# Property prediction and properties-to-microstructure inverse analysis of steels by a machine-learning approach

Zhi-Lei Wang, Yoshitaka Adachi*

Department of Materials Science and Engineering, Nagoya University, Furo-cho, Chikusa-ku, Nagoya 464-8601, Japan

## ABSTRACT

The design of new materials with useful properties is becoming increasingly important. Machine-learning tools Materials Genome Integration System Phase and Property Analysis (MIPHA) and rMIPHA (based on the *R* programming environment) have been independently developed to accelerate the process of materials discovery via a data-driven materials research approach. In the present work, MIPHA and rMIPHA are applied to steel, where machine-learning-based 2D/3D microstructural analysis, direct analysis of property predictions, and properties-to-microstructure inverse analysis were conducted. The results demonstrate that the prediction models deliver satisfactory performance. The inverse exploration of microstructures related to desired target properties (e.g., stress–strain curve, tensile strength, and total elongation) was realized. MIPHA and rMIPHA are still under improvement. The microstructure-to-processing inverse analysis is expected to be realized in the future.

## 1. Introduction

Science is in an exponential world in which the amount of scientific data is doubly increased every year, which drives the evolution of scientific methods from traditional paper notebooks toward enormous online databases [1]. As data volumes increase, the ability to efficiently extract knowledge from the huge amount of data becomes increasingly important. In response to such a data deluge, the highly efficient and systematic use of databases has become an integral part of the scientific process. Machine learning, which is an artificial intelligence approach to analyzing data and making predictions and decisions based on a huge data volume through various models and algorithms [2,3], has already been successfully applied in many scientific fields [4]. Examples include cognitive game theory (e.g., computer chess) [5,6], pattern recognition (e.g., facial or fingerprint recognition) [7–9], and event forecasting [10].

Because of the staggering compositional and configurational degrees of freedom in materials, the chemical space of materials is far from being exhausted; an enormous number of new materials with useful properties are yet to be discovered [11]. In the traditional experimental science, a material is generally designed from a given chemical composition and processing conditions, followed by microstructure analysis and property evaluation, which is high-cost, low-efficiency, and insufficient for designing a novel material with desired

properties. Therefore, machine learning is now attracting increasing attention in the materials research field to explore unknown information about materials and thus accelerate advances in materials discovery [12,13]. One proposed approach is known as materials informatics, which is scientific and technical and seeks to establish structure–property relationships in a high-throughput, statistically robust, and physically meaningful manner using computational science [13].

The application of machine learning to materials research has led to numerous achievements: predictions of phase diagrams [14], crystal structures [15,16], and materials properties [11,17–19]; developments of interatomic potentials [20–22] and energy functionals [23]; and mapping of materials behavior to process variables [24]. However, these applications are mainly restricted to a direct analysis from structure to property under given chemical compositions and processing conditions. An inverse analysis method that starts from the desired property and predicts the required structural features and processing conditions has not yet been developed.

In the most current materials research, microstructures of materials are studied in two dimensions. However, the two-dimensional (2D) approach gives rise to some criticism because real materials are three-dimensional (3D). For example, flow curves in ferrite–martensite dual phase steel has been reported to be underestimated compared with the flow curves obtained by prediction from 2D plane strain modeling and

* Corresponding author.
  *E-mail address:* adachi.yoshitaka@material.nagoya-u.ac.jp (Y. Adachi).

**Fig. 1.** Functions and characteristics of MIPHA.

those obtained experimentally [25]. Thus, 3D microstructural data of materials appear to be necessary for data-driven materials research.

The present work is aimed at developing a new machine-learning tool, Materials Genome Integration System Phase and Property Analysis (MIPHA), which can realize 2D/3D microstructural analysis and direct and inverse analyses simultaneously. Furthermore, a machine-learning program of rMIPHA based on R script has also been developed; it mainly focuses on analysis of the data obtained from MIPHA. The purpose of this work is to provide materials researchers with a new avenue for data-driven materials design and thus accelerate the materials discovery process.

## 2. MIPHA and rMIPHA

This section introduces MIPHA and rMIHPA. Fig. 1 shows the primary functions and characteristics of MIPHA, including image recognition, image processing, 2D/3D analysis, and direct and inverse analyses. Deep learning [26] and Trainable Weka Segmentation (TWS) [27] approaches are adopted for image recognition and processing functions, respectively. Deep learning can extract image features using an artificial neural network (ANN) with multiple layers, acquiring abstractive features that represent the original image features. Here, GoogLeNet model is used for deep learning, and dropout and average image extraction techniques are employed to suppress overfitting and clear image, respectively. TWS is a machine-learning tool that can realize segmentation of large image datasets automatically after training a classifier by a limited number of manual annotations [27]. Twenty alternative training features (e.g., Gaussian blur, mean, min/max, and anisotropic diffusion) are supplied in TWS to ensure training accuracy. The functions of image recognition and image processing are implemented through a free software of Fiji. In 2D and 3D analysis functions, microstructural characteristics (2D: count fraction, area fraction, circularity, solidity, and ferret's diameter/angle; 3D: count fraction (CF: count/total volume), volume fraction (VF), surface area, Gauss curvature, ferret's diameter, sphericity, genus, Euler-Poincare, piercing particle, isolated inner particle, and branching point) can be analyzed and quantified with help of free and commercial software called the Amira, which is designed for high-dimensional data visualization, processing, and analysis [28].

