

# Machine Learning-based estimation and explainable artificial intelligence-supported interpretation of the critical temperature from magnetic *ab initio* Heusler alloys data

Robin Hilgers,<sup>1,2,\*</sup> Daniel Wortmann,<sup>1</sup> and Stefan Blügel<sup>1,2</sup>

<sup>1</sup>*Peter Grünberg Institute and Institute for Advanced Simulation,  
Forschungszentrum Jülich and JARA, 52425 Jülich, Germany*

<sup>2</sup>*Department of Physics, RWTH Aachen University, Aachen, Germany*

(Dated: November 28, 2023)

Machine Learning (ML) has impacted numerous areas of materials science, most prominently improving molecular simulations, where force fields were trained on previously relaxed structures. One natural next step is to predict material properties beyond structure. In this work, we investigate the applicability and explainability of ML methods in the use case of estimating the critical temperature ( $T_c$ ) for magnetic Heusler alloys calculated using *ab initio* methods determined materials-specific magnetic interactions and a subsequent Monte Carlo (MC) approach. We compare the performance of regression and classification models to predict the range of the  $T_c$  of given compounds without performing the MC calculations. Since the MC calculation requires computational resources in the same order of magnitude as the density-functional theory (DFT) calculation, it would be advantageous to replace either step with a less computationally intensive method such as ML. We discuss the necessity to generate the magnetic *ab initio* results, to make a quantitative prediction of the  $T_c$ . We used state-of-the-art explainable artificial intelligence (XAI) methods to extract physical relations and deepen our understanding of patterns learned by our models from the examined data.

## I. INTRODUCTION

Machine-Learning (ML) modeling has shown to yield promising results in various scientific sectors and applications [1–3]. The ability of flexible learning algorithms to recognize patterns, adapt to data properties, and tackle challenges such as regression, classification, and clustering has established an additional scientific paradigm of data-driven science besides the traditional paradigms of experiments, theories, and simulations. Data-driven science essentially shifts scientific problem-solution strategies for predictions from problem-specific models to versatile data-based models [4–6]. This is also the case for a plurality of materials science applications ranging from superconductivity [7], molecular dynamics [8], materials synthesis, and design [9], knowledge discovery through data mining [10], entropy changes [11], and other topics for both properties and materials prediction [5, 12, 13]. For some of the mentioned applications, *e.g.* in some molecular dynamics simulation applications [8], lightweight and computationally inexpensive ML-based approaches were able to virtually replace established techniques, while in other applications ML-based approaches complement existing methodologies [5]. Data mining-related techniques have shown to be powerful tools in the hands of scientists to discover relations within data, even in the materials science community [10].

There are a multitude of magnetic properties to investigate, many of which are traditionally described by complex models based in part on the quantum mechanics of the many-electron problem. Within the set of magnetic

properties, the critical temperature, also known as the Curie temperature in the context of ferromagnetic materials, represents a key characteristic in both fundamental physics and practical applications. It provides valuable insights into the transitions between different magnetic phases and guides the design and optimization of magnetic materials for technological use. For example, in the design of magnetic materials for the energy use sector of the economy [14], *e.g.* electric power generation, conditioning, conversion, transportation, or the information sector of the economy, *e.g.* spintronics [15] or magnetic storage devices (like hard drives), the critical temperature determines the maximum operating temperature where magnetic data storage remains stable. Typical application demands necessitate critical temperature values significantly exceeding room temperature [16]. Hence, in order to conduct application-oriented material screening studies at a high-throughput scale for materials discovery, a lightweight method is required to predict whether the critical temperature of a compound meets the requirements set by the applications. Existing works, mostly focused on the Curie temperature in ferromagnetic materials [17, 18], while the more general concept describing a wide range of magnetic phases, including ferromagnetic, anti-ferromagnetic, ferrimagnetic, and spin-spiral type ordering is the critical temperature of the phase change transition of the ordered magnetic to a non-magnetic state represents the field of interest in this study.

Within the phase space of magnetic materials, the Heusler (and Heusler-like alloys) alloys are known to represent candidate materials for various technical applications, as the material class of Heusler [19, 20] alloys (as *e.g.* the ordered  $L2_1$  phase) and related disordered phases (such as *e.g.*  $A2$  and  $B2$  phases) are known to exhibit many interesting properties including superconductivity [21], piezoelectricity [22], rare-earth free perma-

---

\* robin.hilgers@rwth-aachen.de

nent magnets [23], and half-metallicity [24]. The combination of multiple properties in a single compound such as *e.g.* both half-metallicity and magnetic stability allow for the occurrence of spin-polarized charge currents, which are a topic that is actively investigated by the scientific community for applications in spintronics [25, 26]. By including not only the ordered but also disordered phases and quaternary Heusler alloys, the phase space of possible compounds increases drastically in comparison to existing works like *e.g.* [27], which restricts the phase space to pure transition-metal Heusler alloys. However, as a Heusler alloy’s structure is defined by the individual compound’s lattice site constituents, the lattice constant, and the symmetry group alone, the structural parameters that have to be considered by a model in order to describe such a system are very limited.

In this paper, we aim to demonstrate the advantages offered by ML, replacing traditional  $T_c$  determination using density-functional theory (DFT) and Monte Carlo (MC) simulations. We focus on the prediction of the magnetic critical temperature for ordered (Including the phases L2<sub>1</sub>, C1<sub>b</sub>, Y, and XA) as well as disordered (Including the phases A2 and B2) magnetic Heusler alloys. The critical temperatures were determined in a two-step process of an *ab initio* KKR-GF [28] DFT simulation followed by an MC simulation of the  $T_c$  as depicted in the top path of Fig. 1. As both steps are comparable in computational cost, we apply our modeling for the whole process as well as only the MC step, taking advantage of magnetic results obtained in the *ab initio* step.

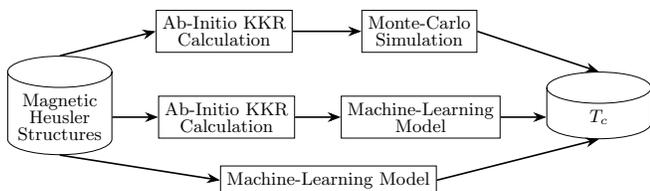


FIG. 1. Schematic depiction of the layered  $T_c$  determination with different ML integration levels.

Beyond that, we discuss the impact of magnetic features for the prediction of high  $T_c$  materials and the usability in high-throughput materials screening applications, which do not include DFT-originated features in the first place. This discussion is heavily assisted by the use of explainable artificial intelligence (XAI) techniques, which we demonstrate to be able to explain model predictions based on materials science data and visualize relations in the training data captured by the ML model [29].

## II. METHODS AND MATERIALS

### A. Data Processing & Cleaning

The examined data was collected at our institute and published as the Jülich-Heusler-magnetic-database

(JuHemd) [30]. It provides not only structural and stoichiometric information on the Heusler compounds, but also magnetic data obtained by DFT and Monte Carlo simulations. The target quantity we want to predict in our modeling is the critical temperature  $T_c$  of the magnetic ordering. While the JuHemd contains experimental values as well as those based on DFT simulations using GGA [31] and LDA [32] exchange-correlation functionals, we restrict our analysis to the GGA-based values as these are provided for most compounds and provide the most homogeneous data quality.

As a first preparation step, we extract the  $T_c$  values together with a set of descriptors for each compound in the database. All information was encoded into a numerical representation and made available for the modeling process. Using the provided metadata to augment the information with additional atomic features, we finally obtain a set of 118 descriptors, as listed in table I. Before any ML modeling is performed, these descriptors  $\{x_i\}$  are then transformed to a standardized form

$$\{z_i\} = \frac{\{x_i\} - \mu_i}{\sigma_i} \quad (1)$$

using the mean  $\mu_i$  and standard deviation  $\sigma_i$  of the  $i$ -th descriptor in the training set.

