



## Full Length Article

# Development of explainable AI-based predictive models for bubbling fluidised bed gasification process

Daya Shankar Pandey<sup>a,\*</sup>, Haider Raza<sup>b</sup>, Saugat Bhattacharyya<sup>c</sup>

<sup>a</sup> School of Physics, Engineering and Computer Science, University of Hertfordshire, Hatfield AL10 9AB, UK

<sup>b</sup> School of Computer Science and Electronic Engineering, University of Essex, Colchester CO4 3SQ, UK

<sup>c</sup> School of Computing and Intelligent Systems, Ulster University, Magee campus, Derry Londonderry, UK

## ARTICLE INFO

## Keywords:

Gasification  
Bubbling fluidised bed  
Machine learning  
Gradient boosting  
Decision tree regression

## ABSTRACT

In this study, seven different types of regression-based predictive modelling techniques are used to predict the product gas composition ( $H_2$ ,  $CO$ ,  $CO_2$ ,  $CH_4$ ) and gas yield (GY) during the gasification of biomass in a fluidised bed reactor. The performance of different regression-based models is compared with the gradient boosting model (GB) to show the relative merits and demerits of the technique. Additionally, SHapley Additive exPlanations (SHAP)-based explainable artificial intelligence (XAI) method was utilised to explain individual predictions. This study demonstrates that the prediction performance of the GB algorithm was the best among other regression-based models i.e. Linear Regression (LR), Multilayer perception (MLP), Ridge Regression (RR), Least-angle regression (LARS), Random Forest (RF) and Bagging (BAG). It was found that at *learning rate* (*lr*) 0.01 and *number of boosting stages* (*est*) 1000 yielded the best result with an average root mean squared error (RMSE) of 0.0597 for all outputs. The outcome of this study indicates that XAI-based methodology can be used as a viable alternative modelling paradigm in predicting the performance of a fluidised bed gasifier for an informed decision-making process.

## 1. Introduction

Optimising the production of clean, economically affordable energy and water with minimum impact on the environment is posing the most significant challenge of the 21<sup>st</sup> century [1]. A rapidly growing world population led to a substantial rise in energy demand. According to the International Energy Outlook report, globally the marketable energy consumption will increase by 48% between 2012 and 2040 [2]. As an alternative to depleting fossil fuel reserves, researchers have been attracted to explore clean and renewable energy resources. In this regard, energy recovery from biomass and waste is important for both developing and developed countries. Biomass has become an attractive source of energy because it not only reduces carbon and greenhouse gas (GHG) emissions to the environment but also improves the air quality [3].

Biomass originates from a range of different sources in a wide variety of forms from untreated biomass (straw, peel kernels, husk *etc.*), treated biomass (wood pellets, olive pruning, softwood pellets *etc.*) to cultivated energy crops, residues and waste-derived fuels (animal manures, dry sewage sludge, municipal solid waste *etc.*). Lignocellulosic biomass is

now considered as a realistic renewable energy source with unique advantages, for instance, low sulphur content, high volatile matter and oxygen content, and is ideally emitting a net-zero carbon into the atmosphere [4].

Thermochemical processes such as combustion, gasification and pyrolysis are widely used for energy recovery from dry lignocellulosic biomass [5]. To convert biomass into bioenergy, gasification is emerged as an alternative to traditional combustion technology, offering distinct advantages such as higher energy efficiency, lower emissions of particulate matter,  $NO_x$  and  $SO_x$  while being compliant with emission standards [6]. In addition, the gasification process can generate electricity from the waste with an efficiency of 34% compared to the incineration process, which has thermal efficiency of around 20% [5]. Gasification provides flexibility to produce heat and power based on a clean biomass-derived product gas or synthesis gas. The product gas is a combination of different gases that evolved during the gasification process, which involves several chemical reactions. The product gases, after being cleaned and filtered, can be used in a variety of applications ranging from internal combustion-based electric generators to combustion for thermal applications or synthesis of chemicals. Furthermore, the by-

\* Corresponding author at: School of Physics, Engineering and Computer Science, University of Hertfordshire, Hatfield AL10 9AB, UK.

E-mail address: [d.pandey@herts.ac.uk](mailto:d.pandey@herts.ac.uk) (D.S. Pandey).

product from the gasifier can be used in the construction industry or in power plants as a supplementary fuel or as a component material for fertiliser [7–9].

Gasification has been identified as a highly efficient thermochemical conversion process but due to its complex nature, this technology is yet to be deployed at an industrial scale. To understand the complexity and optimal operation condition, it is essential to conduct time-consuming and costly large-scale gasification experiments. Therefore, mathematical modelling techniques can be a viable option, which could save both time and money [10]. The advanced machine learning algorithm-based models require fewer system-level details compared to equilibrium, fluid dynamic- or kinetic-based models. ANN-based machine learning techniques have extensively been utilised to predict the product gas compositions, gas yield and the calorific value in gasification processes. Guo *et al.* developed ANN-based models to predict the product gas yield and gas composition in an atmospheric steam-blown biomass fluidised bed gasifier and proposed that the multilayer feed-forward neural network outperformed traditional regression models [11]. The ANN-based model was used in the feasibility study of the municipal solid waste gasification process [12]. Chavan *et al.* proposed two different types of ANN-based data-driven models in order to predict the gas heating value and the gas production rate from a fluidised bed coal gasifier operating in a steady state condition [13]. Similarly, an ANN-based algorithm was used in bubbling and circulating fluidised bed gasifiers to estimate the major product gas compositions (CO, CO<sub>2</sub>, H<sub>2</sub> and CH<sub>4</sub>) and the total gas yield [14]. Pandey *et al.* predicted the lower heating value of the gas and products and gas yield originating from a fluidised bed municipal solid waste gasifier and proposed a robust approach to optimised the network [15]. ANN-based models have also been used in a downdraft fixed bed gasifier to predict gas composition [16]. Machine learning-based regression models have been developed for predicting gas composition and higher heating value. The authors have claimed that the multilayer perceptron and decision tree regression models had comparatively improved outcome compared to polynomial regression and support vector regression [17]. Shahbaz *et al.* illustrated an ANN-model to study the performance of fixed bed steam gasification of palm oil waste using steam to biomass ratio (SBR), temperature, Coal bottom ash and CaO/biomass [18]. Recently, an attempt was made to predict the effect of bed materials in bubbling fluidised bed gasification using an ANN-based algorithm [19]. Serrano *et al.* have also presented an ANN-based algorithm for estimation tar from a bubbling fluidised bed gasifier [20]. Yan *et al.* applied ANN-based machine learning algorithms to estimate the performance of heterogeneous, multi-component materials as oxygen carriers for chemical-looping processes [21].