Property prediction is the main function in direct analysis, where an ANN classifier is used to fit the prediction model. Since an excess of model variables often leads to overfitting [29], a function of data transformation and variable selection [30,31] is installed in the MIPHA. The data of the explanatory variables are first subjected to multiple transformations, such as linear, logarithmic, exponential, square, root,

tangential. A logistic regression is employed to identify the correlation between the explanatory and objective variables, where the transformations that contribute to a high correlation are defined by a gene pattern of 1, while the either ones are defined as 0. Then, a set of efficient transformations that leads to the highest correlation is extracted using a genetic algorithm (GA) [32], which is a metaheuristic inspired by the process of natural selection using for various optimization problems, especially with incomplete or imperfect information or limited computation capacity. The explanatory variables that provide efficient transformations are thereby selected into the input layer of the neural network. In addition, sigmoid function was used as the activation function with the MIPHA in this study. In inverse analysis, a direct analysis model should be established in advance, followed by inverse analysis using the GA, where the population size of 2000, generation of 50, crossover rate of 0.1, and mutation rate of 0.85 were used in this work.

Fig. 2 shows the main functions of rMIPHA that works in the R programming environment, including variable selection, dimension reduction, regression analysis, and direct and inverse analyses. In the variable selection function, the Akaike information criterion (AIC) [33], Bayesian information criterion (BIC) [34], and the least absolute shrinkage and selection operator (LASSO) [35] packages are installed, which are designed for selecting a subset of relevant variables for model construction, so as to simplify the model, shorten the training time, reduce overfitting, as well as make the model easier to interpret. The AIC and BIC are formally as $AIC = 2k – 2ln(L)$ and $BIC = ln(n)k – 2ln(L)$, respectively, where $k$ and $n$ are the number of the explanatory variables and observations estimated by the model, respectively; and $L$ is the maximum value of the likelihood function for the model. The variables that result in the lowest AIC or BIC value for the model are preferred. LASSO is a regression analysis method, which forces regression coefficients of certain variables to be 0, and then effectively chooses a simpler model with those variables whose corresponding absolute values of the coefficient are larger than 0.

In the dimension reduction function, principal component analysis (PCA) [36] and Autoencoder [37] packages are used to convert the high-dimension dataset to a low dimension. PCA normalizes the high-dimension dataset with correlated variables and convert it into a set of linearly uncorrelated vectors that describe the variances of the observations in the dataset, where two principle components PC1 and PC2 are generally used to evaluate the primary variances of the observations. Autoencoder is a type of ANN used to compress a high-dimension dataset into a low-dimension code that can be uncompressed into something closely matching the original dataset.

In the regression analysis function, ANN [38], support vector

**Fig. 2.** Functions and characteristics of rMIPHA.

machine (SVM) [39], random forest (RF) [40], and multiple regression (MR) [41] classifiers with hyper-parameter Bayesian optimization (BO) are installed for fitting data to models. On the basis of the regression models, the property prediction and inverse analysis can be realized using the BO algorithm in direct and inverse analysis functions. In the present work, sigmoid function was used as the activation function in the ANN model. The node number of hidden layer and weight decay, as hyper-parameters, were optimized under a learning rate of 0.01. In the RF model, the hyper-parameters of the numbers of tree and feature in each tree were optimized, where 1000 trees and maximum feature value of 7 were used. In the SVM model, the radial basis function (RBF) was used as the kernel. The penalty coefficient of cost and parameter gamma were optimized in a range of 0.25 ~ 4.

Fig. 3 compares the functions between MIPHA and rMIPHA. MIPHA was developed as dependent software using the Visual Basic language. rMIPHA works in the R language, which is extensively used for statistical computing and data analysis. The functions of image processing and 2D/3D microstructural analysis are unique for MIPHA, whereas rMIHPA shows obvious advantages in regression analysis for its selectable classifiers with hyper-parameter BO. In inverse analysis functions, GA and BO are used for MIPHA and rMIPHA with maximum objective variables of 2 and 3, respectively. In addition, rMIPHA has more options for variable selection and dimension reduction in sparse studies. Further details describing the work of these functions have been introduced in our previous work [42].