Only those compound entries have been included which contain all of the above-mentioned entry labels. Incomplete data points have not been used. Additionally, only magnetic alloys are selected. We chose the magnetic cutoff to be

$$\sum_i |m_i| > 0.1\mu_B \quad (2)$$

where the  $m_i$  denotes the magnetic moment of the atom on site  $i$  in the compound’s molecular formula. Similarly, we did not include compounds with a simulated  $T_c = 0$  K. This leaves us with a final data set size of 408 Heusler compounds.

Since, during the data processing, incomplete data points for Heusler compounds are removed, there are some elements from the periodic table that are contained in the original JuHemd but are not contained anymore in the processed data. The corresponding densities of these atomic numbers, which originate from these removed elements, represent descriptors with zero variance in every compound. Such descriptors are removed before further processing, as they are meaningless for the ML training and evaluation process. In this paper, of the 118 descriptors, there are 11 descriptors in the data set with zero variance, which are hence removed. The whole data order has been randomized in order to avoid the clustering of similar data points due to the alphabetical order. This enforces homogeneity of the data set, which is necessary

TABLE I. List of all features which are contained in the processed data and their corresponding explanation. For all features that were directly derived from the JuHemd, the JuHemd label has been used. Also, JuHemd labels have been included which were used to construct processed quantities even though the original label is not included in the processed data set due to the format, the quantity is given in the JuHemd.

Label	Description
lattice_constant*	Lattice constant of the Heusler
resval*	$T_c$
etotal (Ry) *	Total energy of the compound $E_{Tot}$
formula*	Chemical formula of the compound
Ferro Density†	Fraction of ferromagnetic elements (Fe, Ni, Co) in the Compound
Rare earth Materials Density†	Fraction of rare earth components in the Compound
Symmetry Code†	An integer encoding the Heuslers symmetry group
Individual Magnetic Moments†	Individual magnetic moments $m_i$ of all constituent atoms
Absolute Magnetic Moments †	Individual absolute magnetic moments $ m_i $ of all constituent atoms
Total magnetic moment†	$M = \sum_i m_i$
Sum of absolute magnetic moments†	$M_{Abs} = \sum_i  m_i $
Magnetic State†	Integer encoding the magnetic state (Ferro, AFM, and Spin-Spiral)
Stoichiometry†	5-Digit integer encoding the stoichiometry of the compound
Density by Atomic Number† <sup>1</sup>	Fractional density of each atomic number is encoded by an individual descriptor
Atomic Number†	Atomic number of the constituents $Z_i$
Number of Neutrons†	Number of neutrons of the constituents
Nominal Mass†	Nominal mass of the constituents atoms
Number of Electrons†	Number of electrons of the constituents
Exact Mass†	Exact mass of the constituents atoms
Atomic Radius†	Atomic radii of the constituents atoms
Number of Valence Electrons†	Number of valence electrons of the constituents atoms $e^{val}$
Covalence Radius†	Covalence radius of the constituents atoms
Period†	Period number in the PSE of the constituents atoms
Electronegativity†	Electronegativity of the constituents atoms $\chi_i$
Van der Waals Radius†	Van der Waals radius of the constituents atoms $r_i^{vdw}$
Electron Affinity†	Electron affinity of the constituents atoms $E_{ea\ i}$

\* Available directly from JuHemd

† Constructed descriptors

‡ Added atomic descriptors - most have four entries per compound

<sup>1</sup> This feature has as many entries (columns) as the JuHemd contains a plurality of unique elements from the PSE

for the Cross-Validation (CV) [33] model evaluation to be meaningful.

The code of the data processing script, as well as the code used to generate the following results and figures, is available [34]. This allows *e.g.* to reevaluate the models if more data is added to the JuHemd. Fig. 2 shows the distribution of atomic numbers across different lattice sites in the Heusler compounds. One can see that Manganese, Chromium, and Iron are contained in a large portion of compounds in the data set.

## B. Model Goals & Evaluation

The prediction of  $T_c$  using the descriptors outlined in the previous section leads to a classical regression task. Such regression models aim at predicting  $T_c$  as accurately as possible. Different metrics are available to evaluate their performance. The evaluation method of choice is

also determined by *e.g.* the error which is desired to be minimized and the importance and impact of outliers in the prediction. The metric used for regression models during this work is the coefficient of determination (Denoted as  $R^2$ ) for test sets, as well as the CV scores.  $R^2$  measures how well the describing features explain the change in the target variable. Hence, we can be sure to choose a model which properly links the descriptors to  $T_c$ .

Besides the regression, we also transformed our problem into a classification task. For the critical temperature, this can be done if one is interested in  $T_c$  to be in a certain range. *E.g.* industrial applications [35] as magnetic storage devices typically require magnetic materials to have a  $T_c$  above  $60^\circ\text{C}$  in operating conditions at least. To maintain this comfortably and ensure long-time magnetic stability at those temperatures, we decided on a threshold of 140K above  $60^\circ\text{C}$  as  $T_c$  for a Heusler compound to be considered as “High”  $T_c$  [16]. Classification

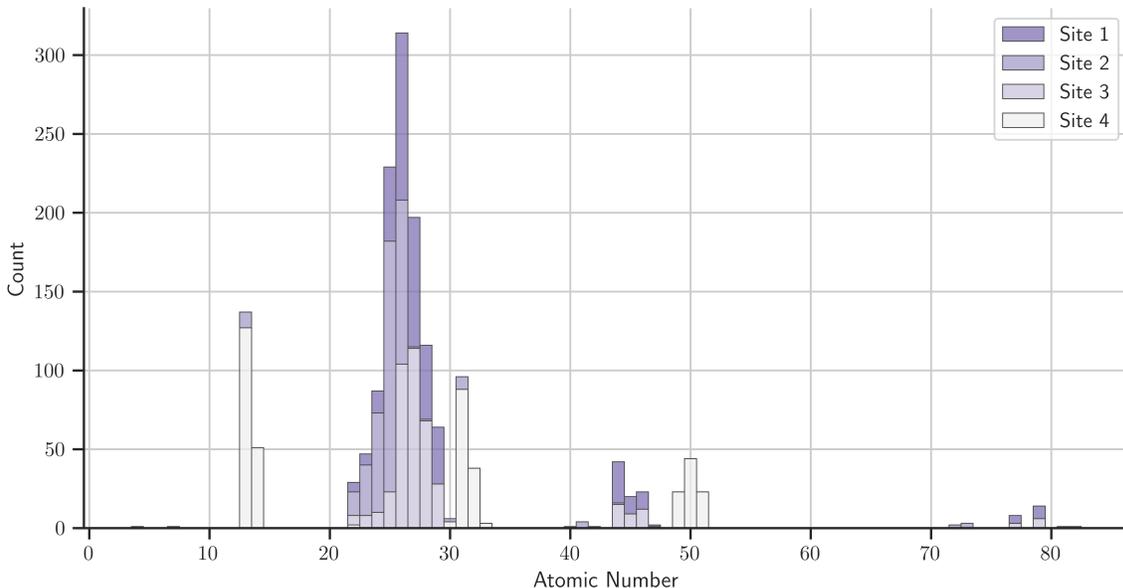


FIG. 2. Distribution of atomic numbers in the GGA data set after processing and cleaning

typically represents an easier modeling task, as the predictive process is less demanding compared to a regression problem. Hence, if one is only interested in magnetic Heusler alloys, which are candidates for an industrial application, but the exact value of  $T_c$  is not of interest in the first place – as the exact value could still be determined in a later step using the established *ab initio* + MC method for the compounds classified as potentially relevant – one can stick to classifying model algorithms. This type of classification model can be used to filter a large number of potential compounds to determine which should be examined further, *e.g.* by a DFT calculation in a high-throughput materials screening context.