Despite the distinctive advantages and the proven prediction capability of advanced machine learning techniques, so far mostly ANN-based models have been used. In the recent past, attempts have been made to derive analytical regression models using genetic programming algorithm and exploiting fuzzy optimisation for a sustainable palm oil-based integrated biorefinery and Bayesian method to quantify the uncertainty in data-driven modelling approached pertaining to the gasification process [10,22–24]. The gradient boosting method was applied to predict the gas yields from the supercritical water gasification process and features were interpreted by the SHapley Additive exPlanations (SHAP) value [25]. However, there have been very few cases of developing explainable artificial intelligence (XAI) models with subsequent validation using experimental results and identifying the key influencing parameters of the gasifier performance reported. The present study is focused on exploiting a range of machine learning models such as linear regression, multilayer perceptron, ridge regression, least angle regressions, random forest, bagging and gradient boosting methods to predict the product gas composition and the total product gas yield. Additionally, SHAP-based XAI method has been utilised to explain individual predictions as well as identify the key influencing parameter. To the best of the author's knowledge, this work is the first of its kind where a gradient-boosting tree-based model has been used under

multiple-input and single-output (MISO) settings in a gasification system with XAI. The key objective of this research is to develop and propose data-driven mathematical models, which can be used to simulate the gasification process with model explainability and improved accuracy. In addition, the proposed model can be useful in identifying, which input parameter has the maximum influence on product gas composition and gas yield, and as a result, it could help to optimise the performance of the gasifier.

## 2. Material and methods

### 2.1. Material

Fig. 1 shows the input and output variables used in this study. Nine input variables: carbon (C), hydrogen (H), oxygen (O), moisture (MC) and ash (Ash) contents for the biomass, and the process parameters such as equivalence ratio (ER), reactor temperature (T), bed material (BM), steam to biomass mass ratio (SBR). These variables are represented by the input vector  $x = [C, H, O, MC, Ash, ER, T, BM, SBR]$ . The input data for the models were chosen based on the biomass properties (C, H, O, MC and ash content) and the process conditions of the gasifier (T, ER and SBR). ER is defined as the ratio between the actual air fed to the gasifier and the necessary amount of air required for stoichiometric combustion of the biomass [10]. The output variables hydrogen (H<sub>2</sub>), carbon monoxide (CO), carbon dioxide (CO<sub>2</sub>), methane (CH<sub>4</sub>) and gas yield (GY) are represented by the output vector  $y = [H_2, CO, CO_2, CH_4, GY]$ . A total of 222 experimental observations from fluidised bed gasifiers were used in this study.

The dataset used in this study were collected from literature and were all from fluidised bed gasifiers operating at steady state condition. Therefore, the application of the proposed models are limited to fluidised bed reactor systems. The experiments were performed at different temperature ( $650 \leq T \leq 1050$  °C) and equivalence ratio ( $0.15 \leq ER \leq 0.49$ ). The input and output variables are in different units, and the mean and standard deviation of values are not in the same order in a numeric sense. Since input and output parameters are not in the same numeric range and can have numerical issues while fitting such kind of dataset in data-driven modelling approaches. To encounter the numerical issues, the dataset was normalised and then the normalised dataset is used for training and testing purposes. The dataset was divided into training and testing sets: 70% data were used to train the model and 30%

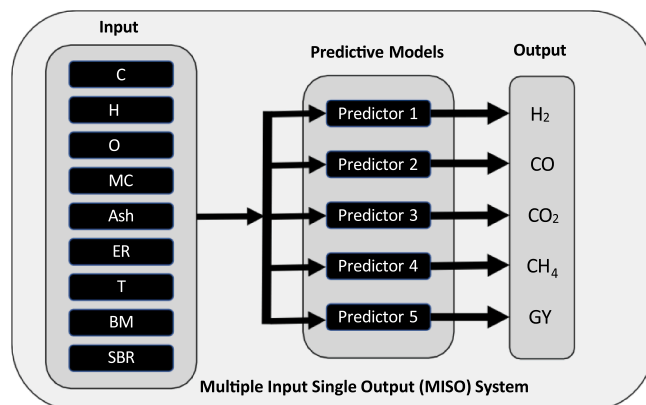


Fig. 1. Multiple Input and Single Output (MISO) systems and Data descriptor: input variables  $x = [C, H, O, MC, Ash, ER, T, BM, SBR]$  and output variables  $y = [H_2, CO, CO_2, CH_4, GY]$ , where carbon (C), hydrogen (H), oxygen (O), moisture (MC) and ash (Ash) contents for the biomass, and the process parameters such as equivalence ratio (ER), reactor temperature (T), bed material (BM), steam to biomass mass ratio (SBR), hydrogen (H<sub>2</sub>), carbon monoxide (CO), carbon dioxide (CO<sub>2</sub>), methane (CH<sub>4</sub>), and gas yield (GY). MISO uses all nine input variables to estimate the composition of each output composition by training five different prediction models (i.e. Predictor 1 to Predictor 5).

was used for testing (to validate the performance of the algorithm for unknown entries). Azadi and Karimi-Jashni [26] provided a detailed explanation of the advantages and disadvantages of different accuracy measures such as mean absolute error (MAE), root mean squared error (RMSE), and normalised root mean squared error (NMSE). Hence, to evaluate the performance of the trained predictive models on the test data, we have selected the RMSE and the coefficient of determination ( $R^2$ ) value to demonstrate the goodness of fit for training and testing data which is a summary of fitness performance in supervised learning.

The experimental data from fluidised bed gasifiers was compiled by [19] from works of literature [4,27–49]. Table 1 reports the means and variances for all the different input and output variables. This information is reported to show the nature of the dataset and to facilitate reproducibility. Details of the datasets are provided in Supplementary information. Additionally, the Python implementation of the proposed work is available via open source in GitHub: <https://github.com/sagihaider/XAI-Gasification-Modelling.git>.

## 2.2. Methods

The objective of this study as mentioned above is to predict the product gas composition ( $H_2$ ,  $CO$ ,  $CO_2$ ,  $CH_4$ ) and the total product gas yield (GY) from the input variables from input vector  $x$ . The solution is to predict the numerical output compositions from numerical input variables. So, for this purpose, we have designed a specialised machine learning model that supports predicting the different outputs from the same set of variables. The model (as shown in Fig. 1) uses all nine input variables to estimate the composition of each output composition by training five different prediction models (Predictor 1 to Predictor 5). In order to gain better insight and interpretability on the influence of the input variables on each trained model, this study also identified the most relevant input variables for each output composition yielded from its corresponding trained model. Additionally, it is well known that AI-based prediction algorithms are not expected to be perfect due to various factors such as dataset shift, data bias, etc. [50]. Their unexplained failures can be adverse.

To overcome this challenge, these days researchers are focusing on developing XAI methods, which have emerged as an effort to gain trust using AI-based predictive algorithms. SHAP interpretability method proposed by Lundberg and Lee [51] is used in this study (a brief description of SHAP is presented in Section 3.1). This approach is distinctive and will help to understand the correlation of the relevant input variables to their corresponding output which could lead to further improvements in future models.