## 3. Application of MIPHA and rMIPHA in steels

Mechanical properties are the foundation of various steels and are highly sensitive to their microstructure. Fig. 4 maps the primary microstructural factors in materials that influence their strength and plasticity. These microstructural factors are classified into first descriptors and second descriptors. The former mainly describes characteristics of the second phase, grain size, crystal orientation, grain boundaries, and dislocations. The latter describes factors derived from the former, such as lattice friction, mobile and immobile dislocation densities, residual stress, elastic anisotropy, and Schmidt factor. To thoroughly understand the relationship between microstructure and properties, estimations of such numerous microstructural factors on the basis of traditional experimental science are insufficient. In addition, the microstructure also strongly depends on chemical compositions and processing conditions of the materials, which makes the estimation more difficult. Thus, machine learning is a powerful approach to exploring the potential relationships among processing conditions, microstructure, and mechanical properties.

In this section, MIPHA and rMIPHA are applied to steels. Direct analysis of property prediction and properties-to-microstructure inverse analysis are carried out. One of the objectives is to study the relationship between microstructure and properties by machine learning; the other objective is to demonstrate the functions of MIPHA and rMIPHA.

### 3.1. Experimental procedure

Cold-rolled (CR) low-carbon steels with different chemical

| Functions | | MIPHA | rMIPHA |
|---|---|---|---|
| Language | | Visual basic | R |
| Image recognition | | O Deep learning | ✕ |
| Image processing | | O Machine-learning-based | ✕ |
| 2D analysis | | O | ✕ |
| 3D analysis | | O | ✕ |
| Direct analysis | Classifier | O ANN | O ANN, SVR, RF, Multiple regression |
| | Hyper-parameter optimization | ✕ | O Bayesian optimization |
| Inverse analysis | Exploration method | O Genetic Algorithm | O Bayesian optimization |
| | Exploration target | O 1-2 objects; sum, product, max, min and specified value exploration | O 1-3 objects; product, max and specified value exploration |
| Sparse study | Variable selection | O Data conversion and selection, ANN according to sensitive coefficient | O AIC, BIC  LASSO, ANN according to weight coefficient, RF according to IncNodePurity |
| | Dimension reduction | ✕ | O Principal component analysis, Auto encoder |

**Fig. 3.** Function comparison between MIPHA and rMIPHA.

**Fig. 4.** Microstructural factors that influence the mechanical properties of materials.

**Table 1**
Chemical compositions and processing conditions of the studied steels.

| Steel | Chemical composition (wt%, N, O: ppm) | Process |
|---|---|---|
| A10-01 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→ annealed at 1000 °C for 5 s→ cooling at 1 °C/s |
| A10-03 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→ annealed at 1000 °C for 5 s→ cooling at 3 °C/s |
| A10-10 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→ annealed at 1000 °C for 5 s→ cooling at 10 °C/s |
| A10-30 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→ annealed at 1000 °C for 5 s→ cooling at 30 °C/s |
| A14-01 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→annealed at 1400 °C for 5 s→cooling to 1000 °C at 50 °C/s→cooling at 1 °C/s |
| A14-03 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→ annealed at 1400 °C for 5 s→ cooling to 1000 °C at 50 °C/s→ cooling at 3 °C/s |
| A14-10 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→ annealed at 1400 °C for 5 s→ cooling to 1000 °C at 50 °C/s→ cooling at 10 °C/s |
| A14–30 | 0.152C-0.015Si-1.51Mn-0.007P-0.0016S-0.027Al-18N-28O | CR→ annealed at 1400 °C for 5 s→ cooling to 1000 °C at 50 °C/s→ cooling at 30 °C/s |
| B10-01 | 0.151C-0.013Si-1.53Mn-0.007P-0.002S-0.193Mo-0.028Al-21N-21O | CR→ annealed at 1000 °C for 5 s→ cooling at 1 °C/s |
| B10-03 | 0.151C-0.013Si-1.53Mn-0.007P-0.002S-0.193Mo-0.028Al-21N-21O | CR→ annealed at 1000 °C for 5 s→ cooling at 3 °C/s |
| B10-10 | 0.151C-0.013Si-1.53Mn-0.007P-0.002S-0.193Mo-0.028Al-21N-21O | CR→ annealed at 1000 °C for 5 s→ cooling at 10 °C/s |
| B10-30 | 0.151C-0.013Si-1.53Mn-0.007P-0.002S-0.193Mo-0.028Al-21N-21O | CR→ annealed at 1000 °C for 5 s→ cooling at 30 °C/s |

compositions and processing conditions were studied in the present work. The CR steel samples were austenitized at 1000 °C or 1400 °C and cooled at 1, 3, 10, or 30 °C/s to room temperature. The chemical compositions of the raw materials and processing parameters are detailed in Table 1.