For the classification task, additional considerations on how to evaluate the model performance have to be made. The number of correctly predicted categories would be called the accuracy. However, the errors made in the classification do not have the same significance. If a compound is classified as a “low  $T_c$ ” but truly has a “high  $T_c$ ” this means the model misses out on a material with a potential industrial application. The other error the model can make is classifying a “low  $T_c$ ” compound as a “high  $T_c$ ” compound. Which in the worst case means a waste of computational resources in the example above. Therefore, the goal for a classification model in this application has to be to minimize data points falsely classified as “low  $T_c$ ” while still keeping the number of falsely as “High  $T_c$ ” classified compounds low, in order not to waste too many computational resources on these false positives. Hence, we decided to continue with the balanced F1-score, which represents a trade-off between precision and recall.

The model performance is determined using 20 % of

our data as a test data set. This test set has been picked randomized out of the whole data set and is used for calculating the test scores only. This gives an insight into how the model would perform on similar but unseen data. 4-fold CV scores were used in the course of this research in order to perform hyperparameter optimization using a grid search algorithm [36, 37]. Hence, for this hyperparameter optimization, we again partition the training data into a 20% validation set for each individual CV fold and use only the remaining 60% for training. After the hyperparameter optimization, the validation set is included to train the model using the best-performing hyperparameters before proceeding with the testing.

The distribution of the  $T_c$  values in the test set is displayed in Fig. 3. The values above 1500 K can be considered as outliers and are hence removed from the data set before the data is used in an ML workflow.

For all shown scores, it holds: The closer the score is to 1.0, the better the model’s predictive performance is.

### C. ML Techniques

The zoo of ML models and techniques continues to grow year by year. It has already grown to such an extent that it is impossible to cover all possibilities and learning algorithms in a single paper. Hence, we limited our analysis to frequently used and established models. It is also worth mentioning that we excluded neural network models (NNMs) from our research on this data set due to the tabular nature of the data [38, 39].

Before training and evaluating models, it is usually not

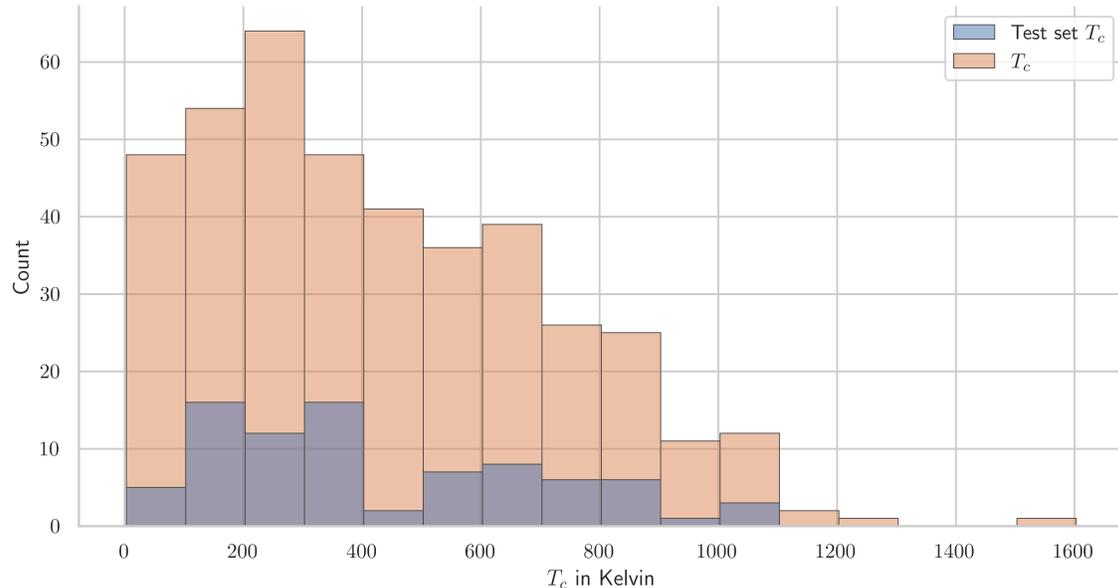


FIG. 3. Distribution of  $T_c$  values in the total data set as well as in the test set only.

possible to anticipate which model will perform best on a given data set. This is commonly referred to as the “no free lunch theorem” [40]. Hence, the regression models we evaluated are depicted in the following table:

Linear	Non-linear	Ensemble
LASSO	K-Nearest Neighbors	Extra Trees
LASSOLars	Decision Tree	Random Forest
Linear Regression		

We have also examined classification models based on similar learning algorithms as some of the regression models depicted in the previous table, as well as a layered indirect classification based on the prediction of the regression models. The indirect classification has been performed to be able to compare the performance of the regression models to their classification counterparts. Since classification is an inherently less complex task than regression, the models would be hard to compare otherwise. The reason underlying this comparison is to determine the best-performing overall model to be used for the feature importance analysis.

#### D. Feature Importance

ML algorithms can be used as black boxes, simply yielding a desired prediction. However, by not applying XAI techniques to understand the model’s prediction, we could miss out on the opportunity to improve our understanding of the underlying physics and validate that the model, indeed, has learned physical key properties and

relations. It is considered best practice to perform feature importance analysis using the model which performs best. This is possible by using the SHAP package [41] including the inbuilt visualization options for the SHAP values. SHAP values represent an ML-specific case of the coalition game theory originated Shapley values [42]. SHAP values can be considered as the estimated average contribution of an individual feature – given a set of features – to the deviation of a predicted value from the mean prediction. Hence, Shapley values can be interpreted as a “driving force” of individual features away from the mean prediction. This allows us to explain the model’s prediction locally for each individual prediction and globally for a set of predictions [41]. The SHAP package is - in principle - model agnostic but has routines optimized for certain model types such as *e.g.* tree-based model [43].

### III. RESULTS & DISCUSSION

In the following, we showcase the scores and results we achieved in training different ML models. In the spirit outlined in the introduction, we investigated the case in which we used descriptors, including results from the DFT simulations, to only learn the results of the Monte-Carlo step first. In a second, independent analysis, we neglected all descriptors that are only available after the DFT simulation and tried to predict  $T_c$  values by using only the atomic data.

For the classification, we will discuss the best-performing model and differences between direct and in-

TABLE II. Regression scores of trained models using the full data set including *ab initio* originated descriptors. The rows show the linear models, the next rows the non-linear models, and the final rows the ensemble models.

	CV Score	Train $R^2$	Test $R^2$
LASSOLars	$\ll 0$	0.77	0.65
LASSO	0.66	0.78	0.66
Linear Reg.	$\ll 0$	0.77	$\ll 0$
Decision Tree Regression	0.59	1.0	0.62
KNN	0.49	0.66	0.57
Extra Trees	0.77	1.0	0.85
Random Forest	0.74	0.97	0.82

direct classification for both the complete descriptor set as well as the reduced descriptor set.

## A. Complete descriptor set

### 1. Regression

A first impression of the predictive performance of two different regression models can be obtained from Fig. 4. For a simple linear model (LASSO) as well as a more complex Extremely Randomized Trees (Extra Trees) regression model, we report the predicted value of  $T_c$  in relation to the value obtained from the full simulation for our test set. While the Lasso results show a systematic error by underestimating the higher values of  $T_c$  while overestimating the critical temperature for the low  $T_c$  Heuslers, this deficiency is substantially reduced in the Extra Trees model. In addition, this model also reproduces the distribution of the values much more accurately and shows less spread around the ideal red line.