## 2.3. Overview of the modelling paradigm

Consider a challenge of function estimation in the classical super-

**Table 1**  
Means and variances for all the different input and output variables.

	Variable	Mean	Variance
Input	C [wt.% db]	47.97	12.43
	H [wt.% db]	6.16	2.36
	O [wt.% db]	34.06	13.16
	MC [wt.% ar]	8.88	5.79
	Ash [wt.% db]	10.53	13.62
	ER [-]	0.27	0.07
	T [°C]	790.30	60.03
	BM	–	–
	SBR [-]	0.21	0.60
	$H_2$ [vol.% $N_2$ free]	22.57	8.98
Output	$CO$ [vol.% $N_2$ free] $CO_2$ [vol.% $N_2$ free]	30.0833.44	7.7110.09
	$CH_4$ [vol.% $N_2$ free]	9.12	3.97
	GY [ $Nm^3/kg$ daf]	2.32	0.97

For bed material: 1: Silica sand; 2: Ofite; 3: Olivine; 4: Alumina, where db, dry basis; ar, as received; daf, dry ash free.

vised learning setting. Let's assume a supervised dataset  $(X, Y)_{i=1}^N$ , where  $X = (x_1, \dots, x_N)$  refers to the explanatory input variables set and  $Y$  is a set of the response variable. The goal is to reconstruct the unknown functional dependence with the estimate  $f(X)$ , such that some specified loss function  $\Psi(Y, f)$  is minimised:

$$\hat{f}(X) = Y$$

$$\hat{f}(X) = \underset{f(X)}{\operatorname{argmin}} \Psi(Y, f(X))$$

The response variable  $Y$  can come from different distributions. This naturally leads to the specification of different loss functions  $\Psi$ . In particular, if the response variable is binary, (i.e.  $Y \in \{0, 1\}$ ), one can consider the binomial loss function. If the response variable is continuous, (i.e.  $Y \in R$ ), one can use any loss function based on the need such as  $L_2$  squared loss function,  $L_1$  absolute loss function, or the robust regression Huber loss.

### 2.3.1. Machine learning methods

We have applied seven perspective-supervised machine learning models (Multivariate Linear Regression (LR), Multilayer perception (MLP), Ridge Regression (RR), Least-angle regression (LARS), Random Forest (RF), Bagging (BAG) and Gradient Boosting (GB)) to predict the product composition and the product gas yield from fluidised bed gasifiers. Details of each supervised machine learning algorithm are given in each section below, respectively.

General Multivariate Linear Regression (LR): The general linear model or general multivariate regression model is a packed way of simultaneously using several multiple linear regression models. Multiple linear regression models can be written as equation (1):

$$\|Y - XB\|_2^2 \quad (1)$$

where  $Y$  is a matrix of multivariate measurements (i.e. dependent variables), and  $X$  is a matrix of observations (i.e. independent variables).  $B$  is a matrix containing parameters that are usually to be estimated by minimising the residual sum of squares between the observed dependent variables in the dataset, and the dependent variables predicted by the linear approximation.

Multilayer Perceptron (MLP): A MLP [52] is a type of artificial neural network (ANN) that consists of multiple layers of interconnected nodes, each with learnable weights and biases. The first layer receives input data, which is then transformed through a series of non-linear transformations in the hidden layers before producing an output in the final layer. MLPs are commonly used for classification and regression tasks and can learn complex non-linear relationships between inputs and outputs. The weights and biases are typically updated using back-propagation, which involves computing gradients and adjusting the parameters to minimise the loss function.

Ridge Regression (RR): RR is a linear least square with  $l_2$  regularisation, which aims at avoiding over-fitting. It minimises the objective function as follows:

$$\|Y - XB\|_2^2 + \lambda \|B\|_2^2 \quad (2)$$

The regularisation term in eq.2 is known as  $l_2$ -norm or Ridge regression penalty or squared penalty.

Least-angle regression (LARS): LARS regression is also a linear least square with  $l_1$  regularisation, which aims at avoiding over-fitting. LARS regression provides an alternate way to train a Lasso regularised linear regression model that adds a penalty to the loss function during training. LARS uses the Akaike information criterion (AIC) and Bayes Information criterion (BIC) to select the value of the regularisation parameter by making a trade-off between the goodness of fit and the complexity of the model. It minimises the objective function as follows:

$$\frac{1}{2 * n\_samples} \|Y - XB\|_2^2 + \lambda \|B\| \tag{3}$$

The regularisation term in eq.3 is known as l1-norm or Lasso regression penalty or absolute penalty, where  $n\_samples$  is the number of samples in the dataset.

Random Forest (RF): RF is a tree-based supervised learning algorithm, which fits a number of decision trees on different subsets of the training dataset to enhance the predictive performance. It also helps in preventing over-fitting. In random forest, tuning of two important hyper-parameters is required (i.e. number of estimators (trees) and the maximum depth of the tree).

Bagging (BAG): BAG is almost the same as the RF. The only difference is that bagging fits a number of decision trees on random subsets of the training dataset to enhance the predictive performance. It also helps in preventing overfitting by reducing the variance. In bagging also tuning of two important hyper parameters is required (i.e. number of estimators (trees) and the maximum depth of the tree).

Gradient Boosting (GB): GB machines are a family of powerful tree-based machine-learning techniques that have shown considerable success in a wide range of practical applications. In the context of the paper, an overview of the GB modelling paradigm is provided as follows:

GB algorithm builds a collection of trees called (an ensemble) of shallow trees, where each tree in the ensemble improves on the previous one. These shallow trees are generally called weak predictive models, which can act as powerful predictive models when appropriately tuned. In GB, boosting supports adding new models to the ensemble in a sequential manner. Additionally, boosting help in handling the bias-variance trade-offs by building a weak learner and then sequentially boosting the performance by adding new trees, where a newly added tree manages the error made by the previously added tree. In other words, the newly added tree considers specific rows from the training dataset, which made the largest prediction error in the previous tree. A schematic diagram illustrating the concept of GB is presented in Fig. 2.

The key elements such as weak learners and their training with respect to errors are also discussed. A weak learner model is one whose prediction performance is slightly better than random guessing, which can be 50%-50% in terms of binary classification settings. Boosting helps to improve the performance in each model by using their experience from the previous model, where it will specifically focus on rows of the training data with the largest error in the previous tree. Focusing on the error from the previous tree, the new tree in the sequential model which enhances its performance as outlined by the pseudo gradient boosting code. The python package *scikit-learn* was used for this purpose.

**Algorithm 1** Gradient Boosting Algorithm

**Input:** X, Y  
**Output:** Y  
**Step 1:** Fit a decision tree to the data:  $f_1(X, Y)$

(continued on next column)

(continued)

**Algorithm 1** Gradient Boosting Algorithm

**Step 2:** Then fit the next decision tree to the error of the previous tree:  $h_1(X) = (Y - f_1(X))$   
**Step 3:** Add new tree to model:  $f_2(X + h_1, Y)$   
**Step 4:** Then fit the next decision tree to the error of the previous tree  $f_2: h_2(X) = (Y - f_2(X))$   
**Step 5:** Add new tree to model:  $f_3(X + h_2, Y)$   
**Step 6:** Repeat this process until the condition is met (i.e. based on ross- validation)

2.3.2. Computational resources

Simulations were performed on a Google Colab Pro cloud-based Python environment, which consists of Intel(R) Xeon(R) CPU @ 2.20 GHz processors with 32 GB of RAM, CPU Core = 2. The computational time cost of training the predictive models are as follows: General Multivariate Linear Regression (LR): 0.00078 s; MLP:  $4 \times 10^{-6}$  s; Ridge Regression (RR): 0.000723 s; Least angle regression (LARS): 0.0175 s; Random Forest (RF): 0.864 s; Bagging (BAG): 1.009 s; and Gradient Boosting (GB): 1.182 s.