Continuous cooling transformation (CCT) curves were measured by the thermal expansion method. Tensile tests were performed to evaluate the mechanical properties of the samples. The microstructures of the samples were observed on their sections parallel to the rolling direction using a proprietary serial-sectional 3D microscope (Genus_3D) [43], where approximately 100 images were serially observed at 0.53 ~ 0.96 μm per interval for each sample. Microstructural analysis was performed by MIPHA. Direct analysis of property prediction and properties-to-microstructure inverse analysis were carried out with MIPHA and rMIPHA.

### 3.2. Machine-learning-based microstructural analysis

Fig. 5 illustrates the 2D and 3D microstructures of sample A10-01.

Fig. 5(a) shows an image as an example observed using Genus_3D. According to the contrast, morphology, and CCT curve, the microstructure was recognized as being composed of four phases: polygon ferrite (PF), Widmanstatten ferrite (WF), pearlite (P), and degenerated pearlite (DP), which were observed as white polygonal, white acicular, dark, and light features, respectively. In addition, in the samples cooled at higher rates of 10 and 30 °C/s, bainite (B) and martensite (M) were observed. Fig. 5(b) shows a cropped image with local contrast normalization corresponding to the area highlighted by the red box in subfigure (a). The cropping and local contrast normalization were carried out to ensure that the subsequent phase segmentation proceeded well. Fig. 5(c) shows a phase-extracted image corresponding to (b), in which the four aforementioned phases are marked in different colors. Fig. 5(d) and (e) show the 3D images reconstructed from the serial images in (b) and (c), which intuitively and proximately present the real microstructure of the sample. Furthermore, Fig. 5(f) shows a 3D image segmented from (e) as an example, clearly displaying the morphology and distribution of the P phase.

In the present work, the aforementioned 2D and 3D microstructure

**Fig. 5.** Microstructures of sample A10-01: (a) an original image observed by Genus_3D; (b) an image with local contrast normalization corresponding to the area highlighted by the red box in (a); (c) a phase-extracted image corresponding to (b); (d) and (e) reconstructed 3D images from the serial images of (b) and (c), respectively; and (f) a 3D image segmented from (e) with P phase only.

characteristics of each phase were analyzed and quantified with help of deep learning and TWS, supplying sufficient topological microstructure information that approaches to a real material. The obtained average information of each microstructure feature was automatically summarized to a CSV file. In addition, the microstructure information of each phase was also statistic to separated CSV files. This indicates that MIPHA has a powerful 2D/3D microstructural analysis function.

### 3.3. Construction of datasets

In order to avoid a complex model and reduce overfitting resulting from excess model variables, BIC estimation was performed to identify the importance of the above 3D microstructure features for the property (stress). The results demonstrated that the CF and VF of the most of phases exhibited high importance for the stress, by which the CF and VF were thereby chosen as the microstructure features for fitting the models in this study.

The quantitative CF and VF of each phase are detailed in Table 2. In addition, the mechanical properties of tensile strength (TS) and total elongation (tEL) estimated from the stress-strain curves are also included (the experimental information of strain and stress is listed in Table S1). These obtained microstructure and property data constitute the "material genomes" used for subsequent direct and inverse analyses. In the present work, two datasets were constructed: one was used for predicting stress-strain curve, and the other was used for inversely exploring a balanced property of TS and tEL. The former contained 111 observations depending on the number of the overall strain/stress data

items of the samples (Table S1) with 14 features in each observation (6 CFs, 6 VFs, strain, and stress). The latter contained 12 observations depending on the number of the studied samples (Table 2) with 14 features in each observation (6 CFs, 6 VFs, TS, and tEL).

In machine learning, overfitting often occurs when a statistical model accurately fits the data at hand but fails to describe the underlying data, which results in inaccurate predictions for the novel material characteristics. One approach to avoiding overfitting is to separate the datasets for training a model and for testing it. Therefore, in the present work, the 75% of the data in each dataset was used for training and the remaining 25% was used for testing. It should be pointed out that the training and testing datasets were split using a round-robin algorithm in MIPHA, while they were randomly split in rMIPHA with a 10-fold cross validation for the training data.

### 3.4. Direct analysis of property prediction

Fig. 6 shows the direct analysis results obtained by MIPHA without variable selection, including the neural network of the prediction model and the predicted stress–strain curve of the A10-01 sample in its plastic deformation period. As shown in Fig. 6(a), in this prediction model, all microstructures (CF and VF) and true strain were used as the explanatory variables (input layer) and the true stress was used as the objective variable (output layer). A hidden layer with nine variables was created between the input and output layers. The correlation coefficients (CCs) of the training and testing datasets were evaluated as 0.98987 and 0.96062 for the present model, indicating a good linear