This can be confirmed by the plot (Fig. 5) of the difference between the simulated (true) value and the prediction, showing a low relative residual error for the whole temperature range except the very low  $T_c$  values. This error for low values arises from the scale of the very low-temperature values, which enlarges the relative error due to its fractional nature.

A more comprehensive overview of the results of various regression models can be obtained from Table (II), in which we report some performance indicators for the linear, non-linear, and ensemble models.

It is clearly visible that the ensemble models are performing best on the test set and in the CV scoring. However, these good predictions are accompanied by a high degree of overfitting to the training data, easily recognizable by the nearly perfect  $R^2$  score on the training set, *i.e.* a very low or even close to vanishing bias [37]. The fact that only the ensemble models exhibit a reasonably low bias indicates that the complexity of other models does not meet the complexity of the quantity to predict and/or the data. In general, the model complexity has

TABLE III. Direct and indirect classification scores of a model selection using the full data set, including *ab initio* originated descriptors. The rows show the linear models, the next rows the single-tree-based model and the final rows the ensemble-based model setups.

	CV Score	Train F1	Test F1	Test Accuracy
Logistic Reg.	0.82	0.91	0.86	0.89
Indirect LASSO	n/a.	0.86	0.81	0.85
Decision Tree	0.74	1.0	0.75	0.77
Extra Trees	0.82	1.0	0.91	0.93
Indirect Extra	n/a.	1.0	0.89	0.92
Trees model				

to be adjusted to the data complexity [46]. It is clear that a simple linear regression, as well as the K-nearest neighbor model, does not meet this requirement in our case. This finding reflects the complexity of the physical processes responsible for the emergence and stability of magnetic phenomena [47].

The typical approach to cope with overfitting is increasing the regularization [48]. However, even by applying regularization, we could not determine models with improved CV scores, which itself indicated that a lack of regularization is not the root cause for the overfitting. Moreover, when dealing with different iterations of the data set over the course of this study we observed an improvement of the model performance, *e.g.* seen in a decreasing variance, with every increase of the total amount of included Heusler compounds. This indicated that a lack of training data causes a high variance for the more complex models. This also explains the higher test score compared to the CV score. The model in the test case had the full training data available, while for the calculations of the CV score each of the four CV scores – which are depicted here – had only 75 % of the training set available for training.

### 2. Classification

Since classification is a significantly easier task than a regression, we expect to see an improved model performance for each classification model compared to the regression case on this data set. In table III, the results of each linear, non-linear, and ensemble classification model, as well as indirect classification models based on a linear model and an ensemble regression model from section III A 1, are displayed. As expected, the CV scores of the classification models are significantly higher than the scores of the regression models, which corresponds to a lower bias. Similarly, the results for the test set are closer to the ideal prediction, as there is less variance occurring than for the regression models. This aligns with our interpretation of the overfitting in the

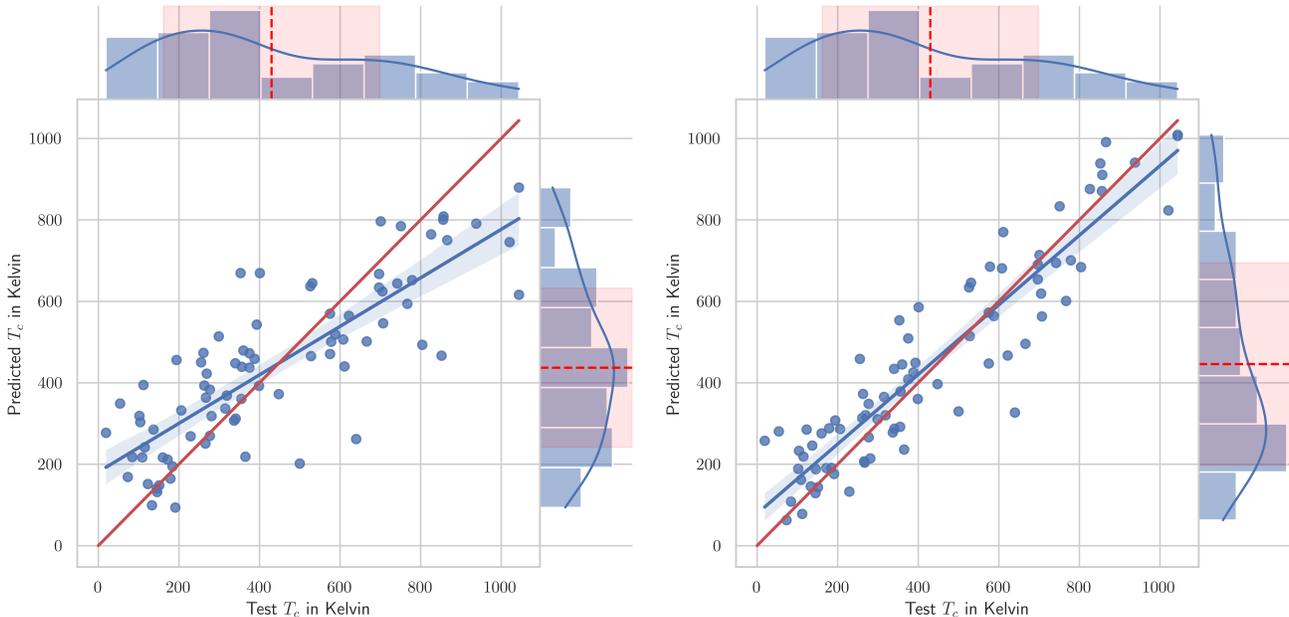


FIG. 4. Prediction result of two models for the test set. On the left, the data for the linear LASSO model is shown, while the right panel shows the data from the Extremely Randomized Trees (Extra Trees) Regression model [44]. The red line indicates a perfect match between the predicted and expected data, the blue line (with shade) a linear regression through the predicted data points (with a 95% CI envelope for the regression, computed using a bootstrap [45] based method). On the side distribution plots of the test sets, true  $T_c$  values and the predictions are added.

regression case due to the fact that classification is an easier task, so there is less training data required to fit classification models to the data as the complexity of the quantity to predict is reduced from a continuous quantity to a binary value. This reduction is only possible as we know which minimum  $T_c$  values are required to be relevant to an industrial application.

### 3. Feature Importance

After searching for a working set of hyperparameters in section III A 1 using the training set, we used the determined hyperparameters and chose the Extra Trees Regression [44] as our best-performing model to conduct a feature importance analysis using the SHAP package [41] and the corresponding SHAP values. The SHAP values have been calculated for the training data set. These values for the most relevant features are shown in Fig. 6. Besides the SHAP values, the color of the data points encodes the relative scale of the feature for each individual data point. Meaning that if there is a clear horizontal color fade visible, this implies a systematic impact of this feature for the predicted quantity.

From Fig. 6 one can see that for the Extra Trees model, the absolute magnetic moment of the compound has the largest impact on the  $T_c$  prediction. All nine most relevant quantities are either magnetic moments or indirectly related to magnetism (*e.g.* the Cobalt density of the

compound), which confirms that the magnetic material-specific properties indeed have the largest impact on the value of the critical temperature. While all the nine quantities are positively correlated to  $T_c$ , *i.e.* have an increasing impact on the  $T_c$  when they increase too. For some of the quantities, this is of course an artifact of our descriptor construction. For example, we encode the magnetic state as an integer, with the smallest possible encoding 000 denoting that the material forms neither a ferromagnet, an anti-ferromagnet, nor a spin-spiral. In contrast, the fact that the model assigns most significance to the nine quantities listed here was obtained without providing any physical knowledge of the system, besides the fact that we included these descriptors in the first place. Thus, the modeling singled out that these parameters indicate the kind of “physical insight” that can be obtained from ML. For example, the high relevance of the absolute magnetic moment of the compound for the  $T_c$  is of course in line with the relation one would obtain from even very basic physics models of magnetism.