3. Results

3.1. Interpretability using SHAP

SHAP computes the contribution of each feature of an instance  $\times$  to explain its prediction. This form of explanation uses the Shapley values from coalition game theory, which provides information on how to fairly distribute the contribution of each prediction among the features. This paper uses a variant of SHAP for a tree-based machine learning model called TreeSHAP [53]. The predictions from the GB model are interpreted both globally and locally. Global interpretability provides insight into the contribution of each feature to the estimated output, which can either be positive or negative. Additionally, each instance or observation is allotted its own SHAP value, which can greatly increase the transparency of a model in terms of explainability. Using the SHAP method, one can easily explain the reasoning behind a single prediction, which further provides details about the impact of the features. This form of interpretability is local in nature and can easily be illustrated with the help of a summary plot. The summary plot in Fig. 3 uses a toy dataset with 8 representative variables ( $X_1$ - $X_8$ ), which combines feature importance with feature effects. Each point on the summary plot is a Shapley value for a feature and an instance. The feature order on the Y-axis is determined by feature importance in decreasing order and the X-axis presents the Shapley value. The colour represents the value of the feature from low to high. Overlapping points are jittered in the Y-axis direction, so we get a sense of the distribution of the Shapley values per feature. The SHAP value plot can further show the positive and negative relationships of the predictors with the target variable.

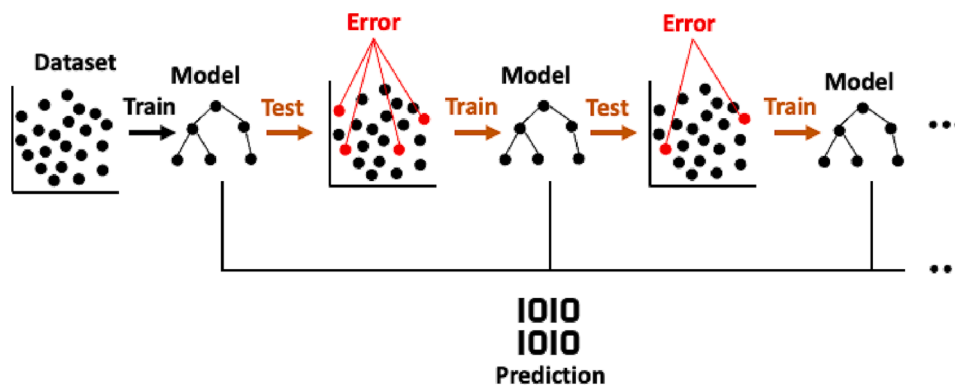


Fig. 2. Gradient boosting-based sequential ensemble approach.

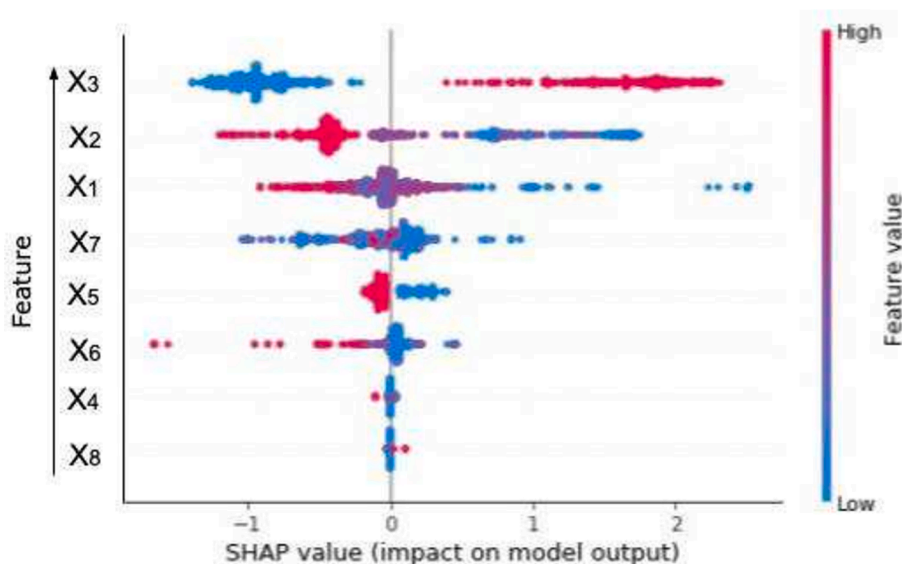


Fig. 3. SHAP summary plot on a toy dataset with 8 representative variables ( $X_1$ - $X_8$ ), where X-axis represents the Shapley value and Y-axis shows the features.

### 3.2. Regression analysis using GB

General multivariate, ridge, least-angle, random forest, bagging and GB regression-based predictive models were exploited, and a thorough grid-search method was used to find the optimal *learning rate* ( $lr$ ) and *number of boosting stages* ( $est$ ) in the training stage. It was found that the GB regression model at  $lr = 0.01$  and  $est = 1000$  yielded the best result in terms of the average root mean square error (RMSE) for all outputs. The following optimal parameters (after a thorough grid search) were employed for each of the comparative algorithms:

- **RR:** Regularisation Strength = 1.0; Tolerance = 1000.
- **MLP:** activation function = Rectified Linear Unit; regularisation parameter = 0.0001; learning rate = 0.01; number of hidden neurons = 2; optimiser = adam; loss function = mean squared logarithmic error loss
- **RF:** Number of Trees = 500; criterion = mean squared error; Bootstrap = True.
- **BAG:** Number of estimators = 500; Bootstrap = True.

It is evident from the literature review that the predictions of ANN-based are generally suitable for complex datasets. Furthermore, classical regression approaches are exploited to fit polynomial models without significant nonlinearity. We have applied six different regression methods i.e. Linear regression (LR), Ridge Regression (RR), Multi-layer Perceptron (MLP), Least-angle regression (LARS), Random Forest regressor with decision Tree as a weak learner (RF) and a Bagging regressor with Decision Tree as a weak learner (BAG) as tabulated in Table 2. It compares the average RMSE across the six outputs of the GB regressor with other standard and ensemble regression methods. The

comparison includes LR, RR, MLP, LARS, RF and BAG. Amongst different combinations for the input variables, the best solution by regression analysis was found using the BAG model containing a mean RMSE of 0.07442 compared to GB which has an RMSE of 0.0604 which is less than the other regression models. Fig. 4 shows the test data prediction of each output composition (in red) using the MISO model and a comparison with its true experimental values (in blue).

It can be seen from Table 2 that the selected GB regressor model is superior to the rest of the regressors. In addition, it is worth mentioning that the reason for choosing the GB simulation paradigm in this study is not only limited to its better performance but also includes its ability to avoid over-fitting, robustness to small changes in the dataset, encouragement of diversity and interpretability of the results.