**Table 2**
Mechanical properties and quantitative microstructures of the samples.

| Steel | YS(MPa) | TS(MPa) | tEL(%) | CFPF | CFP | CFWF | CFDP | CFB | CFM | VFPF | VFP | VFWF | VFDP | VFB | VFM |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| A10-01 | 323 | 481 | 80.6 | 3.30E-05 | 5.16E-05 | 0.000165 | 1.83E-05 | 0 | 0 | 0.4047 | 0.2005 | 0.0845 | 0.3104 | 0 | 0 |
| A10-03 | 308 | 489 | 76.4 | 5.43E-05 | 9.07E-05 | 1.28E-06 | 9.24E-05 | 0 | 0 | 0.2608 | 0.118 | 0.5537 | 0.0674 | 0 | 0 |
| A10-10 | 390 | 591 | 71.1 | 5.17E-05 | 0.000136 | 0.000126 | 0 | 1.54E-06 | 0 | 0.1836 | 0.0452 | 0.1414 | 0 | 0.6297 | 0 |
| A10-30 | 444 | 663 | 63.9 | 8.63E-05 | 0 | 0 | 0.00027 | 9.26E-07 | 4.42E-05 | 0.1576 | 0 | 0.0842 | 0 | 0.5765 | 0.1817 |
| A14-01 | 353 | 516 | 64.4 | 7.67E-05 | 5.15E-05 | 4.04E-06 | 3.21E-05 | 5.30E-06 | 0 | 0.1573 | 0.0212 | 0.3938 | 0.0379 | 0.3897 | 0 |
| A14-03 | 412 | 561 | 67.5 | 4.40E-05 | 4.67E-05 | 3.12E-05 | 5.94E-05 | 5.93E-06 | 0 | 0.0808 | 0.0143 | 0.2572 | 0.1232 | 0.5245 | 0 |
| A14-10 | 521 | 688 | 61.5 | 0 | 2.27E-05 | 0 | 0 | 1.08E-05 | 1.73E-05 | 0 | 0.0094 | 0 | 0 | 0.6249 | 0.3657 |
| A14–30 | 620 | 807 | 60.7 | 0 | 0 | 0 | 0 | 0 | 3.00E-05 | 0 | 0 | 0 | 0 | 0 | 1 |
| B10-01 | 375 | 550 | 70.4 | 0.002244 | 0.000637 | 0.001005 | 0 | 0 | 0 | 0.373652 | 0.06523 | 0.561117 | 0 | 0 | 0 |
| B10-03 | 434 | 600 | 66.3 | 0.00325 | 0.000344 | 0.005683 | 0 | 0.000477 | 0 | 0.109254 | 0.006947 | 0.022508 | 0 | 0.861291 | 0 |
| B10-10 | 483 | 691 | 61.5 | 0.000185 | 0.000215 | 0 | 0 | 2.76E-05 | 0 | 0.118045 | 0.006877 | 0 | 0 | 0.875078 | 0 |
| B10–30 | 489 | 725 | 58.4 | 0 | 0 | 0 | 0 | 7.43E-05 | 1.52E-07 | 0 | 0 | 0 | 0 | 0.160882 | 0.839118 |

correlation between the experimental data and the estimated data. As illustrated in Fig. 6(b), the experimental and predicted curves of sample A10-01 were well fit to each other. In addition, by comparison for all samples, the experimental and predicted curves were still keep a good fitness (Fig. S1). These results suggest satisfactory performance of this model.

As mentioned above, an excess of model variables can also lead to overfitting. Therefore, a prediction model with variable selection was also established by MIPHA. Fig. 7(a) shows the neural network of the fitted model. The variables in the hidden layer were reduced to seven after variable selection; however, this simplified model still had an accuracy approximately equal to that achieved without variable selection (Fig. 6). A comparison of the experimental and predicted stress–strain curves also indicates good performance of this model, as illustrated by the A10-01 sample in Fig. 7(b) (Fig. S2 presents the predicted results of all samples). The aforementioned results demonstrate that variable selection is beneficial in the case of numerous model variables.

The obtained "materials genomes" were also studied by rMIPHA using different classifiers with and without variable selection. Fig. 8 shows the direct analysis results without variable selection. The dataset was pre-estimated by ANN, SVM, and RF classifiers with hyper-parameter BO. Fig. 8(a) shows the performance of the fitted models, as indicated by CC and root-mean-square error (RMSE). The results

demonstrate that the ANN model exhibited the best accuracy on the basis of its high CC and low RMSE. By contrast, substantial overfitting occurred in the SVM model. Fig. 8(b) shows the hyper-parameter BO result for the ANN model under the search conditions of 20 initial points and 10 iterations. The size of 4 and decay of 0.0091 are the best-fit hyper parameters for model, as indicated by the lowest RMSE. Fig. 8(c) shows the neural network of the ANN model, which describes the degree of sensitivity of objective variables to explanatory variables. Red and blue colors express positive and negative sensitivity, respectively, and a wider connection line expresses a larger value. The quantitative degrees of sensitivity of the objective variable to each explanatory variables are listed in Table 3; these values were automatically generated during the model fitting process. The results show that the explanatory variable of true strain was the factor most sensitive to the objective variable of true stress. Fig. 8(d) illustrates the experimental and ANN-predicted true stress–strain curves of the A10-01 sample. The experimental and predicted curves are shown to almost coincide. Similar predictions were also almost observed in the remaining samples (Fig. S3), which indicates excellent model performance resulting from the hyper-parameter BO.