Looking at Fig. 7 one can draw even more conclusions. It can be seen that  $T_c$  is not simply proportional to  $M_{Abs}$ . Instead, Heusler alloys with a higher  $M_{Abs}$  can show a higher  $T_c$ . However, while a high absolute magnetic moment does not guarantee the emergence of a high  $T_c$ , a low  $M_{Abs}$  prevents the occurrence of high  $T_c$  values. Therefore, it is safe to say  $M_{Abs}$  is acting as an upper boundary:

$$T_c \leq CM_{Abs} \quad (3)$$

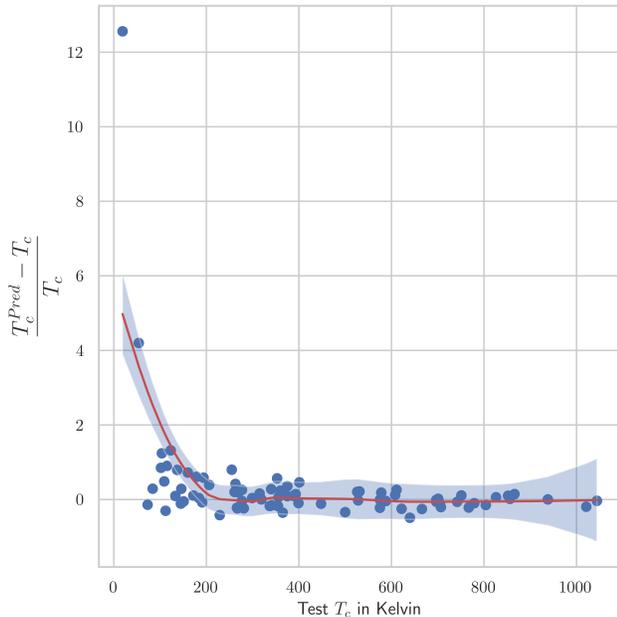


FIG. 5. Residual depiction of the Extra Trees regression shown on the right of Fig. 4 including a LOWESS smoothing applied to the data points with a pointwise 95 % CI envelope.

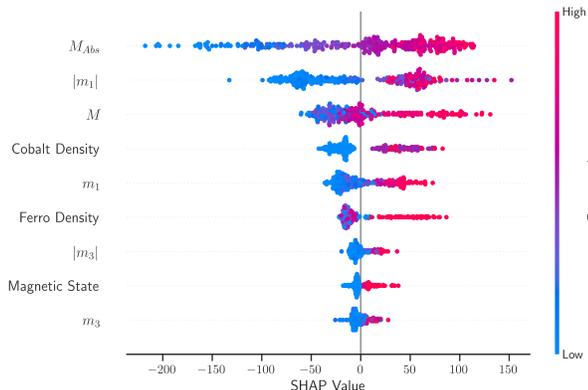


FIG. 6. SHAP beeswarm summary plot of the nine descriptors with the largest SHAP values [41, 42]

## B. Data set without DFT-originated descriptors

### 1. Regression

Retraining the Extra Trees Regression model as well as the LASSO model to the reduced descriptor set and again performing hyperparameter optimization using a grid search algorithm, we achieved the regression scores displayed in table IV.

As expected, one can observe a clear decrease in performance compared to the case where DFT-originated descriptors have been used. In particular, the LASSO results now have huge deviations such that one could

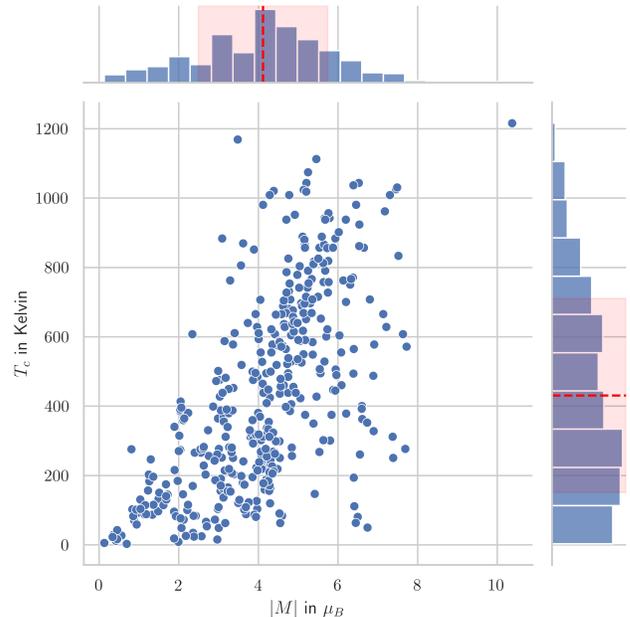


FIG. 7. Relation between absolute magnetic moment of the Heusler compounds to their  $T_c$  values for the whole data set.

TABLE IV. Regression scores of promising models using the reduced data set, excluding *ab initio* originated descriptors.

	CV Score	Train $R^2$	Test $R^2$
Extra Trees	0.52	1.0	0.76
LASSO	0.31	0.58	0.63

question its fitness for any practical application. Therefore, we can already conclude, that a prediction of the critical temperature without the use of the basic magnetic properties predicted by a DFT simulation is not really possible in our scenario. Therefore, we decided not to analyze this further, but to concentrate on the easier classification task.

### 2. Classification

The achieved classification model results using no DFT-originated descriptors at all are displayed in table V. This table contains exactly the same models as seen before in table III. From the test set of 82 compounds, our constructed indirect Extra Trees classification model managed to correctly classify 47 “Low”  $T_c$  and 29 “High”  $T_c$  compounds. Of each class, 3 compounds have been wrongly predicted. We consider falsely classifying a “Low”  $T_c$  compound as a “High”  $T_c$  not so relevant for practical application. The worst outcome in a potential use case is that the model suggests a “High”  $T_c$  compound, and when computing it using a more sophisticated – and hence computationally more intensive – ap-

TABLE V. Direct and indirect classification scores of a model selection using the reduced data set, excluding *ab initio* originated descriptors. The rows show the linear models, the next rows are the single-tree-based model, and the final rows are the ensemble-based model setups.

	CV Score	Train F1	Test F1	Test Accuracy
Logistic Reg.	0.68	0.75	0.75	0.79
Indirect LASSO	n/a.	0.75	0.8	0.88
Decision Tree	0.66	1.0	0.8	0.84
Extra Trees	0.74	1.0	0.84	0.87
Indirect Extra Trees model	n/a.	1.0	0.91	0.93

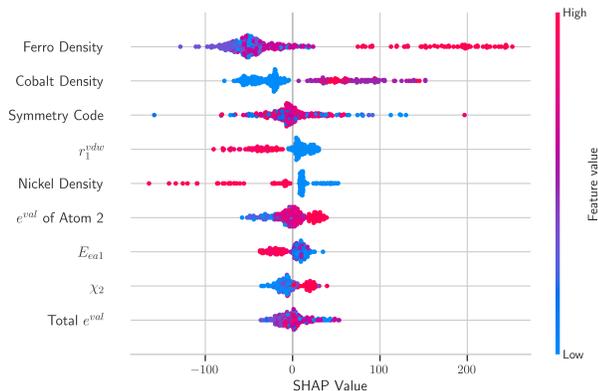


FIG. 8. SHAP beeswarm summary plot of the nine descriptors of the data set without the DFT-originated data with the largest SHAP values.

proach, one finds that the “High”  $T_c$  label has been falsely assigned. However, if a “High”  $T_c$  compound is classified as “Low”  $T_c$  in a high-throughput screening process it will probably never be computed with a more sophisticated approach, which causes this compound to be potentially “lost” for future research. In the case of this model, we saw that this crucial error for falsely classifying a “High”  $T_c$  Heusler as a “Low”  $T_c$  Heusler is below 5% and hence meets typical confidence criteria. This concludes that the indirect Extra Trees classification is capable of classifying the  $T_c$  in “High” and “Low” values even without the DFT-originated data. While “Low” means  $T_c$  is too low to be relevant for industry applications.