The reported RMSE of product gas compositions were 0.0941 for  $H_2$ , 0.0584 for CO, 0.0839 for  $CO_2$ , and 0.0524 for  $CH_4$  whereas, the RMSE for gas yield was 0.0046. It is evident from Table 2 that the predictability of GB based regression model outperformed other models. Although, the  $R^2$  value for CO is slightly lower but is within a whisker limit with RF and BAG methods. Furthermore, the RMSE values positively indicate the accuracy of the GB regression model to the true values of the output composition and vindicate the efficacy and robustness of the proposed predictive approach. The optimal model is selected based on the minimum reported RMSE value in Table 2. The discussion section is focused on identifying the key influencing input parameters on product gas composition and gas yield as well their importance in the gasification process.

### 4. Discussion

Various machine learning mostly ANN-based approaches have been

Table 2  
RMSE and  $R^2$  comparison of GB with other Regression Methods on the test dataset.

Model	$H_2$		CO		$CO_2$		$CH_4$		GY		Mean	
	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$	RMSE	$R^2$
LR	0.1199	0.410	0.1465	0.323	0.1291	0.389	0.1249	0.137	0.0659	0.764	0.1173	0.4046
MLP	0.1289	0.317	0.1498	0.292	0.1130	0.531	0.1344	0.001	0.0682	0.747	0.1188	0.3779
RR	0.1155	0.452	0.1470	0.318	0.1121	0.540	0.1225	0.170	0.0639	0.778	0.1122	0.4516
LARS	0.1164	0.444	0.1465	0.323	0.1285	0.395	0.1257	0.127	0.0646	0.774	0.1163	0.4126
RF	0.0727	0.783	0.0891	0.749	0.0726	0.807	0.0735	0.701	0.0635	0.782	0.0743	0.7644
BAG	0.0734	0.779	0.0893	0.748	0.0721	0.810	0.0733	0.703	0.0628	0.786	0.0742	0.7652
GB	0.0568	0.868	0.0901	0.744	0.0687	0.827	0.0429	0.898	0.0437	0.897	0.0604	0.8468

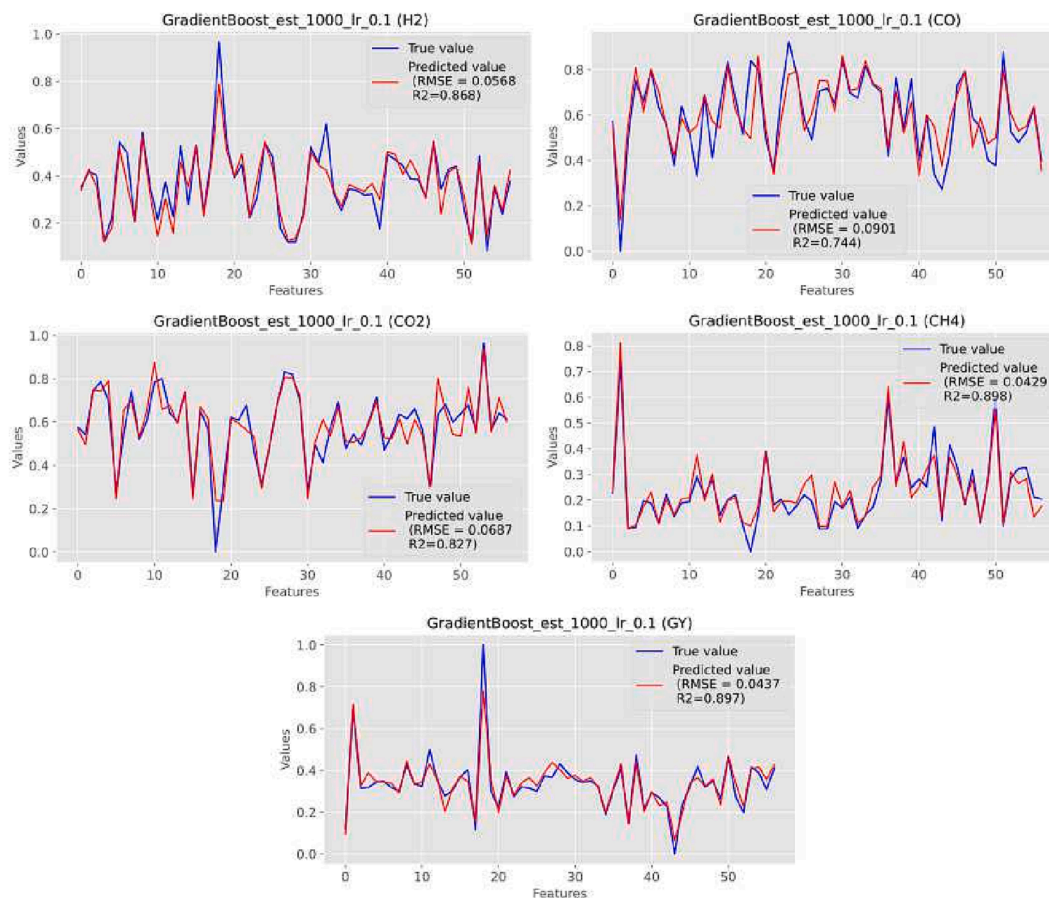


Fig. 4. Gradient boosting-based comparison of individual predicted values with true values of H<sub>2</sub>, CO, CO<sub>2</sub>, CH<sub>4</sub>, and GY and their RMSE values on the test dataset.

exploited as a predictive tool for estimating product gas composition, gas yield, calorific value etc. Moreover, the XAI-based approach is not yet used. In addition, identifying the key influential input parameter and its importance is rarely reported in the literature.

In previous studies, an equation-based approach was used for calculating the relative impact of input variables on the predicted output [54]. In contrast to the equation-based approach, this work used the Shapley value (refer to Section 3.1) to explain the importance and influence of the capacity of each input variable on each output parameter. We have compared the results of the GB model and pointed out the most important and highest impacting features as input. These features were selected based on SHAP value, a feature selection methodology. In the following section, we have explained how input variables are influencing the output and the mechanism behind it.

A summary plot of the global contribution of input variables is presented in Fig. 5. Fig. 5 presents the influencing input parameter on a specific output. It is worth noting that the most influencing input parameter may be negatively affecting the performance of the gasifier. The hydrogen content of the biomass has been identified as the most influencing parameter for hydrogen content in the product gas. It is clearly indicating that the hydrogen composition in the product gas is directly correlated to the hydrogen content in the biomass. Apart from hydrogen content, the temperature of the gasifier and steam to biomass ratio positively influence the hydrogen content in the product gas. This can be linked with the improved water gas shift reaction which favoured in the forward direction at an elevated temperature. It can be evident from Fig. 5 that the ER is the second most influencing parameter however, it is negatively impacting the hydrogen content in the product gas. Increasing ER favours oxidation reactions that would also increase the gasifier temperature and as a result, the equilibrium reactions shift in a forward direction leading to higher CO<sub>2</sub> and H<sub>2</sub>O concentration at the

expense of CO and H<sub>2</sub> [55]. This is also evident by looking at the influence of oxygen content. The effect of the remaining input parameters is negligible as they are jittered on the corresponding X-axis.