Fig. 9 shows the direct analysis results with variable selection. Here, BIC was adopted to evaluate the degree of importance of the explanatory variables. Fig. 9(a) shows the results of BIC variable selection



**Fig. 6.** Direct analysis results without variable selection by MIPHA: (a) neural network of the fitted model and (b) experimental and predicted true stress−strain curves of sample A10-01.

**Fig. 7.** Direct analysis results with variable selection by MIPHA: (a) neural network of the fitted model and (b) experimental and predicted true stress−strain curves of sample A10-01.



**Fig. 8.** Direct analysis results without variable selection by rMIPHA: (a) performance comparison of different models; (b) hyper-parameter BO result for the ANN model; (c) neural network of the ANN model; and (d) experimental and predicted true stress−strain curves of sample A10-01.

**Table 3**
Sensitive degrees of the objective variable to explanatory variables.

| Variable | CFPF | CFP | CFWF | CFDP | CFB | CFM | VFPF | VFP | VFWF | VFDP | VFB | VFM | True strain |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Degree of sensitivity | 0.51599 | 0.54173 | 0.31269 | 0.54976 | 0.50816 | 0.37722 | 0.18469 | 0.98668 | 0.35119 | 0.17412 | 0.44629 | 0.33054 | 1.02794 |



**Fig. 9.** Direct analysis results with variable selection by rMIPHA: (a) the result of BIC variable selection and performance estimation of the BIC variables by ANN; (b) a performance comparison of different models; (c) the degree of importance of explanatory variables evaluated by an RF classifier; and (d) experimental and predicted true stress−strain curves of sample A10-01.

**Table 4**
Inversely explored microstructure related to the explored stress−strain curve in Fig. 10.

| CFPF | CFP | CFWF | CFDP | CFB | CFM | VFPF | VFP | VFWF | VFDP | VFB | VFM |
|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.000975 | 0.000153 | 0.003751 | 0.000085 | 0.000448 | 3.54E-06 | 0.012221 | 0.021191 | 0.160972 | 0.182781 | 0.079276 | 0.543559 |



**Fig. 10.** Comparison of the inversely explored and target stress−strain curves.

(highlighted in the blue box) and a performance estimation of the BIC variables. Notably, this model was fitted by ANN without hyper-parameter optimization. The results show that the model still provided a satisfactory performance indicated by its high CC and low RMSE values

for the training data and testing data. However, when the BIC data was estimated by ANN, SVM, and RF with hyper-parameter BO, RF became the best model, as shown in Fig. 9(b). Fig. 9(c) shows the degree of importance explanatory variables evaluated by an RF classifier. True strain is shown to be the most important variable for true stress in this model, as indicated by its largest IncNodePurity (Increase of Node Purity: an index to express the variable importance). Fig. 9(d) illustrates the experimental and RF-predicted true stress–strain curves for the A10-01 sample. The predictions for all samples are shown in Fig. S4. By comparison, a satisfactory result was still obtained although it was not as good as that achieved without variable selection (Fig. 8).

### 3.5. Properties-to-microstructure inverse analysis

Because of longer training time for finding the best hype-parameters using BO, in this work, inverse analysis was conducted by MIPHA using GA with exploration targets: stress–strain curve, and TS/tEL.

#### 3.5.1. Exploration of a target stress–strain curve

As an example, a target stress–strain curve was arbitrarily written. The stress–strain prediction model in Fig. 6 was inversely analyzed with

**Table 5**
Inversely explored potential TS and tEL as well as the corresponding microstructure.

| CFPF | CFP | CFWF | CFDP | CFB | CFM | VFPF | VFP | VFWF | VFDP | VFB | VFM | TS (MPa) | tEL (%) | TS × tEL |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 0.002015 | 0.000612 | 0.002614 | 1.29E- 05 | 0.000372 | 3.45E- 05 | 0.223769 | 0.113381 | 0.014103 | 0 | 0.03299 | 0.615758 | 780.662 | 79.80517 | 62,300.86 |

a target stress–strain search, which was to explore the microstructure candidate that relates to the written curve. As a result, the microstructure corresponding to the target curve was obtained, as listed in Table 4. Moreover, the explored and target stress–strain curves well fit with each other, as shown in Fig. 10.

