliable description of the underlying system in the investigated design space. A D-optimal experimental design was then carried out using M3, and the design showed that one, single experiment carried out at maximum initial reagent concentrations ( $C_{\text{HMF}}(0) = 0.200 \text{ g L}^{-1}$ ,  $C_{\text{DMF}}(0) = 0.198 \text{ g L}^{-1}$ ), a temperature of 455K and a duration of 594 minutes (approximately 10 hours), with an evenly allocation of 30 sampling points, is sufficiently informative to produce an accurate and precise estimation of the full set of kinetic parameters. Further research will target the application of the proposed optimal design in a laboratory reactor in order to precisely determine the full set of model parameters from experimental data and confirm the results of this analysis. Also it might be interesting for future studies to investigate the validity of M2 in a narrow design space. i.e. at temperatures that can be practically achievable in the lab. The experimental ranking and optimal designs carried out in this investigation suggest that this should be possible without severely reducing the amount of information obtained.

## 6. References

- Asprey, S.P., Naka, Y. (1999). Mathematical problems in fitting kinetic models—some new perspectives. *Journal of Chemical Engineering of Japan*, 32, pp 328–337.
- Buzzi-Ferraris, G. and Manenti, F. (2009) Kinetic models analysis. *Chem Eng Sci*, 64, pp 1061-1074.
- Bortoli, A. (2018) Optimal design of experiments for the identification of kinetic models of HMF hydrogenation. *Master's Thesis*, University of Padua, Italy.
- Bindwal, A. B. and Vaidya, P. D. (2013) Kinetics of Aqueous-Phase Hydrogenation of Levoglucosan over RuC Catalyst. *Industrial & Engineering Chemistry Research*, 52, pp 17781–17789.
- Bindwal, A. B. and Vaidya, P. D. (2014) Reaction Kinetics of Vanillin Hydrogenation in Aqueous Solutions Using a RuC Catalyst. *Energy & Fuels*, 28(5), pp 3357–3362.
- Chis, O., Villaverde, A., Banga, J. and Balsa-Canto, E. (2016) On the relationship between sloppiness and identifiability. *Math Biosc*, 282, pp 147-161.
- Deussen, P. (2019) Optimal design of experiments for the identification of kinetics models of HMF hydrogenation in batch reactor systems. *Master's Thesis*, University College London, England.
- Deussen, P. and Galvanin, F. (2021) On the practical identifiability of kinetic models of hydroxymethylfurfural hydrogenation in batch reaction systems. *Computer Aided Chemical Engineering*, 50, pp 859-865.
- Franceschini, G. and Macchietto, S. (2008) Model-based design of experiments for parameter precision: State of the art. *Chem Eng Sci*, 63, pp 4846-4872.
- Galvanin, F., Enhong, C., Al-Rifai, N., Gavriilidis, A. and Dua, V. (2016) A joint model-based design approach for the identification of kinetic models in continuous flow laboratory systems. *Computers and Chemical Engineering*, 95, pp 202-215.
- Galvanin, F., Macchietto, S. and Bezzo, F. (2006) Model-based design of parallel experiments. *Industrial & Engineering Chemistry Research*, Vol 46, pp 871-882.
- Gawade, A.B., Tiwari, M.S. and Yadav, G.D. (2016) Bio-based green process: Selective hydrogenation of 5-hydroxymethylfurfural to 2,5-dimethylfuran under mild conditions using Pd-Cs<sub>2.5</sub>h<sub>0.5</sub>pw<sub>12</sub>o<sub>40</sub>/K-10 Clay. *ACS Sustainable Chemistry and Engineering*, 4(8), pp 4113-4123.
- Process Systems Enterprise (2019). gPROMS Process Builder version 1.4.0, London, England.

- Grilic, M., Likozar, B. and Levec, J. (2014) Hydrodeoxygenation and hydrocracking of solvolysed lignocellulosic biomass by oxide, reduced and sulphide form of NiMo, Ni, Mo and Pd catalysts. *Applied Catalysis B: Environmental*, 150-151, pp. 275-287.
- Gyngazova, M.S., Negahdar, L., Blumenthal, L.C and Palkovits, R. (2017) Experimental and kinetic analysis of the liquid phase hydrodeoxygenation of 5-hydroxymethylfurfural to 2,5-dimethylfuran over carbon-supported nickel catalysts. *Chemical Engineering Science*, 173, pp 455-464.
- Hu, L., Tang, X., Xu, J., Wu, Z., Lin, L. and Liu, S. (2014) Selective transformation of 5-hydroxymethylfurfural into the liquid fuel 2,5-Dimethylfuran over carbon-supported ruthenium. *Industry and Chemical Engineering Research*, 53, pp 3056-3064.
- Jain, A.B. and Vaidya, P.D. (2016) Kinetics of catalytic hydrogenation of 5-hydroxymethylfurfural to 2,5-bishydroxymethylfuran in aqueous solution over RuC. *International Journal of Chemical Kinetics*, 48(6), pp 318–328.
- Khaled, N. (2020) Latin Hypercube, *MATLAB Central File Exchange*. Retrieved February 24, 2020. [<https://www.mathworks.com/matlabcentral/fileexchange/45793-latin-hypercube>]
- Luo, J., Arroyo-Ramirez, L., Wei, J. and Yun, H. (2015) Comparison of HMF hydrodeoxygenation over different metal catalysts in a continuous flow reactor. *Applied Catalysis A-General*, 508, pp 86–93.
- Margarita, G., Riccomagno, E., White, L.J. (2001) Structural identifiability analysis of some highly structured families of statespace models using differential algebra. *Journal of Mathematical Biology*, Vol 49(5), pp 433-454.
- Matlab version R2018b 9.5.0. MathWorks, Natick, MA, United States.
- Miao, H., Xia, X., Perelson, A.S. and Wu, H. (2011) On identifiability of nonlinear ODE models and applications in viral dynamics. *SIAM Rev Soc Ind & Appl Math*, Vol 53(1), pp 3-39.
- Montgomery, D.C. (2012) Chapter 4: Latin Squares, Randomised Blocks, and Related Designs. *Design and Analysis of Experiments*, 8th edition. Wiley.
- Pohjanpalo, H. (1978) System identifiability based on the power series expansion of the solution. *Mathematical Biosciences*, Vol. 41(1-2), pp 21-33.

Schwaab, M., Lemos, L., Pinto, J., 2008, Optimum reference temperature for reparametrization of the Arrhenius equation Part 2 Problems involving multiple reparametrizations. *Chemical Engineering Science*, 63.

van Putten, R., J. C. Waal, and E. de Jong, (2013) Hydroxymethylfurfural, a versatile platform chemical made from renewable resources. *Chemical Reviews*, 113.

Walter, E. and Pronzato, L. (1996) On the identifiability and distinguishability of nonlinear parametric models. *Mathematics and Computers in Simulation*, 42, pp 125-134.

Journal Pre-proof