### 3. Feature Importance

Computing the SHAP values for the reduced descriptor set and visualizing them as we did before results in the beeswarm plot shown in Fig. 8. As we can see in Fig. 8, removing the DFT-originated descriptors and, therefore, the calculated magnetic moments caused other quantities to become more impactful. As one can expect, these are very closely related to the magnetic moments (*e.g.*

the density of ferromagnetic materials in the compound as well as the Cobalt and Nickel densities). However, now we observe more complex relations than in the previous feature importance plot, demonstrating the lower significance of these quantities for the critical temperatures. We can see a negative correlation between the van der Waals radius of the atom on site one ( $r_1^{vdw}$ ), the Nickel density in the compound, and the electron affinity of the atom on site one ( $E_{ea1}$ ) with a decreasing  $T_c$  as these quantities increase. For the fraction of ferromagnetic atoms, the effect is much less obvious. We can see that very high densities of ferromagnetic atoms in the compound contribute to a largely increased  $T_c$  prediction. However, on the other hand, a low density of ferromagnetic atoms does not lead to an equally decreased prediction of  $T_c$ . Interestingly, this reflects our previous result that a large absolute magnetic moment corresponds to an upper boundary for the  $T_c$ . Since a large amount of ferromagnetic compound constituents is highly correlated with a large magnetic moment. The required and obtained model complexity is also observed in the SHAP values of the symmetry code. Since this is an arbitrary-ordered label for the symmetry group of the compound, there is no clear order of the feature value that correlates with the  $T_c$ . However, the model seems to have learned that some feature values have a larger impact on  $T_c$  than others, which is indeed possible.

As the density of the ferromagnetic atoms, the cobalt and nickel atoms turn out to be relevant quantities in Fig. 8 we investigated their correlation with  $T_c$  in more detail as depicted in Fig. 9. The depicted fractional density histograms confirm the trends we were hinted at by the SHAP beeswarm plot. It is easily visible that a large density of ferromagnetic atoms in the compound is indeed contributing to a larger  $T_c$  value, with one exception: The anti-ferromagnetic case. We can see that there are a few materials that have no ferromagnetic atoms in the compound at all but still a very high  $T_c$ . These are strong anti-ferromagnets. This finding can be related to our previous result for the modeling, including the DFT-based magnetization values, in which we have seen that the SHAP values of  $M_{Abs}$  hint at a larger impact than those of  $M$ . As the anti-ferromagnetic compounds have a vanishing  $M$  but a large  $M_{Abs}$  while resulting in a stable anti-ferromagnetic state with a large  $T_c$ . The same relation is, in principle, true for the cobalt density. However, there are fewer compounds containing cobalt in the data than iron. The Nickel density has – for increasing densities – a negative impact on  $T_c$  according to Fig. 8, and as we can see in Fig. 9, this also emerges from the data. It seems that the presence of Cobalt is not as helpful in contributing to a stable magnetic state as *e.g.* Iron.

## IV. SUMMARY AND OUTLOOK

This work can be seen as a small-scale sandbox-type case study in which lightweight ML algorithms can add

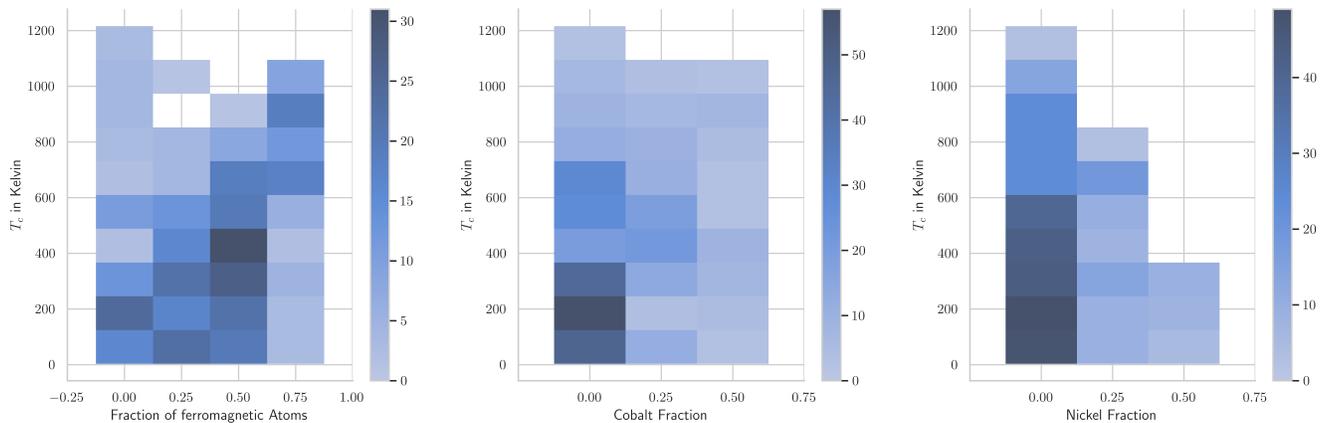


FIG. 9. Relation between  $T_c$  and the fraction of ferromagnetic elements (left), cobalt atoms (center) and nickel (right) shown as a heatmap style histogram. The darker the color, the more compounds can be found in the colored region.

value to existing *ab initio* data and eventually replace costly computational steps in layered calculation workflows in the future.

It was demonstrated that qualitative predictions for material-specific properties are achievable with very small errors, even for the limited data set sizes common in materials science. Also, the expectation that the quantitative prediction is much more difficult and requires descriptors with much higher predictive power, has been confirmed. However, we could also demonstrate that even very simple and readily available descriptors not based on any actual calculation in combination with sufficiently complex models could be utilized in a classification task typically part of any high throughput screening. As demonstrated in this paper, there is a potential use for ML methods in materials science, even in quantitatively predicting properties as complex as the  $T_c$ . It is imaginable to perform similar studies on existing data sets of other material families beyond the Heusler alloys. However, one has to consider that the structural homogeneity of the material class we studied here simplified the complexity of the modeling task. This implies that if one would translate the methodical insights gained from this data set of Heusler alloys to a different material class, that there should either be a similar structural homogeneity or if the structural complexity is increased one also has to choose descriptors with equally elevated descriptive value.

By performing feature importance analysis with XAI techniques – such as SHAP values – we gained physical insights about the relations of the target quantity to the included features, as well as the determining properties of the studied material class given in the examined data set. Such analysis can provide a link between a complex ML process with a hard-to-expose underlying mechanism and true physical insight and the gain of knowledge of the system. In this study, be rediscovered dependencies

expected from simple physical models without actually providing such knowledge to the process.

Finally, we would like to stress that the methodical approach described in this paper is not limited to predicting  $T_c$  or any other magnetic quantity, but that it can be transferred to any other material-specific property. We believe it is even possible to discover that known materials have currently unknown properties using predictive modeling.

## ACKNOWLEDGMENTS

This work was performed as part of the Helmholtz School for Data Science in Life, Earth and Energy (HDS-LEE) and received funding from the Helmholtz Association of German Research Centres. Since parts of the data processing has been performed and the displayed visualizations have been created using dedicated open-source packages, we acknowledge them here [49–57].