In Fig. 5, we can see that the most influencing parameter for CO output is directly linked with the indigenous amount oxygen content of the biomass. Temperature plays a positive impact on CO due to improved Boudouard reaction favouring CO production. Nevertheless, ER can have a detrimental effect on CO because it favours the formation of more stable compound CO<sub>2</sub> resulting from partial oxidation of CO. Steam to biomass ratio (SBR) also affects the final composition of the product gas triggering a water–gas shift reaction. Ash content and bed materials show a positive influence on CO. The effect of moisture content is almost neutral. In conclusion, the unique dynamics of CO output are mainly driven by the carbon and oxygen content of the biomass and the process temperature.

A graphical representation of Shapley value on CO<sub>2</sub> clearly indicates that ER greatly influences the concentration of the product gas (promoting oxidation reaction). It indicates that the maximum CO<sub>2</sub> concentration in the product gas can be linked with the higher ER. Temperature and SBR have a positive influence on the concentration of CO<sub>2</sub> in the product gas linked with equilibrium reaction shift and the water-gas shift reaction. The impact of the oxygen content is almost negligible.

The CH<sub>4</sub> content in the final product gas basically represents the amount of CH<sub>4</sub> formed during the devolatilisation process [56]. However, from Fig. 5 it is evident that the temperature of the reactor is the key influencing parameter for CH<sub>4</sub> followed by the carbon content. The steam to biomass ratio also improves methane formation moreover, the effect of steam to biomass ratio compared to temperature is almost negligible. The CH<sub>4</sub> concentration remains fairly constant since the steam methane reforming reaction does not dominate/play an

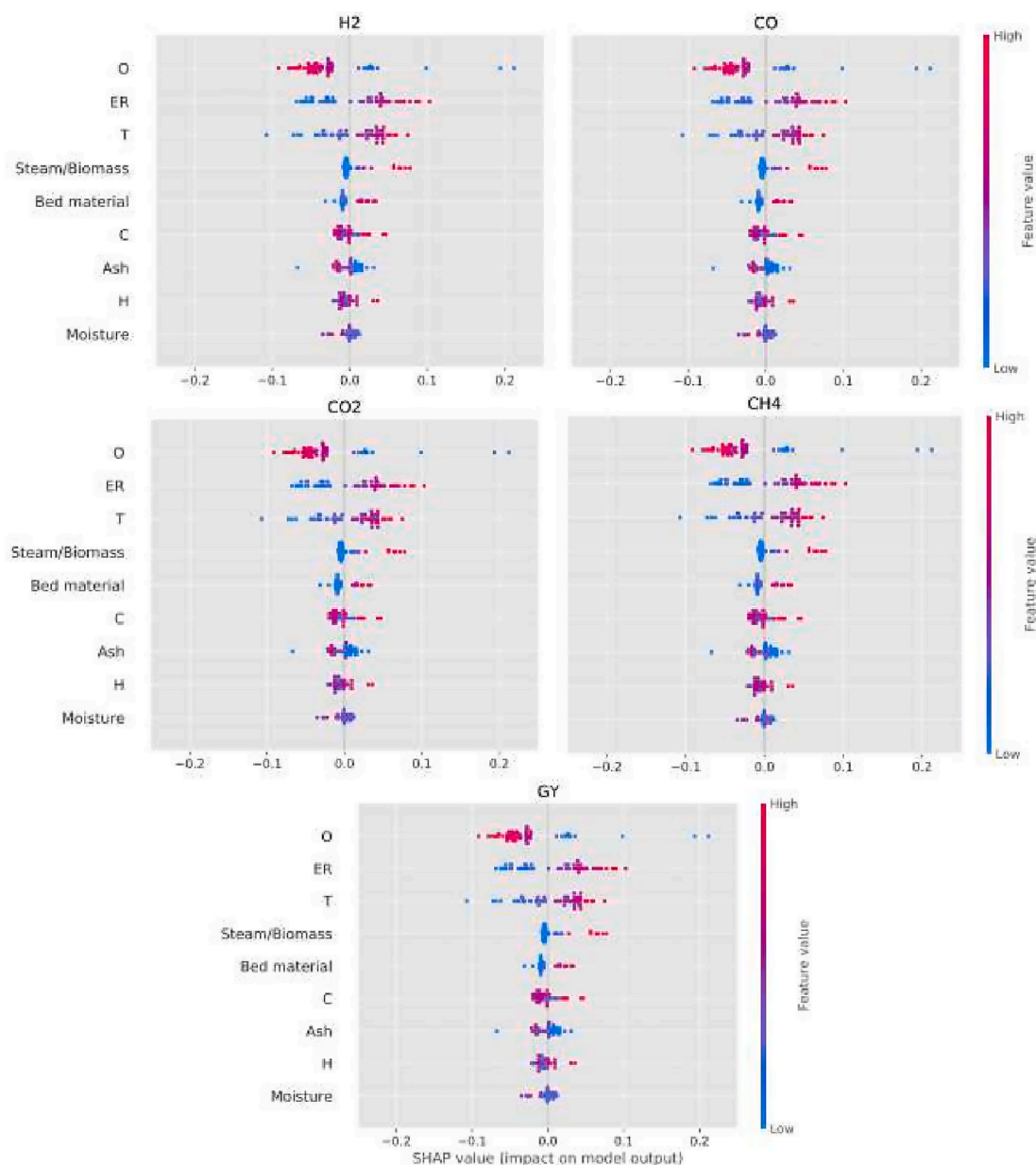


Fig. 5. Summary plot to interpret the global contribution of features on prediction.

influential role below 1000 °C and the CH<sub>4</sub> formation is kinetically controlled [57]. Despite the fact that higher temperature inside the reactor favoured the endothermic reactions, it can be recommended that the temperature of the gasifier should be carefully controlled because it could lead to auto-thermal operational conditions (with an increase in ER) which will increase the CO<sub>2</sub>/H<sub>2</sub>O ratio compared to the desired CO/H<sub>2</sub> ratio.

In addition to comparing the overall performances of proposed machine learning methods, one must investigate the prediction performance for each output variable to explore the viability of machine learning methods for the biomass gasification process more deeply. Lastly, rather than treating the CO and CO<sub>2</sub> outputs as a part of the biomass gasification process, one can create models specifically designed process just for predicting those outputs parameters individually. By using a different set of features and hyper-parameter configurations, it would be possible to obtain greater prediction performance for CO and CO<sub>2</sub>.

In a nutshell, we have shown that GB or other machine learning

paradigms can be helpful in investigating the impact of input variables and can play an important role in optimally designing large-scale gasification experiments with desired product gas yield and composition. Similarly, the process parameters (T, ER and SBR) can be optimised accordingly. Overall, advanced data-driven modelling approaches such as GB has the potential to provide key information for the optimisation of the gasification process and can save time and money. The present model predicted well the impact of the key input and process parameters on the output over the range of tested operating conditions. Nevertheless, it can further be improved if the model can be trained with more data with a wider range of operating conditions.