### 3.5.2. Exploration of targets TS and tEL

Generally, the quality of steel is evaluated by its strength and plasticity, which are characterized by TS and tEL, respectively. Therefore, a direct analysis model with explanatory variables of microstructure and objective variables of TS and tEL was first established by MIPHA; the resultant model showed a CC of 0.95957. This model was then inversely analyzed with a TS × tEL maximum search, which was designed to explore the microstructure that relates to the best balance of strength and plasticity. Table 5 lists the explored potential TS and tEL as well as the corresponding microstructure. The potential TS and tEL are much higher than the experimental results listed in Table 2. The potential TS × tEL can reach 62,300.86, which is 1.27 times larger than the largest experimental result of 48,984.90 (A14–30). In addition, in the explored microstructure, hard phase M and soft phase PF can be considered the primary phases that impart better strength and plasticity to the present steels.

It should be pointed out that the given examples of inverse analysis here explored the microstructures corresponding to desired properties. However, systematical evaluation of the inverse analysis model performance and microstructure-to-processing inverse analysis were not performed in this work restricted by the present functions of MIPHA and rMIPHA, which are still yet under improvement. A properties-to-microstructure-to-processing inverse analysis with evaluation of model performance will be demonstrated in future work. Moreover, the explored results of the properties, and their corresponding microstructure and processing will also be evaluated by both experiment and finite element method [44] in the future.

In this study, data science was applied to steels, which exhibited remarkable advantages compared to the experimental science, such as savings of labor, time and cost, and a more thorough estimation of the relationship between the microstructure and properties. In particular, the proposed properties-to-microstructure inverse analysis explored the potential properties of the studied steels as well as the corresponding microstructure. Since microstructure is a junction connecting the processing and properties, an adequate properties-to-microstructure-to-processing inverse analysis is expected to effectively accelerate the materials discovery process.

## 4. Conclusions

Independently developed machine-learning tools MIPHA and rMIPHA were applied in steels, where machine-learning based microstructural analysis, property prediction, and properties-to-microstructure inverse analysis were conducted. The microstructural components of the samples were quantified, constituting the "materials genomes". Stress–strain curves were predicted on the basis of the materials genomes. The prediction models showed satisfactory accuracies. The microstructures corresponding to desired properties (a target stress–strain curve and target TS/tEL) were inversely explored by MIPHA successfully, where the explored and the target stress–strain curves well matched each other; and the inversely explored potential TS and tEL were much larger than the experimental results. The results presented in this work are expected to provide a new approach in materials design to accelerate the materials discovery process.

## Appendix A. Supporting information

Supplementary data associated with this article can be found in the online version at doi:10.1016/j.msea.2018.12.049.

## References

[1] S. Alexander, J. Gray, 2020 Computing: Science in an exponential world, Nature 440 (2006) 413–414.
[2] R.S. Michalski, J.G. Carbonell, T.M. Mitchell, Machine Learning: an Artificial Intelligence Approach, Springer Science and Business Media, Palo Alto, 2013, pp. 41–81.
[3] N.M. Nasrabadi, Pattern recognition and machine learning, J. Electron. Imaging 16 (2007) 049901.
[4] A.L. Samuel, Computer Games I, Springer, New York, 1988, pp. 335–365.
[5] J.H. Holland, Emergence: from Chaos to Order, OUP, Oxford, 2000, pp. 16–26.
[6] N. Jones, Quiz-playing computer system could revolutionize research, Nat. News (2011), https://doi.org/10.1038/news.2011.95.
[7] N. MacLeod, M. Benfield, P. Culverhouse, Time to automate identification, Nature 467 (2010) 154–155.
[8] J.P. Crutchfield, Between order and chaos, Nat. Phys. 8 (2012) 17–24.
[9] L. Chittka, A. Dyer, Your face looks familiar, Nature 481 (2012) 154–155.
[10] W.C. Hong, Rainfall forecasting by technological machine learning models, Appl. Math. Comput. 200 (2008) 41–57.
[11] G. Pilania, C. Wang, X. Jiang, S. Rajasekaran, R. Ramprasad, Accelerating materials property predictions using machine learning, Sci. Rep. 3 (2013) 2810.
[12] R. LeSar, Materials informatics: an emerging technology for materials development, Stat. Anal. Data Min. 1 (2009) 372–374.
[13] K. Rajan, Materials informatics, Mater. Today 8 (2005) 38–45.
[14] C.J. Long, J.H. Simpers, M. Murakami, R.C. Srivastava, I. Takeuchi, V.L. Karen, X. Li, Rapid structural mapping of ternary metallic alloy systems using the combinatorial approach and cluster analysis, Rev. Sci. Instrum. 78 (2007) 072217.
[15] G. Hautier, C.C. Fischer, A. Jain, T. Mueller, G. Ceder, Finding nature's missing ternary oxide compounds using machine learning and density functional theory, Chem. Mater. 22 (2010) 3762–3767.
[16] D. Morgan, S. Curtarolo, K. Persson, J. Rodgers, G. Ceder, Predicting crystal structures with data mining of quantum calculations, Phys. Rev. Lett. 91 (2003) 135503.
[17] K. Hansen, G. Montavon, F. Biegler, S. Fazli, M. Rupp, M. Scheffler, O.A. von Lilienfeld, A. Tkatchenko, K.R. Müller, Assessment and validation of machine learning methods for predicting molecular atomization energies, J. Chem. Theory Comput. 9 (2013) 3404–3419.
[18] K. Hansen, F. Biegler, R. Ramakrishnan, W. Pronobis, O.A. von Lilienfeld, K.R. Müller, A. Tkatchenko, Machine learning predictions of molecular properties: accurate many-body potentials and nonlocality in chemical space, J. Phys. Chem. Lett. 6 (2015) 2326–2331.
[19] T.D. Huan, A.M. Kanakkithodi, R. Ramprasad, Accelerated materials property predictions and design using motif-based fingerprints, Phys. Rev. B 92 (2015) 014106.
[20] T. Morawietz, J. Behler, A density-functional theory-based neural network potential for water clusters including van der Waals corrections, J. Phys. Chem. A 117 (2013) 7356–7366.
[21] J. Behler, Neural network potential-energy surfaces in chemistry: a tool for large-scale simulations, Phys. Chem. Chem. Phys. 13 (2011) 17930–17955.
[22] A.P. Bartók, M.C. Payne, R. Kondor, G. Csányi, Gaussian approximation potentials: the accuracy of quantum mechanics without the electrons, Phys. Rev. Lett. 104 (2010) 136403.
[23] J.C. Snyder, M. Rupp, K. Hansen, K.R. Müller, K. Burke, Finding density functionals with machine learning, Phys. Rev. Lett. 108 (2012) 253002.
[24] H.K.D.H. Bhadeshia, Neural networks and information in materials science, Stat. Anal. Data Min. 1 (2009) 296–305.
[25] C. Thomser, Modelling of the Mechanical Properties of Dual Phase Steels Based on Microstructure (Ph.D. Thesis), RWTH-Aachen, Germany, 2009.
[26] Y. LeCun, Y. Bengio, G. Hinton G, Deep learning, Nature 521 (2015) 436–444.