We acknowledge Stefano Sanvito for the inspiration to represent the fractional density of each atomic number in the data set as a standalone descriptor, which we gathered from one of his talks.

The authors thank Fabian Lux for initially hinting us at the Extra Trees regression model. We acknowledge Dirk Witthaut for pointing out the advantages of SHAP values in XAI in comparison to model bound feature importance methods. The authors thank Roman Kováčik for discussing the structure and contents of the data published in Ref. 30. Hence, we acknowledge the computing time granted by the RWTH Aachen University (Project: jara0182) which was required in order to collect the data. As the original database as well as the ML-ready data was hosted by Materials Cloud, we acknowledge them [58].

- [1] H. Shi, S. Liu, J. Chen, X. Li, Q. Ma, and B. Yu, Predicting drug-target interactions using lasso with random forest based on evolutionary information and chemical structure, *Genomics* **111**, 1839 (2019).
- [2] W. L. Bi, A. Hosny, M. B. Schabath, M. L. Giger, N. J. Birkbak, A. Mehrtash, T. Allison, O. Arnaout, C. Abbosh, I. F. Dunn, R. H. Mak, R. M. Tamimi, C. M. Tempany, C. Swanton, U. Hoffmann, L. H. Schwartz, R. J. Gillies, R. Y. Huang, and H. J. W. L. Aerts, Artificial intelligence in cancer imaging: Clinical challenges and applications, CA: A Cancer Journal for Clinicians 10.3322/caac.21552 (2019).
- [3] R. V. Shah, G. Grennan, M. Zafar-Khan, F. Alim, S. Dey, D. Ramanathan, and J. Mishra, Personalized machine learning of depressed mood using wearables, *Translational Psychiatry* **11**, 10.1038/s41398-021-01445-0 (2021).
- [4] T. Zhou, Z. Song, and K. Sundmacher, Big data creates new opportunities for materials research: A review on methods and applications of machine learning for materials design, *Engineering* **5**, 1017 (2019).
- [5] H. J. Kulik, T. Hammerschmidt, J. Schmidt, S. Botti, M. A. L. Marques, M. Boley, M. Scheffler, M. Todorović, P. Rinke, C. Oses, A. Smolyanyuk, S. Curtarolo, A. Tkatchenko, A. P. Bartók, S. Manzhos, M. Ihara, T. Carrington, J. Behler, O. Isayev, M. Veit, A. Grisafi, J. Nigam, M. Ceriotti, K. T. Schütt, J. Westermayr, M. Gastegger, R. J. Maurer, B. Kalita, K. Burke, R. Nagai, R. Akashi, O. Sugino, J. Hermann, F. Noé, S. Pilati, C. Draxl, M. Kuban, S. Rigamonti, M. Scheidgen, M. Esters, D. Hicks, C. Toher, P. V. Balachandran, I. Tamblin, S. Whitelam, C. Bellinger, and L. M. Ghiringhelli, Roadmap on machine learning in electronic structure, *Electronic Structure* **4**, 023004 (2022).
- [6] Y. Igarashi, K. Nagata, T. Kuwatani, T. Omori, Y. Nakanishi-Ohno, and M. Okada, Three levels of data-driven science, *Journal of Physics: Conference Series* **699**, 012001 (2016).
- [7] P. J. Garcia Nieto, E. García-Gonzalo, and J. Paredes-Sánchez, Prediction of the critical temperature of a superconductor by using the woa/mars, ridge, lasso and elastic-net machine learning techniques, *Neural Computing and Applications* **33**, 1 (2021).
- [8] V. L. Deringer, N. Bernstein, A. P. Bartók, M. J. Cliffe, R. N. Kerber, L. E. Marbella, C. P. Grey, S. R. Elliott, and G. Csányi, Realistic atomistic structure of amorphous silicon from machine-learning-driven molecular dynamics, *The Journal of Physical Chemistry Letters* **9**, 2879 (2018).
- [9] K. Takahashi and Y. Tanaka, Material synthesis and design from first principle calculations and machine learning, *Computational Materials Science* **112**, 364 (2016).
- [10] A. Jain, G. Hautier, S. P. Ong, and K. Persson, New opportunities for materials informatics: Resources and data mining techniques for uncovering hidden relationships, *Journal of Materials Research* **31**, 977 (2016).
- [11] H. Ucar, D. Paudyal, and K. Choudhary, Machine learning predicted magnetic entropy change using chemical descriptors across a large compositional landscape, *Computational Materials Science* **209**, 111414 (2022).
- [12] R. Ramprasad, R. Batra, G. Pilania, A. Mannodi-Kanakithodi, and C. Kim, Machine learning in materials informatics: recent applications and prospects, *npj Computational Materials* **3**, 10.1038/s41524-017-0056-5 (2017).
- [13] J. Schmidt, M. R. G. Marques, S. Botti, and M. A. L. Marques, Recent advances and applications of machine learning in solid-state materials science, *npj Computational Materials* **5**, 10.1038/s41524-019-0221-0 (2019).
- [14] O. Gutfleisch, M. A. Willard, E. Brück, C. H. Chen, S. G. Sankar, and J. P. Liu, Magnetic materials and devices for the 21st century: Stronger, lighter, and more energy efficient, *Advanced Materials* **23**, 821 (2011), <https://onlinelibrary.wiley.com/doi/pdf/10.1002/adma.201002180>.
- [15] V. Laliena, G. Albalade, and J. Campo, Stability of the skyrmion lattice near the critical temperature in cubic helimagnets, *Phys. Rev. B* **98**, 224407 (2018).
- [16] T. Shang, E. Canévet, M. Morin, D. Sheptyakov, M. T. Fernández-Díaz, E. Pomjakushina, and M. Medarde, Design of magnetic spirals in layered perovskites: Extending the stability range far beyond room temperature, *Science Advances* **4**, 10.1126/sciadv.aau6386 (2018), <https://www.science.org/doi/pdf/10.1126/sciadv.aau6386>.
- [17] J. F. Belot, V. Taufour, S. Sanvito, and G. L. W. Hart, Machine learning predictions of high-curie-temperature materials (2023), arXiv:2307.06879 [cond-mat.mtrl-sci].
- [18] J. Nelson and S. Sanvito, Predicting the curie temperature of ferromagnets using machine learning, *Physical Review Materials* **3**, 104405 (2019).
- [19] F. Heusler, W. Starck, and E. Haupt, Magnetisch-chemische studien, *Verh. Dtsch. Phys. Ges* **5**, 219 (1903).
- [20] F. Heusler and E. Take, The nature of the heusler alloys, *Transactions of the Faraday Society* **8**, 169 (1912).
- [21] H. Uzunok, E. Karaca, S. Bağcı, and H. Tütüncü, Physical properties and superconductivity of heusler compound  $\text{liga}_2\text{rh}$ : A first-principles calculation, *Solid State Communications* **311**, 113859 (2020).
- [22] A. Roy, J. W. Bennett, K. M. Rabe, and D. Vanderbilt, Half-heusler semiconductors as piezoelectrics, *Physical Review Letters* **109**, 10.1103/physrevlett.109.037602 (2012).
- [23] Q. Gao, I. Opahle, O. Gutfleisch, and H. Zhang, Designing rare-earth free permanent magnets in heusler alloys via interstitial doping, *Acta Materialia* **186** (2020).
- [24] S. Idrissi, S. Ziti, H. Labrim, and L. Bahmad, Half-metallicity and magnetism in the full heusler alloy  $\text{fe}_2\text{mnsn}$  with  $\text{l}21$  and  $\text{XA}$  stability ordering phases, *Journal of Low Temperature Physics* **202**, 343 (2021).
- [25] A. Davidson, V. P. Amin, W. S. Aljuaid, P. M. Haney, and X. Fan, Perspectives of electrically generated spin currents in ferromagnetic materials, *Physics Letters A* **384**, 126228 (2020).
- [26] A. Hirohata and D. C. Lloyd, Heusler alloys for metal spintronics, *MRS Bulletin* **47**, 593 (2022).
- [27] S. Sanvito, C. Oses, J. Xue, A. Tiwari, M. Zic, T. Archer, P. Tozman, M. Venkatesan, M. Coey, and S. Curtarolo, Accelerated discovery of new magnets in the heusler alloy family, *Science Advances* **3**, e1602241 (2017), <https://www.science.org/doi/pdf/10.1126/sciadv.1602241>.
- [28] P. Rüßmann, Be-Zimmermann, G. Géranton, C. Oran, and E. Rabel, *Judftteam/jukkr: v3.6* (2022).
- [29] X. Zhong, B. Gallagher, S. Liu, B. Kailkhura, A. Hisz-