## 5. Conclusion

In this work, a GB-based data-driven modelling approach was exploited to predict the impact of input and process parameters on product gas composition and the total product gas yield in a fluidised bed gasifier. Compared to classical regressions (LR, RR), Least-angle

regression (LARS), Multilayer perception (MLP), Random Forest (RF) and Bagging (BAG), the Gradient Boost (GB) model shows the best predictive capacity with the lowest RMSE (0.0604) and highest  $R^2$  (0.8468). Furthermore, we discussed and illustrated the impact of influencing input variables using a graphical representation called SHAP value.

The proposed GB-based regression model not only predicted the outcomes of product gas composition and the total gas yield from fluidised bed gasifiers accurately but also provided greater insight into the relative influence of the input variables (physical) on the final product. The developed model predicted product gas composition and the gas yield well and can be employed in learning and prediction of nonlinear complex mapping of gasification yields. Since, these data-driven simulation paradigms offer distinctive advantages over the first principle (mass and energy balance) models such as it does not require domain-specific knowledge which makes it easy to exploit to simulate complex thermochemical processes. Therefore, it is expected that data-driven and machine-learning approaches will complement the design of experiments, complex process optimisation, characterisation and sensitivity analysis of input parameters and their impacts on the final product.

### CRedit authorship contribution statement

**Daya Shankar Pandey:** Conceptualization, Data curation, Investigation, Methodology, Visualization, Formal analysis, Writing – original draft, Writing – review & editing. **Haider Raza:** Data curation, Formal analysis, Methodology, Visualization, Writing – original draft, Writing – review & editing. **Saugat Bhattacharyya:** Data curation, Formal analysis, Methodology, Visualization, Writing – original draft, Writing – review & editing.

### Declaration of Competing Interest

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

### Data availability

Data will be made available on request.

### Acknowledgments

H.R. was supported by the Economic and Social Research Council (ESRC) funded Business and Local Government Data Research Centre under Grant ES/S007156/1.

### Appendix A. Supplementary data

Supplementary data to this article can be found online at <https://doi.org/10.1016/j.fuel.2023.128971>.

### References

- Omenn GS. Grand challenges and great opportunities in science, technology, and public policy. *Science* 2006;314(5806):1696–704.
- J. Conti, P. Holtberg, J. Diefenderfer, A. LaRose, J. T. Turnure, L. Westfall, International energy outlook 2016 with projections to 2040, Tech. rep., US-DOE Energy Information Administration (EIA), Washington, DC (United States) (2016).
- Tester JW, Drake EM, Driscoll MJ, Golay MW, Peters WA. Sustainable Energy: Choosing Among Options. MIT press; 2012.
- Katsaros G, Pandey DS, Horvat A, Almansa GA, Fryda LE, Leahy JJ, et al. Experimental investigation of poultry litter gasification and co-gasification with beech wood in a bubbling fluidised bed reactor—effect of equivalence ratio on process performance and tar evolution. *Fuel* 2020;262:116660.
- Murphy JD, McKeogh E. Technical, economic and environmental analysis of energy production from municipal solid waste. *Renew Energy* 2004;29(7):1043–57.
- Arena U. Process and technological aspects of municipal solid waste gasification. A review. *Waste Manage* 2012;32(4):625–39.
- Serrano D, Kwapinska M, Sanchez-Delgado S, Leahy JJ. Fly ash characterization from cynara cardunculus l. gasification. *Energy Fuels* 2018;32(5):5901–9.
- Pandey DS, Yazhenskikh E, Müller M, Ziegner M, Trubetskaya A, Leahy JJ, Kwapinska M. Transformation of inorganic matter in poultry litter during fluidised bed gasification. *Fuel Process Technol* 2021;221:106918.
- Pandey DS, Kwapinska M, Leahy JJ, Kwapinski W. Fly ash from poultry litter gasification—can it be utilised in agriculture systems as a fertiliser? *Energy Procedia* 2019;161:38–46.
- Pandey DS, Pan I, Das S, Leahy JJ, Kwapinski W. Multi-gene genetic programming based predictive models for municipal solid waste gasification in a fluidized bed gasifier. *Bioresour Technol* 2015;179:524–33.
- Guo B, Li D, Cheng C, Lü Z-a, Shen Y. Simulation of biomass gasification with a hybrid neural network model. *Bioresour Technol* 2001;76(2):77–83.
- Xiao G, Ni M-J, Chi Y, Jin B-S, Xiao R, Zhong Z-P, et al. Gasification characteristics of MSW and an ANN prediction model. *Waste Manage* 2009;29(1):240–4.
- Chavan P, Sharma T, Mall B, Rajurkar B, Tambe S, Sharma B, et al. Development of data-driven models for fluidized-bed coal gasification process. *Fuel* 2012;93:44–51.
- Puig-Arnau M, Hernández JA, Bruno JC, Coronas A. Artificial neural network models for biomass gasification in fluidized bed gasifiers. *Biomass Bioenergy* 2013;49:279–89.
- Pandey DS, Das S, Pan I, Leahy JJ, Kwapinski W. Artificial neural network based modelling approach for municipal solid waste gasification in a fluidized bed reactor. *Waste Manage* 2016;58:202–13.
- Baruah D, Baruah D, Hazarika M. Artificial neural network based modeling of biomass gasification in fixed bed downdraft gasifiers. *Biomass Bioenergy* 2017;98:264–71.
- Elmaz F, Yuçel Ö, Mutlu AY. Predictive modeling of biomass gasification with machine learning-based regression methods. *Energy* 2020;191:116541.
- Shahbaz M, Taqvi SA, Loy ACM, Inayat A, Uddin F, Bokhari A, et al. Artificial neural network approach for the steam gasification of palm oil waste using bottom ash and cao. *Renew Energy* 2019;132:243–54.
- Serrano D, Golpour I, Sanchez-Delgado S. Predicting the effect of bed materials in bubbling fluidized bed gasification using artificial neural networks (ANNs) modeling approach. *Fuel* 2020;266:117021.
- Serrano D, Castello D. Tar prediction in bubbling fluidized bed gasification through artificial neural networks. *Chem Eng J* 2020;402:126229.
- Yan Y, Mattisson T, Moldenhauer P, Anthony EJ, Clough PT. Applying machine learning algorithms in estimating the performance of heterogeneous, multi-component materials as oxygen carriers for chemical-looping processes. *Chem Eng J* 2020;387:124072.
- Pan I, Pandey DS. Incorporating uncertainty in data driven regression models of fluidized bed gasification: a bayesian approach. *Fuel Process Technol* 2016;142:305–14.
- Kasivivanathan H, Ng RT, Tay DH, Ng DK. Fuzzy optimization for retrofitting a palm oil mill into a sustainable palm oil-based integrate biorefinery. *Chem Eng J* 2012;200:694–709.
- Patil-Shinde V, Kulkarni T, Kulkarni R, Chavan PD, Sharma T, Sharma BK, et al. Artificial intelligence-based modeling of high ash coal gasification in a pilot plant scale fluidized bed gasifier. *Ind Eng Chem Res* 2014;53(49):18678–89.
- Liu S, Yang Y, Yu L, Zhu F, Cao Y, Liu X, et al. Predicting gas production by supercritical water gasification of coal using machine learning. *Fuel* 2022;329:125478.
- Azadi S, Karimi-Jashni A. Verifying the performance of artificial neural network and multiple linear regression in predicting the mean seasonal municipal solid waste generation rate: a case study of fars province, Iran. *Waste Manage* 2016;48:14–23.
- Arena U, Di Gregorio F. Gasification of a solid recovered fuel in a pilot scale fluidized bed reactor. *Fuel* 2014;117:528–36.
- Baratieri M, Pieratti E, Nordgreen T, Grigiante M. Biomass gasification with dolomite as catalyst in a small fluidized bed experimental and modelling analysis. *Waste Biomass Valoriz* 2010;1(3):283–91.
- Gomez-Barea A, Arjona R, Ollero P. Pilot-plant gasification of olive stone: a technical assessment. *Energy Fuels* 2005;19(2):598–605.
- Serrano D, Sánchez-Delgado S, Horvat A. Effect of sepiolite bed material on gas composition and tar mitigation during c. cardunculus l. gasification. *Chem Eng J* 2017;317:1037–46.
- Serrano D, Kwapinska M, Horvat A, Sanchez-Delgado S, Leahy JJ. Cynara cardunculus l. gasification in a bubbling fluidized bed: The effect of magnesite and olivine on product gas, tar and gasification performance. *Fuel* 2016;173:247–59.
- Arena U, Zaccariello L, Mastellone ML. Fluidized bed gasification of waste-derived fuels. *Waste Manage* 2010;30(7):1212–9.
- M. Campoy Naranjo, Biomass and waste gasification in fluidised bed: pilot plant studies, University of Seville.
- Mansaray K, Ghaly A, Al-Taweel A, Hamdullahpur F, Ugursal V. Air gasification of rice husk in a dual distributor type fluidized bed gasifier. *Biomass Bioenergy* 1999;17(4):315–32.
- Miccio F, Piriou B, Ruoppolo G, Chirone R. Biomass gasification in a catalytic fluidized reactor with beds of different materials. *Chem Eng J* 2009;154(1–3):369–74.
- Narvaez I, Orio A, Aznar MP, Corella J. Biomass gasification with air in an atmospheric bubbling fluidized bed. effect of six operational variables on the quality of the produced raw gas. *Ind Eng Chem Res* 1996;35(7):2110–20.