[27] I. Arganda-Carreras, V. Kaynig, C. Rueden, K.W. Eliceiri, J. Schindelin, A. Cardona, H. Sebastian Seung, Trainable weka segmentation: a machine learning tool for microscopy pixel classification, Bioinformatics 33 (2017) 2424–2426.

[28] D. Stalling, M. Westerhoff, H.C. Hege, Amira: a highly interactive system for visual data analysis, The Visualization Handbook 38 (2005), pp. 749–767.

[29] N. Wagner, J.M. Rondinelli, Theory-guided machine learning in materials science, Front. Mater. 3 (2016) 28.

[30] J.W. Tukey, Exploratory Data Analysis, Addison-Wesley, Indianapolis, 1977.

[31] S. Chatterjee, B. Price, Regression Analysis by Example, John Wiley & Sons, New Jersey, 1991.

[32] M. Mitchell, An Introduction to Genetic Algorithms, MIT Press, Cambridge, 1998.

[33] H. Akaike, A new look at the statistical model identification, IEEE Trans. Automat. Control 19 (1974) 716–723.

[34] H.S. Bhat, N. Kumar, On the Derivation of the Bayesian information criterion, School of Natural Sciences, University of California, 2010.

[35] R. Tibshirani, Regression shrinkage and selection via the lasso, J. R. Stat. Soc. B 58 (1996) 267–288.

[36] H. Hotelling, Analysis of a complex of statistical variables into principal components, J. Educ. Psychol. 24 (1993) 417–441.

[37] G. Hinton, R. Salakhutdinov, Reducing the dimensionality of data with neural networks, Science 313 (2006) 504–507.

[38] R.J. Schalkoff, Artificial Neural Networks 1 McGraw-Hill, New York, 1997.

[39] M.A. Hearst, S.T. Dumais, E. Osuna, J. Platt, B. Scholkopf, Support vector machines, IEEE Intell. Syst. Appl. 13 (1998) 18–28.

[40] A. Liaw, W. Matthew, Classification and regression by random forest, R News 2 (2002) 18–22.

[41] J. Cohen, P. Cohen, S.G. West, L.S. Aiken, Applied Multiple Regression/correlation Analysis for the Behavioral Sciences, Routledge, London, 2013.

[42] Y. Adachi, Cutting edge of steel informatics and future prospects, ISIJ Newslett. 23 (6) (2018), ⟨https://www.researchgate.net/publication/325596024_Cutting_Edge_of_Steel_Informatics_and_Future_Prospects⟩.

[43] Y. Adachi, N. Sato, M. Ojima, M. Nakayama, Y.T. Wang, Development of fully automated serial-sectioning 3D microscope and topological approach to pearlite and dual-phase microstructure in steels, Proc. First Int. Conf. 3D Mater. Sci. (2012) 37–42.

[44] O.C. Zienkiewicz, R.L. Taylor, The Finite Element Method, McGraw-hill, London, 1977.