- panski, and T. Y.-J. Han, Explainable machine learning in materials science, *npj Computational Materials* **8**, 10.1038/s41524-022-00884-7 (2022).
- [30] R. Kováčik, P. Mavropoulos, and S. Blügel, The juhemd (jülich-heusler-magnetic-database) of the monte carlo simulated critical temperatures of the magnetic phase transition for experimentally reported heusler and heusler-like materials (2022).
- [31] J. P. Perdew, K. Burke, and Y. Wang, Generalized gradient approximation for the exchange-correlation hole of a many-electron system, *Phys. Rev. B* **54**, 16533 (1996).
- [32] B. Fricke, W.-D. Sepp, T. Bastug, S. Varga, K. Schulze, J. Anton, and V. Pershina, Use of the DV  $\alpha$ -method in the field of superheavy atoms, in *Advances in Quantum Chemistry* (Elsevier, 1998) pp. 109–121.
- [33] T. Hastie, R. Tibshirani, and J. Friedman, *The elements of statistical learning: data mining, inference and prediction*, 2nd ed. (Springer, 2009).
- [34] R. Hilgers, D. Wortmann, and S. Blügel, Data processing for the juhemd database and ml-training and evaluation scripts (2022).
- [35] C. Felser, L. Wollmann, S. Chadov, G. H. Fecher, and S. S. P. Parkin, Basics and prospective of magnetic heusler compounds, *APL Materials* **3**, 041518 (2015).
- [36] F. Hutter, L. Kotthoff, and J. Vanschoren, eds., *Automatic Machine Learning: Methods, Systems, Challenges* (Springer, 2019).
- [37] G. James, D. Witten, T. Hastie, and R. Tibshirani, *An Introduction to Statistical Learning: with Applications in R*, 2nd ed. (Springer, 2021).
- [38] R. Shwartz-Ziv and A. Armon, Tabular data: Deep learning is not all you need, *Information Fusion* **81**, 84 (2022).
- [39] L. Grinsztajn, E. Oyallon, and G. Varoquaux, Why do tree-based models still outperform deep learning on typical tabular data?, *Advances in Neural Information Processing Systems* **35**, 507 (2022).
- [40] D. Wolpert and W. Macready, No free lunch theorems for optimization. *IEEE Transactions on Evolutionary Computation* **1**(1), 67-82, *Evolutionary Computation*, IEEE Transactions on **1**, 67 (1997).
- [41] S. M. Lundberg and S.-I. Lee, A unified approach to interpreting model predictions, in *Advances in Neural Information Processing Systems 30*, edited by I. Guyon, U. V. Luxburg, S. Bengio, H. Wallach, R. Fergus, S. Vishwanathan, and R. Garnett (Curran Associates, Inc., 2017) pp. 4765–4774.
- [42] L. S. Shapley, A value for n-person games, in *Contributions to the Theory of Games II*, edited by H. W. Kuhn and A. W. Tucker (Princeton University Press, Princeton, 1953) pp. 307–317.
- [43] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee, From local explanations to global understanding with explainable AI for trees, *Nature Machine Intelligence* **2**, 56 (2020).
- [44] P. Geurts, D. Ernst, and L. Wehenkel, Extremely randomized trees, *Machine Learning* **63**, 3 (2006).
- [45] B. Efron, Bootstrap Methods: Another Look at the Jackknife, *The Annals of Statistics* **7**, 1 (1979).
- [46] L. Li and Y. S. Abu-Mostafa, Data complexity in machine learning, *Computer Science Technical Reports* 10.7907/Z9319SW2 (2006).
- [47] D. Gatteschi and L. Bogani, Complexity in molecular magnetism, in *Complexity in Chemistry and Beyond: Interplay Theory and Experiment*, edited by C. Hill and D. G. Musaev (Springer Netherlands, Dordrecht, 2012) pp. 49–72.
- [48] X. Ying, An overview of overfitting and its solutions, *Journal of Physics: Conference Series* **1168**, 022022 (2019).
- [49] C. R. Harris, K. J. Millman, S. J. van der Walt, R. Gommers, P. Virtanen, D. Cournapeau, E. Wieser, J. Taylor, S. Berg, N. J. Smith, R. Kern, M. Picus, S. Hoyer, M. H. van Kerkwijk, M. Brett, A. Haldane, J. F. del Río, M. Wiebe, P. Peterson, P. Gérard-Marchant, K. Sheppard, T. Reddy, W. Weckesser, H. Abbasi, C. Gohlke, and T. E. Oliphant, Array programming with NumPy, *Nature* **585**, 357 (2020).
- [50] Lukasz. Mentel, mendelev – a python resource for properties of chemical elements, ions and isotopes (2014).
- [51] T. Tantau, *The TikZ and PGF Packages* (2013).
- [52] J. D. Hunter, Matplotlib: A 2d graphics environment, *Computing in Science & Engineering* **9**, 90 (2007).
- [53] T. A. Caswell, M. Droettboom, A. Lee, E. S. De Andrade, T. Hoffmann, J. Hunter, J. Klymak, E. Firing, D. Stansby, N. Varoquaux, J. H. Nielsen, B. Root, R. May, P. Elson, J. K. Seppänen, D. Dale, Jae-Joon Lee, D. McDougall, A. Straw, P. Hobson, , Hannah, C. Gohlke, T. S. Yu, E. Ma, A. F. Vincent, S. Silvester, C. Moad, N. Kniazev, E. Ernest, and P. Ivanov, matplotlib/matplotlib: Rel: v3.4.3 (2021).
- [54] M. L. Waskom, seaborn: statistical data visualization, *Journal of Open Source Software* **6**, 3021 (2021).
- [55] W. S. Cleveland, Robust locally weighted regression and smoothing scatterplots, *Journal of the American Statistical Association* **74**, 829 (1979).
- [56] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay, Scikit-learn: Machine learning in Python, *Journal of Machine Learning Research* **12**, 2825 (2011).
- [57] S. Seabold and J. Perktold, statsmodels: Econometric and statistical modeling with python, in *9th Python in Science Conference* (2010).
- [58] L. Talirz, S. Kumbhar, E. Passaro, A. V. Yakutovich, V. Granata, F. Gargiulo, M. Borelli, M. Uhrin, S. P. Huber, S. Zoupanos, C. S. Adorf, C. W. Andersen, O. Schütt, C. A. Pignedoli, D. Passerone, J. VandeVondele, T. C. Schulthess, B. Smit, G. Pizzi, and N. Marzari, Materials cloud, a platform for open computational science, *Scientific Data* **7**, 10.1038/s41597-020-00637-5 (2020).