- [37] Pandey DS, Kwapinska M, Gomez-Barea A, Horvat A, Fryda LE, Rabou LP, Leahy JJ, Kwapinski W. Poultry litter gasification in a fluidized bed reactor: effects of gasifying agent and limestone addition. *Energy Fuels* 2016;30(4):3085–96.
- [38] Roche E, de Andrés JM, Narros A, Rodríguez ME. Air and air-steam gasification of sewage sludge. the influence of dolomite and throughput in tar production and composition. *Fuel* 2014;115:54–61.
- [39] Christodoulou C, Grimekis D, Panopoulos K, Pachatouridou E, Iliopoulou E, Kakaras E. Comparing calcined and un-treated olivine as bed materials for tar reduction in fluidized bed gasification. *Fuel Process Technol* 2014;124:275–85.
- [40] De Andres JM, Narros A, Rodriguez ME. Behaviour of dolomite, olivine and alumina as primary catalysts in air-steam gasification of sewage sludge. *Fuel* 2011;90(2):521–7.
- [41] De Andres JM, Narros A, Rodriguez ME. Air-steam gasification of sewage sludge in a bubbling bed reactor: Effect of alumina as a primary catalyst. *Fuel Process Technol* 2011;92(3):433–40.
- [42] Van Huynh C, Kong S-C. Performance characteristics of a pilot-scale biomass gasifier using oxygen-enriched air and steam. *Fuel* 2013;103:987–96.
- [43] Kaewluan S, Pipatmanomai S. Potential of synthesis gas production from rubber wood chip gasification in a bubbling fluidised bed gasifier. *Energy Convers Manage* 2011;52(1):75–84.
- [44] Kaewluan S, Pipatmanomai S. Gasification of high moisture rubber woodchip with rubber waste in a bubbling fluidized bed. *Fuel Process Technol* 2011;92(3):671–7.
- [45] Lahijani P, Zainal ZA. Gasification of palm empty fruit bunch in a bubbling fluidized bed: a performance and agglomeration study. *Bioresour Technol* 2011;102(2):2068–76.
- [46] Lan W, Chen G, Zhu X, Wang X, Wang X, Xu B. Research on the characteristics of biomass gasification in a fluidized bed. *J Energy Inst* 2019;92(3):613–20.
- [47] Loha C, Chattopadhyay H, Chatterjee PK. Energy generation from fluidized bed gasification of rice husk. *J Renew Sustain Energy* 2013;5(4):043111.
- [48] Lv P, Xiong Z, Chang J, Wu C, Chen Y, Zhu J. An experimental study biomass air-steam gasification in a fluidized bed. *Bioresour Technol* 2004;95(1):95–101.
- [49] Katsaros G, Pandey DS, Horvat A, Almansa GA, Fryda LE, Leahy JJ, et al. Gasification of poultry litter in a lab-scale bubbling fluidised bed reactor: Impact of process parameters on gasifier performance and special focus on tar evolution. *Waste Manage* 2019;100:336–45.
- [50] Raza H, Prasad G, Li Y. EWMA model based shift-detection methods for detecting covariate shifts in non-stationary environments. *Pattern Recogn* 2015;48(3):659–69.
- [51] S. M. Lundberg, S.-I. Lee, A unified approach to interpreting model predictions, in: I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, R. Garnett (Eds.), *Advances in Neural Information Processing Systems*, Vol. 30, Curran Associates, Inc., 2017. URL <https://proceedings.neurips.cc/paper/2017/file/8a20a8621978632d76c43dfd28b67767-Paper.pdf>.
- [52] T. Hastie, R. Tibshirani, J. H. Friedman, J. H. Friedman, *The elements of statistical learning: data mining, inference, and prediction*, Vol. 2, Springer, 2009.
- [53] S. M. Lundberg, G. G. Erion, S.-I. Lee, Consistent individualized feature attribution for tree ensemblesarXiv:1802.03888.
- [54] D. G. Garson, *Interpreting neural network connection weights*.
- [55] Niu M, Huang Y, Jin B, Wang X. Simulation of syngas production from municipal solid waste gasification in a bubbling fluidized bed using aspen plus. *Ind Eng Chem Res* 2013;52(42):14768–75.
- [56] Gómez-Barea A, Leckner B. Modeling of biomass gasification in fluidized bed. *Prog Energy Combust Sci* 2010;36(4):444–509.
- [57] Dupont C, Boissonnet G, Seiler JM, Gauthier P, Schweich D. Study about the kinetic processes of biomass steam gasification. *Fuel* 2007;86(1–2):32–40.